## Retrieval Augmented Visual Question Answering with Outside Knowledge

#### Weizhe Lin

Department of Engineering University of Cambridge United Kingdom wl356@cam.ac.uk

## Abstract

Outside-Knowledge Visual Question Answering (OK-VQA) is a challenging VQA task that requires retrieval of external knowledge to answer questions about images. Recent OK-VQA systems use Dense Passage Retrieval (DPR) to retrieve documents from external knowledge bases, such as Wikipedia, but with DPR trained separately from answer generation, introducing a potential limit on the overall system performance. Instead, we propose a joint training scheme which includes differentiable DPR integrated with answer generation so that the system can be trained in an end-to-end fashion. Our experiments show that our scheme outperforms recent OK-VQA systems with strong DPR for retrieval. We also introduce new diagnostic metrics to analyze how retrieval and generation interact. The strong retrieval ability of our model significantly reduces the number of retrieved documents needed in training, yielding significant benefits in answer quality and computation required for training.

#### 1 Introduction

Visual Question Answering (VQA) is a challenging problem that lies at the intersection of Computer Vision, Natural Language Processing, and Information Retrieval. The objective in VQA is to read an image and provide an answer to an accompanying question about the image content. Current approaches to VQA employ deep-learning-based systems to jointly understand images and text.

VQA is particularly challenging when the answer to the question is not directly available in the image. In *Knowledge-based* VQA (KB-VQA), the VQA system must access external knowledge sources to find a correct and complete answer. The Ouside-Knowledge VQA task (OK-VQA) (Marino et al., 2019) consists of questions that requires general knowledge and simple inference to answer (Fig. 1). Such questions are even hard for humans.

## Bill Byrne

Department of Engineering
University of Cambridge
United Kingdom
bill.byrne@eng.cam.ac.uk

Unlike other KB-VQA datasets (e.g. FVQA (Wang et al., 2017)) which provide an associated knowledge base, OK-VQA encourages using any outside knowledge in answering questions.



Question: How many teeth does this animal use to have? Answer: 26

Figure 1: OK-VQA contains questions whose answer cannot be found within the image.

The need to adapt and refresh knowledge sources motivates the study of KB-VQA systems that can extract knowledge from both structured (e.g. ConceptNet (Speer et al., 2017)) and unstructured knowledge representations (e.g. Wikipedia passages). Recent designs (Luo et al., 2021; Gao et al., 2022) approach VQA in two distinct steps: (1) Knowledge Retrieval extracts documents from a large knowledge base; (2) Answer Generation produces an answer from these documents. Knowledge Retrieval can be done via Dense Passage Retrieval (DPR) (Karpukhin et al., 2020), which consists of a question encoder and a document encoder (both Transformer-based) that encode questions and documents into separate dense representations. The DPR system is trained to assign higher scores to documents intended to be helpful in answering questions, so that document sets can be retrieved and passed to Answer Generation.

Knowledge Retrieval based on DPR is powerful but has some readily observed limitations, particularly in model training. Firstly, whether a retrieved document is useful in answering a question cannot be easily determined, even if an answer is provided. Prior work (Qu et al., 2021; Luo et al., 2021) has addressed this problem using "Pseudo Relevance Labels" which are based on whether a document

contains a given answer. However, these are only a weak signal of potential relevance and may encourage DPR to retrieve misleading documents. Secondly, the document retriever and answer generator are trained separately. To ensure that the answer generator sees relevant documents in training, systems can retrieve large numbers of documents ( $\sim$ 50+) (Gao et al., 2022; Gui et al., 2021), but at the cost of slower training and more GPU usage, and also possibly presenting misleading material to the answer generator.

Joint training of the retriever and answer generator offers a solution to these problems. The aim is twofold: (1) to improve the retrieval of documents truly relevant to providing a given answer; and (2) to reject documents with pseudo relevance but not actual relevance.

Retrieval Augmented Generation (RAG) (Lewis et al., 2020) has shown that end-to-end joint training of a DPR-based QA system can outperform baseline two-step systems. A notable feature of RAG is a loss function that incorporates marginalized likelihoods over retrieved documents such that the training score of a document is increased whenever it improves prediction.

However, in preliminary OK-VQA experiments we found that RAG did not perform well. Our investigations found that a good portion of OK-VQA training questions are answerable in closed-book form (i.e. using pre-trained models such as T5 (Raffel et al., 2020)) with information extracted only from the image, with the unintended consequence that the RAG loss function awards credit to documents that did not actually contribute to answering a question. We also found that difficult questions that are unanswerable with the knowledge available to retrieval were more prevalent in OK-VQA than in the Open QA datasets (e.g. Natural Questions (Kwiatkowski et al., 2019)) on which RAG was developed. In both of these scenarios, the RAG loss function leads to counter-intuitive adjustments to the document scores used in training the retrieval model, leading to decreased VQA performance.

Motivated by these findings, we propose a novel neural-retrieval-in-the-loop framework for joint training of the retriever and the answer generator. We formulate a loss function that avoids sending misleading signals to the retrieval model in the presence of irrelevant documents. This formalism combines both pseudo relevance labels and model predictions to refine document scores in training.

