# RAVEN: Multitask Retrieval Augmented Vision-Language Learning

**Varun Nagaraj Rao**[1][*]**, Siddharth Choudhary**[2]**, Aditya Deshpande**[3][†]**, Ravi Kumar Satzoda**[2]**,
Srikar Appalaraju**[2][‡]

[1]Princeton University, [2]AWS AI Labs, [3]Apple

## Abstract

The scaling of large language models to encode all the world's knowledge in model parameters is unsustainable and has exacerbated resource barriers. Retrieval-Augmented Generation (RAG) presents a potential solution, yet its application to vision-language models (VLMs) is underexplored. Existing methods focus on models designed for single tasks. Furthermore, they're limited by the need for resource intensive pretraining, additional parameter requirements, unaddressed modality prioritization and lack of clear benefit over non-retrieval baselines. This paper introduces RAVEN, a multitask retrieval augmented VLM framework that enhances base VLMs through efficient, task-specific fine-tuning. By integrating retrieval augmented samples without the need for additional retrieval-specific parameters, we show that the model acquires retrieval properties that are effective across multiple tasks. Our results and extensive ablations across retrieved modalities for the image captioning and VQA tasks indicate significant performance improvements compared to non retrieved baselines – +1 CIDEr on MSCOCO, +4 CIDEr on NoCaps, and nearly a +3% accuracy on specific VQA question types. This underscores the efficacy of applying RAG approaches to VLMs, marking a stride toward more efficient and accessible multimodal learning.

## 1 Introduction

The rapid growth in model sizes in NLP, as highlighted by OpenAI's LLM progression from GPT-2's 1.5 billion parameters (Radford et al., 2019) to GPT-3's 175 billion (Brown et al., 2020), and further to over a trillion in GPT-4 (OpenAI, 2023), is a source of increasing concern. This trend requires more data and computational power, leading to higher carbon emissions and presenting significant obstacles for less-resourced researchers (Strubell et al., 2019). In response, the field is pivoting to approaches like Retrieval-Augmented Generation (RAG) (Lewis et al., 2020), which incorporates external non-parametric world knowledge into a pretrained language model, removing the necessity of encoding all information directly into the model's parameters. However, this strategy is not yet widely applied in vision-language models (VLMs) (Li et al., 2022; Wang et al., 2021; Alayrac et al., 2022; Chen et al., 2022c; Radford et al., 2021; Wang et al., 2022b), which process both image and textual data, and are typically more resource-intensive. Moreover, VLMs often rely on massive datasets like LAION-5B (Schuhmann et al., 2022), presenting a significant opportunity for performance gains through retrieval augmentation.

The scant prior work exploring retrieval augmentation applied to VLMs, although promising, is beset with several limitations. Most importantly, they rely on pretraining with retrieval specific parameters (Hu et al., 2023; Ramos et al., 2023b; Yang et al., 2023); as a result the performance improvement over non-retrieval baselines cannot be established and the benefit due to retrieval augmentation cannot be independently discerned. Next, model architectures are suited to only a single task, and therefore, experimental evaluation is also only presented on a single task e.g. on image captioning (Ramos et al., 2023b,a; Yasunaga et al., 2023); other image-to-text tasks like VQA are ignored. Further, the decision on which modality to prioritize during retrieval - textual, visual, or a combination of both - is not established. Some works (Yasunaga et al., 2023; Chen et al., 2022a) retrieve and concatenate both image and text, while others (Ramos et al., 2023a,b; Yang et al., 2023) only retrieve text, even though they all evaluate on image-to-text tasks. Finally, we also observe that overlaps

---

between the retrieval and pre-training/fine-tuning datasets exist; for example, Ramos et al. (2023a,b) pretrain and retrieve from MSCOCO. This can confound the benefits attributed to the RAG approach, underscoring the need for a larger and non-overlapping external memory.

In this paper, we present RAVEN (see Figure 1), a multitask retrieval augmented framework adaptable to any multitask base VLM. The framework does not rely on pretraining with retrieval specific parameters, and is suitable to a variety of tasks. Importantly, the design of RAVEN allows for a comprehensive investigation of the performance benefits over non-retrieval baselines, and implications of retrieving and using different modalities. Specifically, our key contributions are as follows:

1. We are the first to design a multitask retrieval augmented VLM framework (RAVEN), which relies on only fine-tuning, no retrieval specific trainable parameters and is adaptable to any multitask base VLM.

2. Our method allows for comprehensive ablations which examine the trade-offs between retrieval modalities and their advantages relative to non-retrieval baselines while using a non-overlapping and larger external memory.

3. We demonstrate the benefits and limitations of our approach on Image Captioning and VQA through quantitative and qualitative analysis. Our results achieve a new state-of-the-art performance improvement compared to non retrieved baselines: +1 CIDEr on MSCOCO, +4 CIDEr on NoCaps (using magnitudes of fewer parameters than prior works), and nearly a +3% accuracy on specific VQA question types.

Broadly, our work expands the empirical knowledge on RAG techniques and contributes to the rapidly growing body of work focusing on their applications to multitask VLMs. Ultimately, this work establishes a clearer understanding of the role of retrieval augmentation in VLMs, paving the way for more efficient and sustainable approaches in the field.