We find significantly better performance on OK-VQA compared to RAG. In this paper:

- We present a novel joint training framework Retrieval Augmented Visual Question Answering (RA-VQA) for Knowledge Retrieval and Answer Generation that improves over RAG and two-step baseline systems based on DPR (Karpukhin et al., 2020).
- We investigate visually grounded features transformed into 'language space' and assess their contribution to OK-VQA performance.
- We study the role of document retrieval in KB-VQA and evaluate its interaction with retrieval-augmented generation. We also show that retrieval becomes more efficient in joint training, requiring retrieval of relatively few (~ 5) documents in training.

#### 2 Related Work

Open-domain QA systems. These QA systems are designed to answer questions from datasets such as Natural Questions (Kwiatkowski et al., 2019). The knowledge needed to answer questions can be in pre-trained models (Roberts et al., 2020), knowledge-graphs (KGs) (Lin et al., 2019; Feng et al., 2020; Lv et al., 2020; Saffari et al., 2021) or document collections (Chen et al., 2017; Izacard and Grave, 2021; Guu et al., 2020; Lee et al., 2019; Lewis et al., 2020). In retrieval-based systems, differential retrieval can be combined with extractive question answering, as in REALM (Guu et al., 2020) and ORQA (Lee et al., 2019), as well as with generative answer generation, as in RAG (Lewis et al., 2020).

VQA Systems. Modelling vision and language is central to VQA. Models can aggregate visual and textual features via cross-modality fusion (Yu et al., 2018; Singh et al., 2019; Yu et al., 2019; Jiang et al., 2020; Guo et al., 2021). Systems can also be pre-trained on large vision-and-language collections (Jia et al., 2021) and then fine-tuned for VQA tasks (Tan and Bansal, 2019; Chen et al., 2020; Gan et al., 2020; Li et al., 2020b; Wang et al., 2022; Zhang et al., 2021; Li et al., 2021) with VQA datasets such as VQA 2.0 (Antol et al., 2015).

Knowledge-based VQA Systems. KB-VQA can access both structured data, such as ConceptNet and other KGs (Narasimhan et al., 2018a; Garderes et al., 2020; Li et al., 2020a; Wu et al., 2022; Marino et al., 2021), as well as unstructured data such as Wikipedia passages (Wu et al., 2022; Gao

et al., 2022; Gui et al., 2021). A variety of multimodal approaches have been explored to access external knowledge. ConceptBERT (Garderes et al., 2020) uses attention to aggregate graph node embeddings from ConceptNet. KRISP (Marino et al., 2021) uses a "symbolic knowledge module" to match ConceptNet KG entities with language/visual elements in questions. MAVEx (Wu et al., 2022) uses multiple information sources (Google Images, Wikipedia sentences, and ConceptNet) to validate promising answer candidates. VRR (Luo et al., 2021) uses Google Search in a retriever-reader pipeline to perform open-ended answer generation.

We also note unpublished contemporaneous work on OK-VQA at the time of submission. TRiG (Gao et al., 2022) shows that it is feasible to transform images into textual features for VQA. The features used are similar to those presented here, although without an emphasis on the role of knowledge retrieval. PICa (Yang et al., 2022) 'prompts' GPT-3 with descriptive captions generated from images, and KAT (Gui et al., 2021) exploits an ensemble of DPR, T5, and GPT-3 to improve OK-VQA performance.

## 3 Methodology

We present our RA-VQA framework that consists of: (1) Vision-to-Language Transformation (Sec. 3.1); (2) Weakly-supervised Dense Passage Retrieval (Sec. 3.2); (3) Joint Training of Retrieval and Answer Generation (Sec. 3.3).

## 3.1 Vision-to-Language Transformation

Prior work has established that images can be transformed into text such that large pre-trained language-based Transformers (e.g. BERT (Devlin et al., 2019), GPT-2 (Radford et al., 2019), and T5) can be applied to VQA tasks (Luo et al., 2021; Yang et al., 2022). Systems can be based on straightforward image caption, but we have found improvements by introducing additional visually-grounded features. In RA-VQA, each image is represented by visual objects and their attributes, image caption, and any text strings detected within the image. We use an object detection model VinVL (Zhang et al., 2021) that was pre-trained on large object detection datasets to extract visual elements and their attributes (e.g. color and material).

Formally, for an image I we use VinVL to extract a set of visual objects  $\{o_i\}$ , along with a set of

text attributes for each visual object  $\{a_{i,j}\}$ . Visual objects and their attributes are extracted by VinVL at confidence thresholds 0.8 and 0.6, respectively.

Image captioning is performed to extract relationships and interactions among visual elements such as "a woman holding a knife cuts a cake". The pre-trained captioning model Oscar+ (Zhang et al., 2021) is applied to process visual features extracted from the VinVL model to generate a caption for the image. To answer questions related to text strings in images (e.g. "which language is the book written in?"), Google OCR (Optical Character Recognition) APIs are used to extract text strings from each image.