## 2   Related Work

### 2.1   Vision Language Models

Vision language models are an emerging type of multi-modal AI system that can process both vi-

sual and textual data (Appalaraju et al., 2024, 2021) They build upon recent advances in computer vision and natural language processing to generate textual descriptions of images, answer visual questions, and perform other vision-and-language tasks. Earlier works in this direction unified multiple tasks like image captioning, image classification etc. using a simple sequence-to-sequence framework. Some notable examples include OFA (Wang et al., 2022b), GIT (Wang et al., 2022a), SimVLM (Wang et al., 2021). Recent vision-language models (Biten et al., 2022) augment pre-trained large language models with visual encoder. For example, Frozen (Tsimpoukelli et al., 2021), Flamingo (Alayrac et al., 2022), BLIP (Li et al., 2022), InstructBLIP (Dai et al., 2023), LLaVA (Liu et al., 2023), MiniGPT-4 (Zhu et al., 2023), Kosmos-1 (Huang et al., 2023), Pali (Chen et al., 2022c). In this work, we use OFA (Wang et al., 2022b) as the baseline rather than using VLMs augmented with pretrained LLMs. This choice allows us to remove the effects of in-context learning abilities of the pretrained language models from the resulting enhancement brought by retrieval-augmented vision-language modeling.

### 2.2   Retrieval Augmented Generation in NLP

Retrieval augmentation has become an important technique for improving natural language processing models. One of the first works in this area was kNN-LM by Khandelwal et al. (Khandelwal et al., 2020) who showed how interpolating over nearest neighbors from any text collection could improve generalization. This was followed by RETRO (Borgeaud et al., 2021), which scaled up the retrieval corpus to trillions of tokens. Another line of work has focused on integrating Wikipedia passages directly into models like REALM (Guu et al., 2020), RAG (Lewis et al., 2020), and FiD (Izacard and Grave, 2021). By retrieving and conditioning on relevant Wikipedia passages, these models can better perform knowledge-intensive downstream tasks like question answering. Overall, retrieval augmentation has proven to be a highly effective way of injecting knowledge into language models to improve their capabilities. The techniques have progressed from simple corpus retrieval to integrated and scalable architectures that retrieve from large knowledge bases like Wikipedia.
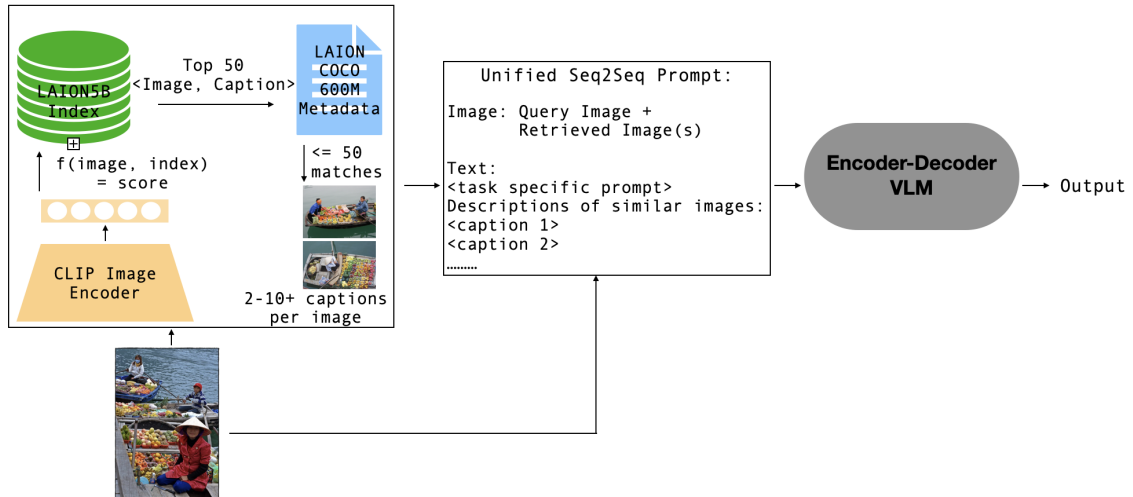
Figure 1: Illustration of our RAVEN framework. Given an input image, we retrieve image-text pairs from an external memory. Subsequently, we use a multitask pretrained base vision-language model (VLM) to encode the retrieved samples along with the query and decode to generate an output by attending over both the query and retrieved samples.

## 2.3 Retrieval Augmented Generation in VLMs

Recent years have seen significant progress in extending retrieval-augmented generation to vision-language models. One of the earliest works is Multimodal Retrieval-Augmented Transformer (MuRAG) which utilizes non-parametric multimodal memory for language generation improvement (Chen et al., 2022a). In image-to-text generation, Smallcap (Ramos et al., 2023b), exhibits competitive performance on COCO and other domains through retrieval from target-domain data. Sarto et al.(Sarto et al., 2022) use kNN memory for image captioning, enhancing knowledge retrieval from external corpora. Re-ViLM (Yang et al., 2023), built upon the Flamingo (Alayrac et al., 2022), and supports retrieving the relevant knowledge from the external database for zero and in-context few-shot image-to-text generations. Recently, Iscen et al. (Iscen et al., 2023) proposed to equip contrastive vision-text models with the ability to refine their embedding with cross-modal retrieved information from a memory at inference time, which greatly improved their zero-shot predictions. Hu et al. (Hu et al., 2023) presented REVEAL that learns to encode world knowledge into a large-scale memory, and to retrieve from it to answer knowledge-intensive queries, and achieves state-of-the-art results on visual question answering and image captioning. In text-to-image generation, Chen et al. (Chen et al., 2022b) presented Re-Imagen that uses retrieved information to produce high-fidelity

and faithful images, even for rare or unseen entities. RA-CM3 is the first multimodal model that can retrieve and generate mixtures of text and images and exhibits novel capabilities such as knowledge-intensive image generation and multi-modal in-context learning (Yasunaga et al., 2023).

Our multitask framework, RAVEN, extends beyond RA-CM3 by supporting both captioning and VQA, and it diverges from REVEAL (Hu et al., 2023) by attaining retrieval capabilities solely through fine-tuning, eliminating the need for pretraining and additional retrieval-specific parameters; and is adaptable to any base VLM.

## 3 Proposed Approach

### 3.1 RAVEN Framework

Our framework, RAVEN, is illustrated in Figure 1. At a high level, given a multimodal input consisting of images and text, we use a **retriever** to retrieve relevant image-text pairs from a large external memory. Subsequently, we use a pretrained multitask encoder-decoder **VLM** which refers to the retrieved context in addition to the multimodal query and generates a textual output. Importantly, we demonstrate that through short, but efficient, task specific fine-tuning of the base VLM, with concatenated retrieval augmented samples and no additional retrieval-specific trainable parameters, the model acquires retrieval properties which generalizes to multiple tasks. We now describe both these components in detail.

## 3.2 Multimodal Retriever

Our semantic search based retrieval system, relies on the Facebook AI Similarity Search (FAISS) library (Douze et al., 2017). FAISS enables high-dimensional vector indexing within an external memory and facilitates efficient search through an approximate nearest neighbor approach based on a specified similarity measure, such as dot-product similarity. We utilize the publicly available Laion-5B (Schuhmann et al., 2022) image-based index which consists of 5 billion images and corresponding alt text.

To describe the retrieval steps in detail, we first encode the query image using a CLIP-based image encoder (Radford et al., 2021) into a dense vector. Next, we follow the Dense Retrieval method outlined in Karpukhin et al. (2020) to retrieve the top 'k' (k can be specified by the user) image-text pairs by scoring the query (image) and memory data as follows:

$$\text{score(query, memory)} = E(\text{query})^T E(\text{memory}) \tag{1}$$

where E is the CLIP-based image encoder. Finally, we perform Maximum Inner Product Search (MIPS) over the memory to obtain the top 'k' candidate image-text pairs sorted according to the score.

Our retrieval approach ensures that the retrieved samples, which are provided as additional context to the model, along with the query image, are *relevant*, *diverse* and in the *style* of our target datasets. Relevance is easily ensured through sampling based on the top similarity score. However, simply sampling based on relevance score can result in exact or near duplicates resulting in poor performance. To avoid this redundancy and enhance diversity, we exclude near duplicate images. Finally, to use COCO-style captions rather than the noisy image alt text in Laion-5B, we map the retrieved samples from Laion-5B down to the Laion-COCO 600M [1] subset, whose captions are synthetically generated using a BLIP model trained on COCO-style captions. This can result in some missing data due to lack of matches with LAION-COCO 600M and also due to failure of LAION-COCO 600M raw image downloads. Our approach is robust to these missing samples.

## 3.3 Base Vision-Language Model (VLM)

RAVEN relies on a multitask, multimodal encoder-decoder base VLM which can easily leverage additional multimodal context from an external memory.

**Architecture.** For image encoding, we use a ResNet, and for text encoding we use a byte-pair encoding (BPE) to convert the text sequence into a subword sequences, and then embed them into features. We adopt a unified vocabulary encompassing linguistic and visual tokens, incorporating subwords, image codes, and location tokens. The base architecture is the transformer; this serves as the backbone for the encoder-decoder framework. To enhance stability and hasten convergence, the model uses head scaling for self-attention, post-attention layer normalization (LN), and LN following the first layer of FFN. For positional information, separate absolute position embeddings are used for text and images. Notably, we decouple position correlation from token embeddings and patch embeddings, while employing 1D relative position bias for text and 2D relative position bias for images.

**VL Tasks.** All cross-modal tasks are cast as Seq2Seq generation. We focus on 2 popular image-to-text tasks, image captioning and visual question answering (VQA). For image captioning, the model adeptly adopts the Seq2Seq format, generating captions based on both the provided image and the input textual prompt, "What does the image describe?". For VQA, the model takes in the image and the question as inputs, learning to generate accurate responses.

**Need for Retrieval in VL tasks.** Retrieval can benefit performance in VL tasks as contextual information can be crucial for guiding models to accurate answers. Moreover, the retrieval mechanism can mitigate bias by sourcing information from diverse datasets, countering the influence of biased training data. Specifically, in VQA, image content, such as object attributes, strongly correlates with questions and answers, making captions valuable auxiliary information while similar/retrieved images are less informative (Gur et al., 2021). In captioning, additional textual context resembles few-shot inference (Yasunaga et al., 2023).