Hence, a VQA training set  $\{(I, q, S)\}$ , where S is a set of answers to a question q about I, can be transformed into a text-only training set  $\mathcal{T} = \{(x, S)\}$  that we use for RA-VQA. The string x contains all the text features extracted from the image (the question, the textual attributes for each identified visual object, the generated caption, and any OCR'd text), with special tokens marking the start and end of each type of feature (Fig. 2).

## 3.2 Weakly-supervised Dense Passage Retrieval

Dense Passage Retrieval in RA-VQA consists of a query encoder  $\mathcal{F}_q$  and a document encoder  $\mathcal{F}_d$ , both as Transformer-like encoders. The aim is to retrieve K documents from an external knowledge database  $\mathcal{Z} = \{z_i\}_{i=1}^{N^d}$  (e.g. Wikipedia passages) that are expected to be useful for answering a question. DPR encodes questions and documents separately into dense feature vectors  $\mathcal{F}_q(x) \in \mathbf{R}^h$  and  $\mathcal{F}_d(z) \in \mathbf{R}^h$ . A scoring function is used to retrieve documents for each question as the inner product between the representations of x and z

$$r(x,z) = \mathcal{F}_q^{\top}(x)\mathcal{F}_d(z) \tag{1}$$

RA-VQA training aims to maximize r(x, z) when document z is relevant to answering the question. As discussed in Sec. 1, the relevance between q and z cannot be easily obtained and "pseudo relevance labels" serve as a proxy. We use a pseudo relevance function  $H(z, \mathcal{S})$  which is 1 if z contains an answer in  $\mathcal{S}$  (by string match), and 0 otherwise.

For each question-answer pair  $(x, \mathcal{S})$  one positive document  $z^+(x)$  is extracted for training. Inbatch negative sampling is used: all documents in a training batch other than  $z^+(x)$  are considered to be negative for  $(x, \mathcal{S})$  (Karpukhin et al., 2020).

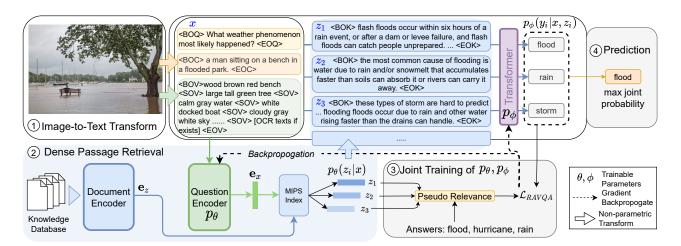


Figure 2: Model overview. (1) Using object detection/image captioning/Optical Character Recognition to transform visual signals into language space. (2) Dense Passage Retrieval retrieves documents that are expected to be helpful from the knowledge database; (3) Training the retriever  $p_{\theta}$  and the answer generator  $p_{\phi}$  together using our proposed RA-VQA loss. (4) The answer with highest joint probability  $p_{\theta}(z_i|x)p_{\phi}(y_i|x,z_i)$  is selected.

Denoting the negative documents as  $\mathcal{N}(x,\mathcal{S})$  and the score of the positive document as  $\widehat{r}^+(x)$  leads to the DPR loss  $\mathcal{L}_{DPR}$ :

$$-\sum_{(x,\mathcal{S})\in\mathcal{T}}\log\frac{\exp\left(\widehat{r}^{+}(x)\right)}{\exp\left(\widehat{r}^{+}(x)\right)+\sum_{z\in\mathcal{N}(x,\mathcal{S})}\exp\left(\widehat{r}(x,z)\right)}$$

# 3.3 RA-VQA: Joint Training of Document Retrieval and Answer Generation

Given a full query string x extracted from the image-question pair (I,q), DPR returns the K highest scoring documents  $\{z_k\}_{k=1}^K$ . The score assigned by the document retriever  $p_{\theta}(\cdot|x)$  to a retrieved document is

$$p_{\theta}(z_k|x) = \frac{\exp(\widehat{r}(x, z_k))}{\sum_{j=1}^{K} \exp(\widehat{r}(x, z_j))}$$
(3)

Open-ended answer generation for each retrieved document  $z_k$  is performed with a generative model, such as T5, with parameters  $\phi$ :

$$y_k = \operatorname*{argmax}_{y} p_{\phi}(y|x, z_k) \tag{4}$$

For each document  $z_k$  retrieved for a training item  $(x, \mathcal{S})$ , we train the answer generator to produce the answer string  $s_k^*$  from the concatenation of x and  $z_k$  (as shown in Fig. 2). We select the most popular human response  $s_k^*$  from  $\mathcal{S}$  such that  $s_k^*$  is contained in  $z_k$ ; in the case that  $z_k$  does not contain any answer, the most popular answer  $s^* \in \mathcal{S}$  is selected  $s_k^* = s^*$ . Through this design, we customize

the generation target  $s_k^*$  for each retrieved document instead of training all  $(x, z_k)$  pairs towards the most popular human response  $s^*$ . This has been proved to improve the system performance (Appendix B.1).