**Reasons for OFA(Wang et al., 2022b) as a VLM backbone.** We list 4 reasons for choosing OFA rather than alternates like Beit-3 (Wang et al., 2023) and Open Flamingo (Awadalla et al., 2023): *First*, OFA is naturally suited to our approach as it unifies multiple modalities and tasks into a single Seq2Seq model; the multitask backbone is a deliberate design choice that underscores the versatility

---

[1] https://laion.ai/blog/laion-coco/

of our approach and is a foundational element crucial to our model's architecture. *Second*, we can easily endow the model retrieval augmented capabilities through short, but efficient, task specific fine-tuning with no additional trainable parameters. Moreover, we intentionally avoided recent MLLM models like LLaVa or Flamingo which contain an LM to not add additional trainable parameters, remove their in-context learning ability and isolate retrieval capabilities within an encoder-decoder backbone, a first in the field. *Third*, the codebase is open source, modular and easy to extend. *Finally*, the base OFA model is not very large (182M parameters) given our compute and finance limitations, but sufficient to demonstrate the benefits of our framework.

## 4 Experiments

In this section, we evaluate the performance of our approach under the fine-tuning setting on various image captioning and VQA benchmarks. We aim to demonstrate the benefits of retrieval augmentation on the generated captions and answers through retrieving relevant knowledge from a large external non-overlapping database with the fine-tuning datasets. Our experiments show clear benefits of our approach compared to non-retrieval baselines. Furthermore, the performance is competitive with similarly sized models, and even exceeds the performance of existing widely used captioning and VQA models several magnitudes larger.

### 4.1 Training Setup

#### 4.1.1 Data

We make use of an external memory and task specific fine-tuning datasets in our implementation. For captioning, we use the MSCOCO 2014 Karpathy Splits for fine-tuning and NoCaps for a zero-shot evaluation. For VQA, we use the VQA v2 dataset augmented with VG-QA questions during fine-tuning. We use Laion-5B index as our external memory and map down to Laion-COCO 600M subset to retrieve image-caption pairs. The datasets are summarized in Table 1 and 2. Notably, unlike prior work, we ensure the fine-tuning datasets and external memory do not have any overlap, to realize the true benefits of retrieval augmentation in practical settings.

**Missing Samples:** Retrieved data can be missing for 2 reasons: (1) lack of matches of the Laion-5B retrieved samples with the Laion-COCO 600M

| Dataset | Split | # of images (original) | # of images (caption) | # of images (caption + image) | Size - w or w/o retrieval |
|---|---|---|---|---|---|
| MSCOCO | train | 113287 | 108780 | 107800 | 37G / 64G |
| Karpathy | val | 5000 | 4776 | 4725 | 330M / 573M |
| Split (2014) | test | 5000 | 4817 | 4778 | 329M / 576M |
| NoCaps | val | 4500 | 4275 | 4239 | 295M / 512M |

Table 1: Captioning dataset summary

| Split | # of samples | # of images (original) | # of images (caption) | # of images (caption + image) | Size - w or w/o retrieval |
|---|---|---|---|---|---|
| train | 1,358,769 | 121,277 | 116,439 | 115,387 | 106G / 151G |
| val | 10,402 | 2,000 | 1,924 | 1,906 | 653M / 1.2G |
| test-dev | 107,394 | 36,807 | 35,107 | 34,760 | 28G / 50G |
| test-std | 447,793 | 81,434 | 77,856 | 77,098 | 28G / 50G |

Table 2: VQA v2 dataset summary

subset, and (2) raw image download failure. For captioning, we only work on the subset of samples which have both retrieved captions and images. We validate that augmentation with images is not useful, and subsequently decide to only use retrieved captions for augmentation. For VQA, we retain the original dataset, and missing captions are handled with an empty string. This allows us to evaluate our results on the VQA evaluation server. Importantly, the model learns to be robust to samples which may not have corresponding retrieved context at inference; a scenario common in practice.

#### 4.1.2 Implementation

Our retriever uses the off-the-shelf CLIP image encoder (Radford et al., 2021) for both the query and memory encoders. We use FAISS (Douze et al., 2017) to index the external Laion-5B image-based memory and perform MIPS-based top-50 retrieval. We then map down to the Laion-COCO 600M subset ensuring to select, when it exists, the top-1 image (excluding exact or near duplicates), and all associated metadata, including the top caption, all captions and alt text. The retrieved samples are concatenated with the original samples in the TSV file provided as input during the fine-tuning process.

We ensure our fine-tuning process is able to operate in resource constrained settings. We use a lightweight OFA-base (Wang et al., 2022b) model checkpoint of 182M parameters as our multitask VLM. The maximum sequence length is 1024. We fine-tune the model for 8-12 hours, upto 10 epochs, on 4 V100 32GB GPU's. Our implementation is in PyTorch. We increase the max source length from 80 upto 600 to account for the retrieved samples. Otherwise, we rely on the task-specific default hyperparameters in the OFA-base run scripts.

Following the OFA implementation, we optimize

the model with the standard cross-entropy loss. Given an input image i, a prompt t, and an output y, we minimize the loss $L = -\sum_{j=i}^{|y|} \log P_\theta(y_j | y < j, i, t)$ where $\theta$ refers to the model parameters. For inference, we decode using beam search, to enhance the quality of generation. For the VQA task, we employ a trie-based search to only search over a bounded set of vocabulary (top 3129 VQA v2 answers) to prevent labels out of the closed label set during inference.

## 4.2 Evaluation Setup

### 4.2.1 Baselines

We establish baselines to gauge the performance of RAVEN in comparison to various configurations: **Captioning.** (1) Retrieval Only: This baseline involves using the top caption retrieved from the memory as the generated output. It serves as a benchmark to assess the additional benefits gained through fine-tuning the OFA-base model. (2) Zero Shot In-Context Retrieval: During inference, this baseline directly concatenates the retrieved top caption and all captions with the prompt. The objective is to evaluate the model's capacity to leverage retrieved context without any pretraining or fine-tuning. (3) No Retrieved Samples: In this scenario, the model undergoes fine-tuning solely on the target dataset without incorporating any retrieved context. This baseline helps establish a performance reference point.

**VQA.** No Retrieved Samples: Similar to the captioning task, this baseline involves fine-tuning the model exclusively on the target dataset without incorporating any retrieved context.

In all cases, we report performance gains relative to the "No Retrieved Samples" baselines to highlight the efficacy of our proposed approach. Notably, most prior work fail to report this baseline making it challenging to assess the benefits of retrieval augmentation.

Additionally, we provide a comparative analysis by reporting recent baselines and the current State-of-the-Art (SOTA) for both captioning and VQA tasks. This comparative assessment considers performance metrics and the number of parameters, offering a comprehensive view of the landscape and positioning our model within the current state-of-the-art research.

### 4.2.2 Metrics

In evaluating the performance of RAVEN for captioning, we employ two key metrics: BLEU@4 and CIDEr. BLEU@4 measures the quality of generated captions by assessing the overlap of n-grams (in this case, four-grams) between the generated caption and reference captions. Meanwhile, the CIDEr metric gauges the diversity and distinctiveness of generated captions by considering consensus across multiple reference captions.

For the VQA task, we utilize accuracy as the evaluation metric. This measure is computed using the `Eval.ai` server.

### 4.2.3 Ablations

We explore three distinct sets of ablations for both captioning and VQA: text-only, image-only, and combined image and text. To the best of our knowledge, we are the first to comprehensively discern the impact of text and image modalities in retrieval augmented VLMs, providing valuable insights to model practitioners.

**Captioning.** For the text-only ablation, we experiment with various combinations, concatenating one or more of the top caption, all captions, and image alt text. This helps us discern the impact of textual information in isolation. In the image-only ablation, we alter the patch size, doubling it, and employ a horizontal concatenation strategy. If a retrieved image is present, we concatenate it with the query image. In cases where the retrieved image is absent, we duplicate the query image. This analysis provides valuable insights into the model's reliance on visual information alone. For the combined image and text ablation, we adopt a similar approach to the image-only case for processing images. Simultaneously, we concatenate the top caption and all captions to the text prompt. This exploration allows us to understand the synergistic effects of both modalities.

**VQA.** Building on insights gained from the captioning task, where naive image fusion through concatenation proved less useful (see Table 3), we hypothesize that captions serve as good auxiliary information in image-to-text tasks, while similar/retrieved images are less informative, since the content of the image and the objects contained is often very correlated with the question and answer. Therefore, in the VQA ablations, we exclusively consider text concatenation scenarios. This involves combining one or more of the top caption, all captions, and alt text when available. In instances where the retrieved sample is missing, we concatenate with an empty string.

| Retrieval Modality | | # of Parameters | Ablation Description | MSCOCO | | NoCaps |
|---|---|---|---|---|---|---|
| Image | Text | | | BLEU@4 | CIDEr | CIDEr |
| **Our Approach (Image, Text, Image+Text Retrieval)** | | | | | | |
| - | - | - | retrieval only | 0.1905 | 74.98 | 71.68 |
| - | - | 182M | zero shot in-context retrieval with top caption + all captions | 0.3777 | 128.91 | 103.99 |
| - | - | 182M | no retrieved samples | 0.4102 | **137.25** | **106.69** |
| - | ✓ | 182M | top caption | 0.4102 | **138.23* (+0.98)** | 109.76 |
| - | ✓ | 182M | alt text | 0.4125 | 137.19 | 106.81 |
| - | ✓ | 182M | all captions concatenated | 0.4057 | 137.70 | 109.72 |
| - | ✓ | 182M | top caption + all captions | 0.4108 | **138.17* (+0.92)** | **111.00 (+ 4.31)** |
| - | ✓ | 182M | top caption + all captions + alttext | 0.4104 | 138.03 | 109.88 |
| ✓ | - | 182M | image | 0.4087 | 136.95 | 106.22 |
| ✓ | ✓ | 182M | image + top caption + all captions | 0.4081 | 136.85 | 107.28 |
| **Image Captioning Baselines (Fine-tuning)** | | | | | | |
| - | - | 420M | Re-ViLM (base, (Yang et al., 2023)) | 0.378 | 129.1 | 105.2 |
| - | - | 364M | Flamingo (base, re-implementation from (Yang et al., 2023)) | 0.370 | 128.0 | 102.8 |
| - | - | 252M | BLIP$_{CapFilt-L}$ (Li et al., 2022) | 0.404 | 136.7 | 113.2 |
| - | - | 172M | VL-T5 (Cho et al., 2021) | 0.346 | 116.1 | 4.4 |
| - | - | 1.4B | SimVLM (huge, (Wang et al., 2021)) | 0.406 | 143.3 | 110.3 |
| - | - | 5.1B | GIT2 (current SOTA (Wang et al., 2022a)) | 0.432 | 146.4 | 126.9 |