We identify two subsets of the retrieved documents  $\{z_k\}_{k=1}^K$  based on pseudo relevance labels and model predictions:

$$\mathcal{P}^{+}(x,\mathcal{S}) = \{k : y_k = s_k^* \land H(z_k,\mathcal{S}) = 1\}; \\ \mathcal{P}^{-}(x,\mathcal{S}) = \{k : y_k \neq s_k^* \land H(z_k,\mathcal{S}) = 0\}.$$
 (5)

 $\mathcal{P}^+$  are indices of pseudo relevant documents that also help the model generate popular answers whereas  $\mathcal{P}^-$  identifies documents not expected to benefit answer generation. In joint training, we intend to increase the scores of documents in  $\mathcal{P}^+$  while decreasing the scores for those in  $\mathcal{P}^-$ .  $z_k$  will be put into the negative set if it does not contain any answer  $(H(z_k,S)=0)$  and the generation is incorrect  $(y_k \neq s_k^*)$ . This is motivated by our intention to reduce scores for those documents that contain no answers and fail to answer questions.

Formally, joint training of retrieval and answer generation is achieved with a loss  $\mathcal{L}_{RA-VQA}$  that reflects both model predictions and pseudo relevance:

$$-\sum_{(x,\mathcal{S})\in\mathcal{T}} \left( \sum_{k=1}^{K} \log p_{\phi}(s_{k}^{*}|x, z_{k}) + \sum_{k\in\mathcal{P}^{+}(x,\mathcal{S})} \log p_{\theta}(z_{k}|x) - \sum_{k\in\mathcal{P}^{-}(x,\mathcal{S})} \log p_{\theta}(z_{k}|x) \right)$$
(6)

<sup>&</sup>lt;sup>1</sup>There are 5 annotators for each OKVQA question. The popularity of an answer is measured by the number of annotators who voted for it.

<sup>&</sup>lt;sup>2</sup>Note that in this case  $H(z_k, S) = 0$  already implies that  $z_k$  does not contain any answer and thus  $s_k^* = s^*$ .

The first term in the loss improves answer generation from queries and retrieved documents, taken together. The remaining terms affect document retrieval: the second term encourages retrieval of documents that are not only pseudo relevant but also lead to production of correct answers, while the third term works to remove irrelevant items from the top ranked retrieved documents. The in-

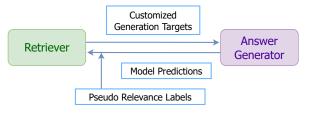


Figure 3: Information flow between the retriever and the answer generator.

formation flow is demonstrated in Fig. 3. Retrieval and generation complement each other in training: pseudo relevance labels and model predictions provide positive and negative signals to improve retrieval, and the improved retrieval leads to improved answer generation by training towards  $s_k^*$ , a customized generation target for each retrieved document  $z_k$ .

### 3.4 RA-VQA Generation

Given an image query (I, q), a full query x is created (Sec. 3.1) and answer generation searches for the answer with the highest joint probability:

$$\{z_k\}_{k=1}^K = \underset{z}{\operatorname{argmax}}^K p_{\theta}(z|x)$$

$$\widehat{y}, \widehat{z} = \underset{y, z_k}{\operatorname{argmax}} p_{\phi}(y|x, z_k) \ p_{\theta}(z_k|x)$$
(7)

Answers reflect both generation and retrieval models and retrieval confidence plays a strong role, unlike some prior work such as Luo et al. (2021).

### 3.5 Pre-Computed FAISS Document Indices

Since repeated computation of embeddings for all documents is costly, we follow Lewis et al. (2020) who find that it is enough to train only the question encoder  $\mathcal{F}_q$  and leave document encoder  $\mathcal{F}_d$  fixed. As shown in Fig. 2, document embeddings are pre-extracted with a pre-trained DPR document encoder. The FAISS system (Johnson et al., 2019) is used to index all document embeddings which enables fast nearest neighbour search with sub-linear time complexity. In training, question embeddings are generated dynamically and documents with

highest scores are retrieved using the pre-computed index.

### 4 Experiments

#### 4.1 Datasets and RA-VQA Configurations

OK-VQA (Marino et al., 2019) is currently the largest knowledge-based VQA dataset. It consists of 14,031 images and 14,055 questions. These questions are split into a training set (9,009 questions) and a test set (5046 questions). In addition to understanding images and questions, external knowledge sources are needed to answer questions.

As outside knowledge we use the knowledge corpus collected by Luo et al. (2021) from Google Search. We use the corpus GS-full which consists of 168,306 documents covering training and test questions. In Appendix B.1 we also report on GS-train, which contains documents relevant to OK-VQA training set questions only.

*Pre-training:* We start with pre-trained versions of BERT-base and T5-large as the document retriever and the answer generator, respectively. The retriever was refined by training it on GS-full under the DPR loss (Equation 2) with pseudo relevance labels released by Luo et al. (2021). The already strong retriever serves as a good starting point for all DPR-based models presented in this paper (including RA-VQA and our replication of baselines in the literature).