*Gain with respect to the non retrieved baseline is comparable to the only prior work which reported it for the MSCOCO captioning task (Sarto et al., 2022)

Table 3: Fine-tuning evaluation results using cross-entropy optimization on MSCOCO, and NoCaps benchmarks, compared with different image captioning baselines. For NoCaps, we finetune on MSCOCO karpathy train following prior works (Li et al., 2022), and perform zero-shot evaluation. We use the Laion-5B image index mapped down to the Laion-COCO 600M subset as our external memory. We report BLEU@4 and CIDEr scores for different methods and show the gain in the best performing models compared to the non-retrieved baseline.

| # of Parameters | Ablation Description | Test-Dev Accuracy % | | | |
|---|---|---|---|---|---|
| | | number | other | yes/no | overall |
| **Our Approach (Text Retrieval)** | | | | | |
| 182M | no retrieved samples | **58.55** | **67.47** | **90.12** | **75.89** |
| 182M | alttext | 61.10 | 67.94 | 90.10 | 76.29 |
| 182M | alttext + all captions | 57.84 | 67.92 | 90.46 | 76.06 |
| 182M | top caption + all captions | **61.33* (+ 2.78%)** | **68.27* (+0.80%)** | **90.54* (+0.42%)** | **76.75* (+0.86%)** |
| **VQA Baselines (Fine-tuning)** | | | | | |
| 122M | UnifiedVLP (Zhou et al., 2020) | 52.10 | 60.30 | 87.20 | 70.50 |
| 252M | BLIP$_{CapFilt-L}$ (Li et al., 2022) | - | - | - | 78.25 |
| 1.4B | SimVLM (huge, (Wang et al., 2021)) | - | - | - | 80.30 |
| 80B | Flamingo (Alayrac et al., 2022) | - | - | - | 82.00 |
| 55B | PaLI-X (2023) - current SOTA (Chen et al., 2022c) | - | - | - | 86.10 |

*Gain with respect to the non retrieved baseline surpasses that of the only prior work which reported it for the VQA v2 task (Gur et al., 2021)

Table 4: Finetuning evaluation results on VQA v2 benchmarks, compared with the non retrieval VQA baseline. We finetune our method on VQA v2 train split using a subset of the OFA dataset. We report Test-Dev accuracy % from the eval.ai server for different methods.

# 5 Results

## 5.1 Quantitative Analysis

**Captioning.** The results for image captioning, presented in Table 3, reveal notable insights. Baseline comparisons indicate that both the retrieval-only and zero-shot in-context retrieval fall short of the no-retrieved samples baseline, underscoring the value of fine-tuning on the target dataset. The absence of zero-shot in-context retrieval ca-

pabilities may be attributed to the absence of a language model in the transformer-based encoder-decoder VLM architecture. In the text-only ablation, concatenating with the top caption and/or all captions yields optimal performance, demonstrating a gain of nearly 1 CIDEr point on MSCOCO and up to 4 CIDEr points on zero-shot NoCaps. The gain with respect to the non retrieved baseline is comparable to the only prior work which reported it (+1.2 CIDEr score) for the MSCOCO

Figure 2: Examples of the retriever output given a query image.

captioning task (Sarto et al., 2022). This empha-
sizes the valuable contextual information provided
by retrieved captions. However, concatenating with
alt text proves less effective due to its inherent
noise. Both image-only and combined image and
text concatenation exhibit performance below the
non-retrieved baseline, suggesting that retrieved im-
ages and naive concatenation introduce noise rather
than relevant context. In fine-tuning settings, our
model performs competitively with similar-sized
models such as BLIP. Notably, in the zero-shot set-
ting on NoCaps, our model surpasses SimVLM
(1.4B vs 182M parameters), achieving a CIDEr
score of 111.0 compared to 110.3.

**VQA.** Given the limited efficacy observed in
the use of retrieved image for captioning (see Ta-
ble 3), we exclusively explore text augmentation
strategies for VQA. The results, presented in Table
4, align with the captioning outcomes, affirming
the efficacy of text-only augmentation. Notably,
across all question categories, text-only augmenta-
tion yields improvements in accuracy ranging from
0.42% to 2.78%. The gain with respect to the non
retrieved baseline surpasses that of the only prior
work which reported it (+0.36% accuracy) for the
VQA v2 task (Gur et al., 2021). The highest perfor-
mance is achieved through concatenating the top
caption and all captions with the question, while
the addition of alt text introduces noise, resulting
in lower performance. The overall performance of
our model in VQA remains competitive and com-
parable to similar-sized models, underscoring its
robustness in leveraging textual information for
accurate question answering.

## 5.2 Qualitative Analysis

In this section, we present qualitative examples
that elucidate the efficacy and limitations of our
approach.



Figure 3: Examples where RAVEN succeeds in generat-
ing the correct answer.

**Retriever Output.** Figure 2 illustrates the out-
put of the retriever for a given query image. The
retrieved images align with the query image, em-
phasizing relevance. However, Laion-5B's image
alt text is observed to be noisy and differs from
the required COCO-style captions. Mapping down
to synthetically generated BLIP captions from the
LAION-COCO 600M subset, mitigates the style
issue by mimicking the COCO caption style, and
offers more valuable context to the model.

**Incorporating World Knowledge.** Figure
3 demonstrates VQA outputs leveraging world
knowledge. The model adeptly utilizes entity-
rich captions from the retriever to disambiguate
between entities, as seen in the bear image distin-
guishing logs from rocks. Additionally, the model
accurately identifies nuanced details, such as a boy
squatting while playing baseball, by leveraging

relevant context in the captions, such as the term "crouches."

**Retriever Failures.** Despite successes, retrieved context may not consistently contribute to specific questions, particularly when inquiries concern entities not prominently featured in the image. This issue is more pronounced in tasks such as VQA, rather than in captioning, where general knowledge about the image is often sufficient to generate high quality and diverse captions. Illustrated in Figure 4, failure cases for VQA depict relevant but insufficiently informative captions. For instance, captions for an elephant image focus on the foreground elephant, neglecting details about the background mountains and forest. Similarly, captions for a cake image lack information about the cake lifter in the corner.
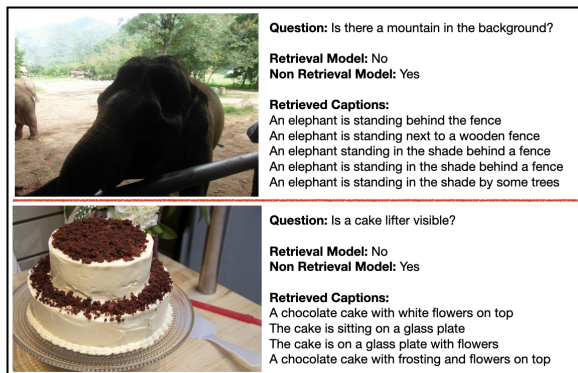


Figure 4: Examples where RAVEN fails in generating the correct answer.

**Multimodal Query Embedding.** Considering scenarios where retrieved context may lack specificity, we propose the joint use of image and text modalities as input to the retriever, when available. Figure 5 demonstrates an example where creating a multimodal query embedding by averaging image and text embeddings separately results in relevant captions addressing both the image and the question. Comprehensive exploration of scenarios where specific entity properties lack corresponding captions is deferred to future work.

## 6 Conclusion and Future Work

To address escalating model size and computational demands, we propose a retrieval augmentation framework, an alternative to storing extensive world knowledge within model parameters. Our contributions introduce a multitask, multimodal retrieval-augmented vision-language model, demonstrating adaptability across multi-
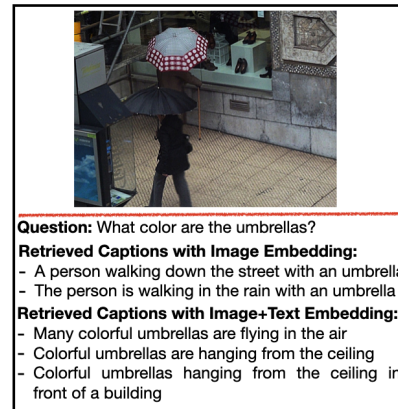


Figure 5: An example depicting the benefits of using a multimodal query embedding (average of image and question embedding). This results in the retrieval of captions relevant to both the image and question.

ple tasks through computationally efficient task-specific fine-tuning. Utilizing concatenated multimodal retrieval-augmented samples from an external non-overlapping memory, without additional trainable parameters, our single model acquires robust retrieval properties. This showcases benefits in both captioning and VQA tasks using a unified approach. Notably, extensive ablations across text, image, and image-text modalities, systematically compared against non-retrieved baselines, provide valuable insights. Our findings underscore that retrieval augmentation, particularly with text in image-to-text tasks, optimally enhances performance, especially in the zero-shot setting.

Future directions involve refining sampling strategies for enhanced diversity, exploring alternative image fusion approaches, and investigating a mixture of experts to afford the model flexibility in leveraging retrieved context. Additionally, we propose extending retrieval over a composite index (image+text) to further optimize performance.

## 7 Limitations

We use a relatively small model to demonstrate performance on 2 tasks. While we acknowledge the demonstrating our approach on more tasks and larger models would be beneficial, we defer this to future work due to compute and financial constraints. RAVEN's current capability to handle diverse tasks like image captioning and VQA within a single model framework already stands as a significant advancement; and is sufficient to demonstrate the benefit of our framework.

# References

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736.

Srikar Appalaraju, Bhavan Jasani, Bhargava Urala Kota, Yusheng Xie, and R. Manmatha. 2021. Docformer: End-to-end transformer for document understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 993–1003.