*OK-VQA Fine-tuning:* Our **RA-VQA** framework trains the answer generator and the retriever jointly under Equation 6.

We also report on variants of RA-VQA, to investigate the contribution of various model components to overall performance:

**RA-VQA-NoDPR** omits retrieval entirely so that answers are generated by the fine-tuned T5 alone. RA-VQA generation in Equation 7 simplifies to

$$\widehat{y}_{NoDPR} = \operatorname*{argmax}_{y} p_{\phi}(y|x) \tag{8}$$

**RA-VQA-FrDPR** leaves the retriever frozen after pre-training and fine-tunes the generator only.

**RA-VQA-NoPR** is a version of RA-VQA in which document retrieval is trained only with model predictions. The loss function is as Equation 6, but with positive and negative document sets defined

$$\mathcal{P}_{NoPR}^{+}(x,\mathcal{S}) = \{k : y_k = s_k^*\}; \mathcal{P}_{NoPR}^{-}(x,\mathcal{S}) = \{k : y_k \neq s_k^*\}.$$
 (9)

**RA-VQA-NoCT** replaces the customized generation targets by the single most popular response

 $(s_k^*)$  becomes  $s^*$  in Equation 6) so that the generator is trained to produce the same answer from every retrieved document.

#### 4.2 Evaluation

The following metrics are applied to assess the quality of individual answers generated and documents retrieved. Average scores are then computed over the evaluation set. The average of 3 runs with different seeds is reported.

#### 4.2.1 Answer Evaluation

**VQA Score**: We follow Marino et al. (2019) to compute VQA Scores using pre-processed human annotations S:

$$VQAScore(y, S) = \min\left(\frac{\#_{S}(y)}{3}, 1\right), \quad (10)$$

where  $\#_{\mathcal{S}}(y)$  is the number of annotators who answered y. This score ensures that a model is partially rewarded even if it generates one of the less popular answers from amongst the human responses.

**Exact Match (EM)** treats annotated answers equally:  $EM(y, S) = \min(\#_S(y), 1)$ .

#### 4.2.2 Retrieval Evaluation

Following Luo et al. (2021), we use pseudo relevance to ascertain whether the retrieved documents are relevant to the response. It concerns pseudo relevance instead of actual relevance but is still a reasonable metric for retrieval evaluation.

**Pseudo Relevance Recall (PRRecall)** @ $\mathbf{K}$  measures how likely the retrieved K documents contains at least one positive document:

$$PRRecall@K = \min \left( \sum_{k=1}^{K} H(z_k, S), 1 \right). \quad (11)$$

#### 4.2.3 Integrated System Evaluation

The above methods evaluate retrieval and answer generation as separate processes. We propose additional metrics that assess how the two processes behave in an integrated VQA system.

The **Hit Success Ratio** (**HSR**) counts questions that require external knowledge to answer:

$$HSR = \mathbb{1}\{\widehat{y} \in \mathcal{S} \land \widehat{y}_{NoDPR} \notin \mathcal{S}\}.$$
 (12)

HSR reflects the value of incorporating external documents into answer generation.

By contrast, **Free Success Rate (FSR)** counts questions that can be answered without external knowledge.

$$FSR = \mathbb{1}\{\widehat{y} \in \mathcal{S} \land \widehat{y}_{NoDPR} \in \mathcal{S}\}. \tag{13}$$

A high FSR suggests a model can generate correct answers 'freely' without being distracted by retrieved documents if they are not needed.

We also assess performance as a function of the number of documents retrieved during training and testing,  $\mathbf{K}_{\text{train}}$  and  $\mathbf{K}_{\text{test}}$ . In practice,  $K_{\text{train}}$  has the greater effect on GPU usage, since a large  $K_{\text{train}}$  requires at least  $K_{\text{train}}$  forward passes for each question and an Adam-like optimizer must compute and store the associated gradients (Kingma and Ba, 2015). In contrast, GPU memory required during testing is significantly less, as there is no optimizer involved. We are in particular interested in the ability of knowledge-augmented systems that can be robustly trained with small  $K_{\text{train}}$  while yielding improved performance with large  $K_{\text{test}}$ .

## 4.3 Baseline Systems

## 4.3.1 Retrieval Augmented Generation

RAG (Lewis et al., 2020) is based on DPR and an answer generator that are trained jointly by approximately marginalizing the probability of y over the retrieved documents. In the notation of Sec. 3:

$$p_{RAG}(y|x) \approx \sum_{k=1}^{K} p_{\phi}(y|x, z_k) p_{\theta}(z_k|x) \qquad (14)$$

The answer generator and the retriever are jointly trained by optimizing the RAG loss:  $-\sum_{(x,\mathcal{S})\in\mathcal{T}}\log\left(p_{RAG}(s^*|x)\right).$  The rationale is that  $p_{\theta}(z_k|x)$  will increase if  $z_k$  has a positive impact on answer generation (Lewis et al., 2020). We consider RAG as an important baseline and have carefully replicated its published implementation<sup>3</sup>.