Srikar Appalaraju, Peng Tang, Qi Dong, Nishant Sankaran, Yichu Zhou, and R. Manmatha. 2024. Docformerv2: Local features for document understanding. *AAAI*, abs/2306.01733.

Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. 2023. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*.

Ali Furkan Biten, Ron Litman, Yusheng Xie, Srikar Appalaraju, and R Manmatha. 2022. Latr: Layout-aware transformer for scene-text vqa. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16548–16558.

Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, T. W. Hennigan, Saffron Huang, Lorenzo Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, and L. Sifre. 2021. Improving language models by retrieving from trillions of tokens. In *International Conference on Machine Learning*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Wenhu Chen, Hexiang Hu, Xi Chen, Pat Verga, and William Cohen. 2022a. Murag: Multimodal retrieval-augmented generator for open question answering over images and text. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5558–5570.

Wenhu Chen, Hexiang Hu, Chitwan Saharia, and William W. Cohen. 2022b. Re-imagen: Retrieval-augmented text-to-image generator. *ArXiv*, abs/2209.14491.

Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. 2022c. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*.

Jaemin Cho, Jie Lei, Hao Tan, and Mohit Bansal. 2021. Unifying vision-and-language tasks via text generation. In *International Conference on Machine Learning*, pages 1931–1942. PMLR.

Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. Instructblip: Towards general-purpose vision-language models with instruction tuning. *Preprint*, arXiv:2305.06500.

Matthijs Douze, Jeff Johnson, and Hervé Jegou. 2017. Faiss: A library for efficient similarity search.

Shir Gur, Natalia Neverova, Chris Stauffer, Ser-Nam Lim, Douwe Kiela, and Austin Reiter. 2021. Cross-modal retrieval augmentation for multi-modal classification. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 111–123.

Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. REALM: Retrieval-augmented language model pre-training. *arXiv preprint arXiv:2002.08909*.

Ziniu Hu, Ahmet Iscen, Chen Sun, Zirui Wang, Kai-Wei Chang, Yizhou Sun, Cordelia Schmid, David A Ross, and Alireza Fathi. 2023. Reveal: Retrieval-augmented visual-language pre-training with multi-source multimodal knowledge memory. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23369–23379.

Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Qiang Liu, Kriti Aggarwal, Zewen Chi, Johan Bjorck, Vishrav Chaudhary, Subhojit Som, Xia Song, and Furu Wei. 2023. Language is not all you need: Aligning perception with language models. *ArXiv*, abs/2302.14045.

Ahmet Iscen, Alireza Fathi, and Cordelia Schmid. 2023. Improving image recognition by retrieving from web-scale image-text data. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19295–19304.

Gautier Izacard and Edouard Grave. 2021. Leveraging passage retrieval with generative models for open domain question answering. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 874–880, Online. Association for Computational Linguistics.

Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and

Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics.

Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. Generalization through Memorization: Nearest Neighbor Language Models. In *International Conference on Learning Representations (ICLR)*.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.

Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pretraining for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning.

R OpenAI. 2023. Gpt-4 technical report. *arXiv*, pages 2303–08774.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Rita Ramos, Desmond Elliott, and Bruno Martins. 2023a. Retrieval-augmented image captioning. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 3648–3663.

Rita Ramos, Bruno Martins, Desmond Elliott, and Yova Kementchedjhieva. 2023b. Smallcap: lightweight image captioning prompted with retrieval augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2840–2849.

Sara Sarto, Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. 2022. Retrieval-augmented transformer for image captioning. In *Proceedings of the 19th International Conference on Content-based Multimedia Indexing*, pages 1–7.

Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. 2022. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294.

Emma Strubell, Ananya Ganesh, and Andrew McCallum. 2019. Energy and policy considerations for deep learning in NLP. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3645–3650, Florence, Italy. Association for Computational Linguistics.

Maria Tsimpoukelli, Jacob Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. 2021. Multimodal few-shot learning with frozen language models. *Proc. Neural Information Processing Systems*.

Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and Lijuan Wang. 2022a. Git: A generative image-to-text transformer for vision and language. *arXiv preprint arXiv:2205.14100*.

Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022b. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International Conference on Machine Learning*, pages 23318–23340. PMLR.

Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, and Furu Wei. 2023. Image as a foreign language: BEiT pretraining for vision and vision-language tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.

Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. 2021. Simvlm: Simple visual language model pretraining with weak supervision. *arXiv preprint arXiv:2108.10904*.

Zhuolin Yang, Wei Ping, Zihan Liu, Vijay Korthikanti, Weili Nie, De-An Huang, Linxi Fan, Zhiding Yu, Shiyi Lan, Bo Li, et al. 2023. Re-vilm: Retrieval-augmented visual language model for zero and few-shot image captioning. *arXiv preprint arXiv:2302.04858*.

Michihiro Yasunaga, Armen Aghajanyan, Weijia Shi, Rich James, Jure Leskovec, Percy Liang, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2023. Retrieval-augmented multimodal language modeling. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org.

Luowei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason Corso, and Jianfeng Gao. 2020. Unified vision-language pre-training for image captioning and vqa. In *Proceedings of the AAAI conference on artificial intelligence*, pages 13041–13049.

Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing

vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*.