#### **4.3.2** Baseline Systems in the Literature

We compare against the published OK-VQA results from systems described in Sec. 2: **ConceptBERT**, **KRISP**, **MAVEx**, and **VRR**. We also report performance against unpublished (non peer-reviewed) systems **TRiG**<sup>4</sup>, **PICa**, and **KAT**. **TRiG** uses a similar image-to-text transform as this work, so to enable fair comparison with our model we replicate their knowledge fusing method with our features. Baseline results are reported in Table 1; baseline results marked \* are our own. **TRiG**\*, our own implementation of TRiG, concatenates *K* encoder outputs for the decoder to use in generation.

We make some particular observations. Our TRiG\* improves over the results released in its

<sup>&</sup>lt;sup>3</sup>The authors released RAG in huggingface.

<sup>&</sup>lt;sup>4</sup>At the time of submission, TRiG has not been published.

Model	T5	GPT-3	$K_{ m train}$	$K_{ m test}$	Knowl. Src.	PRRecall	HSR / FSR	H/F	EM	VQA
ConceptBERT	×	×	-	-	С					33.66
KRISP	X	×	-	-	C + W					38.35
VRR	×	×	100	100	GS					45.08
MAVEx	×	×	-	-	W + C + GI					39.40
KAT-T5	$\checkmark$	×	40	40	W					44.25
TRiG	$\checkmark$	×	5	5	W				49.21	45.51
TRiG	$\checkmark$	×	100	100	W				53.59	49.35
TRiG-Ensemble	$\checkmark$	×	100	100	W				54.73	50.50
TRiG*	$\checkmark$	×	5	5	GS				52.79	48.32
RAG*	$\checkmark$	×	5	5	GS	82.34	12.28 / 40.24	0.31	52.52	48.22
RA-VQA (Ours)	$\checkmark$	×	5	5	GS	82.84	16.75 / 41.97	0.40	58.72	53.81
RA-VQA (Ours)	$\checkmark$	×	5	50	GS	96.55	17.32 / 42.09	0.41	59.41	54.48
Ablation Study										
RA-VQA-FrDPR	<b>√</b>	×	5	5	GS	81.25	15.01 / 40.76	0.37	55.77	51.22
RA-VQA-NoPR	$\checkmark$	×	5	5	GS	77.67	15.97 / 41.83	0.38	57.80	52.98
RA-VQA-NoCT	$\checkmark$	×	5	5	GS	83.77	14.55 / 42.96	0.33	57.51	52.67
GPT-3-based Systems	(>17.	5 Billion F	Paramete	rs)						
PICa	×	<b>√</b>	-	-	GPT-3					48.00
KAT-Knowledge-T5	$\checkmark$	$\checkmark$	40	40	W + GPT-3					51.97
KAT-Ensemble	$\checkmark$	$\checkmark$	40	40	W + GPT-3					54.41

Table 1: RA-VQA vs. Baseline Systems. Knowledge Sources: ConceptNet; Wikipedia; Google Search; Google Images; GPT-3 closed book knowledge. H/F: HSR to FSR ratio. PRRecall, HSR, FSR, and EM are reported in percentage (%). PRRecall is reported at the corresponding  $K_{\text{test}}$ .

paper (VQA Score of 48.32 vs 45.51) at  $K_{\text{train}} = K_{\text{test}} = 5$ ; TRiG and TRiG Ensemble both benefit from more retrieved documents in training and testing ( $K_{\text{train}} = K_{\text{test}} = 100$ ), although at great computational cost. Best performance with KAT-T5 and VRR similarly requires large document collections in training and in test.

We include results from GPT-3 based systems because they are amongst the best in the literature, but we note that GPT-3 is so much bigger than T5 (175 billion parameters in GPT-3 v.s. 770 million in T5-large) that simply switching a system implementation from T5 to GPT-3 can give significant improvements: KAT-T5 achieved a 44.25 VQA Score while ensembling it with GPT-3 yields 54.41; and GPT-3 alone already achieved good performance with prompting (PICa with 48.00 VQA Score). Our RA-VQA system is based on T5, but we still find competitive results even in comparison to systems incorporating GPT-3 (54.48 vs 54.41 of KAT-Ensemble).

## 4.4 RA-VQA Performance Analysis

We find that RA-VQA matches or improves over all baseline systems with a VQA Score of 54.48. This is with a configuration of  $K_{\text{train}} = 5$  and  $K_{\text{test}} = 50$ , thus validating our claim that RA-VQA can use a large number of retrieved documents in testing (50) while using relatively few retrieved documents

in training (5). We find that reducing the number of retrieved documents in test ( $K_{\rm test}=5$ ) reduces the VQA Score, but still yields performance better than all baselines except the KAT ensemble.

We also find that RA-VQA performs well relative to GPT-3 baselines. RA-VQA yields a higher VQA score than KAT-Knowledge-T5 (54.48 vs. 51.97) and matches the KAT-Ensemble system. We emphasize that RA-VQA is significantly smaller in terms of parameters (and in model pre-training data) than these GPT-3 based systems and that training RA-VQA requires much less memory ( $K_{\text{train}} = 5 \text{ vs } K_{\text{train}} = 40$ ).

4.4.1 Contributions of Query Features and DPR to Overall Performance

Model	Q	О	A	С	T	VQA Score
RA-VQA-NoDPR	<b>√</b>	×	×	×	×	28.05
RA-VQA-NoDPR	$\checkmark$	$\checkmark$	×	×	×	40.95
RA-VQA-NoDPR	$\checkmark$	$\checkmark$	$\checkmark$	×	×	42.14
RA-VQA-NoDPR	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	×	45.31
RA-VQA-NoDPR	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	46.16
RA-VQA-FrDPR	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	51.22
RA-VQA	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	53.81

Table 2: Ablation study on input features and system configurations:  $\underline{\mathbf{Q}}$  uestions;  $\underline{\mathbf{Q}}$  bjects;  $\underline{\mathbf{A}}$  ttributes associated with objects;  $\underline{\mathbf{C}}$  aptions; visible  $\underline{\mathbf{T}}$  ext from OCR. K=5 in RA-VQA and RA-VQA-FrDPR.

A detailed ablation study on input features is presented in Table 2. As shown, the T5 model

fine-tuned on OK-VQA achieves a 28.05 VQA Score. The VQA Score increases to 46.16 as objects, object attributes, image captions, and OCR texts are incorporated into RA-VQA-NoDPR. With 5 retrieved documents, RA-VQA-FrDPR yields a 51.22 VQA Score, with a further improvement (53.81 VQA Score) in full training of retrieval and answer generation, confirming that outside knowledge is needed to answer OK-VQA questions.

## 4.4.2 Benefits of Integrated Training

Joint training is a key benefit of our proposed RA-VQA framework: model predictions combine with pseudo relevance labels to improve retrieval, and the resulting improved retrieval in turn provides customized answer generation targets. To quantify these effects, we take RA-VQA-FrDPR as a starting point (Table 1). Comparing it with other RA-VQA models suggests that DPR training in itself is necessary, as using only pre-trained DPR (RA-VQA-FrDPR) leads to weaker VQA Score (51.22). Using model predictions alone in joint DPR training (RA-VQA-NoPR) leads to a higher VQA Score (52.98 vs 51.22), but a significantly lower PRRecall (77.67% vs 81.25%). The model decides to remove some pseudo relevant documents but achieves better performance. This points to a potential problem that can arise. Pseudo relevance is only an imperfect indication of true relevance and so is not an ideal criteria on its own. Training DPR to retrieve pseudo relevant documents could result in misleading documents being used in answer generation.

Using both pseudo relevance labels and model predictions in DPR training (RA-VQA) improves VQA Score to 53.81 and notably improves PRRecall to 82.84%. Including pseudo relevance ensures that potentially useful documents are retained, even while the generator is still learning to use them.

We also note that when generation targets are not customized for each retrieved document (RA-VQA-NoCT), VQA Score drops by 1.14 relative to RA-VQA, showing that customized generation targets play an important role in the overall system: by training the model to extract the reliable answers available in retrieved documents, answer generation and retrieval are both improved.

#### 4.4.3 Interaction of Retrieval and Generation

Table 1 also reports our investigation into the interaction between document retrieval and answer gen-

eration. In comparing RA-VQA-FrDPR (frozen DPR) to RA-VQA, we see that joint training of DPR yields not only improved EM but also significantly higher HSR (+1.74%) and FSR (+1.21%). Manual inspection of OK-VQA reveals that there are many general knowledge questions. For example, document retrieval is not needed to answer the question "Is this television working?" in reference to a picture of a broken television lying in a field. A high FSR indicates good performance on such questions. By contrast, a high HSR reflects the ability to use document retrieval to answer the questions that truly require external documents.

Both EM and HSR are further improved for  $K_{\rm test} = 50$  in RA-VQA, with little change in FSR. The increased HSR to FSR ratio (0.41 vs. 0.40) indicates that RA-VQA is using these additional retrieved documents to answer the questions that need outside knowledge.

HSR and FSR also explain the relatively weak performance of RAG\*. We see that although RAG\* and RA-VQA-FrDPR have similar FSRs, RAG\* has higher PRRecall but lower HSR (by -2.73%). This suggests RAG\*'s DPR model is not well matched to its answer generator. The result is that retrieved documents remain unexploited. In manual examination of gradients of document scores in training, we find anecdotally that adjustments to document scores are often counter-intuitive: documents that do not contain answers can still have their scores upvoted if the answer generator happens to find a correct answer by relying only on the ability of T5 model. This works against a model's ability to find answers in retrieved documents even when those documents are relevant.

#### **4.4.4** Effects of $K_{\text{train}}$

As noted, retrieving a large collection of documents in training is costly (large  $K_{\rm train}$ ). Fig. 4 shows that RA-VQA can be trained with relatively few retrieved documents ( $K_{\rm train}=5$ ). We gradually increase  $K_{\rm train}$  while fixing  $K_{\rm test}=K_{\rm train}$  (dash lines) and  $K_{\rm test}=50$  (solid lines). RA-VQA achieves consistent performance ( $\sim$ 54.4 VQA Score) at  $K_{\rm train}\geq 5$  and  $K_{\rm test}=50$ , which suggests that our joint training scheme is able to gather most useful knowledge into a top-50 list even when the model is trained to retrieve fewer documents. This is not the case for the frozen DPR systems which require increasing  $K_{\rm train}$  to obtain best performance. RA-VQA's superior performance shows that joint training of retrieval and generation yields clear ben-

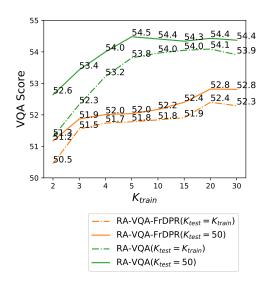


Figure 4: VQA Scores against  $K_{\text{train}}$ . Dashed line:  $K_{\text{test}} = K_{\text{train}}$ ; solid line:  $K_{\text{test}} = 50$ . Our proposed model achieves the best performance when additional documents are retrieved in test ( $K_{\text{test}} = 50$ ). This holds even for models trained to retrieve fewer documents.

efits in computation and answer quality.

We also note that  $K_{\text{train}}=5$  is an optimal design choice that strikes a balance between training computational cost and system performance. The green curves at  $K_{\text{train}}<5$  also suggest that it is beneficial to include at least 5 documents in training. This is because  $K_{\text{train}}\geq 5$  offers a good PRRecall (over 80%), which provides documents of higher quality for training the system.

## 4.5 Additional Analyses

We redirect readers to other interesting analyses (e.g. an analysis of computational cost) and case studies in Appendices B-E.

In addition, in Appendix F, we evaluate our proposed framework on another popular Knowledge-based VQA dataset, FVQA (Fact-based VQA) (Wang et al., 2017). Similarly, the RA-VQA framework with joint training achieves better results over the baseline systems with a frozen DPR component, showing the generalizability of our proposed framework.

#### 5 Conclusion

Retrieval-Augmented Visual Question Answering is a novel modelling framework for integrated training of DPR and answer generation. We have evaluated RA-VQA on the OK-VQA task and we find significantly better performance than the independent training of component system. Through diagnostic metrics such as HSR and FSR we analysed the interaction between retrieval and generation,

and have also shown how RA-VQA's gains arise relative to other approaches, such as RAG. As a further practical benefit, we found that RA-VQA can be used with larger numbers of retrieved documents than were used in system training, yielding computational savings without sacrificing performance.

The code for this paper will be released at https://github.com/LinWeizheDragon/Retrieval-Augmented-Visual-Question-Answering.

## 6 Acknowledgement

W. Lin was supported by a Research Studentship funded by Toyota Motor Europe (RG92562(24020)). We thank our colleagues, Daniel Olmeda Reino (Toyota Motor Europe) and Jonas Ambeck (Toyota Motor Europe), who provided insight and expertise in this project.

We thank Zhilin Wang (University of Washington) and Alexandru Coca (University of Cambridge) for comments that greatly improved the manuscript. We would also like to thank all the reviewers for their knowledgeable reviews.

#### 7 Limitations

One possible limitation is that some relevant information (such as relative positioning of objects in the image) could be lost in transforming images independently of the information being sought. Extracting visual features based on queries could be a natural next step, although query-specific processing of the image collection would be computationally expensive.

We selected the Google Search corpus (Luo et al., 2021) as the knowledge base for our question answering system. Its advantages are that it is large, openly available, and can be readily used to replicate the results in this paper. However some visual question types (e.g. 'Is the athlete right or left handed?') could plausibly require both complex reasoning and more closely relevant documents from additional knowledge sources (such as Wikipedia). Our system may be limited with respect to these considerations.

#### References

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433.

- Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. 2007. Dbpedia: A nucleus for a web of open data. In *The semantic web*, pages 722–735. Springer.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading Wikipedia to answer opendomain questions. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1870–1879, Vancouver, Canada. Association for Computational Linguistics.
- Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. Uniter: Universal image-text representation learning. In *European conference on computer vision*, pages 104–120. Springer.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics. ACL.
- Yanlin Feng, Xinyue Chen, Bill Yuchen Lin, Peifeng Wang, Jun Yan, and Xiang Ren. 2020. Scalable multi-hop relational reasoning for knowledge-aware question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1295–1309, Online. Association for Computational Linguistics.
- Zhe Gan, Yen-Chun Chen, Linjie Li, Chen Zhu, Yu Cheng, and Jingjing Liu. 2020. Large-scale adversarial training for vision-and-language representation learning. *Advances in Neural Information Processing Systems*, 33:6616–6628.
- Feng Gao, Qing Ping, Govind Thattai, Aishwarya Reganti, Ying Nian Wu, and Prem Natarajan. 2022. Transform-retrieve-generate: Natural language-centric outside-knowledge visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5067–5077.
- François Garderes, Maryam Ziaeefard, Baptiste Abeloos, and Freddy Lecue. 2020. Conceptbert: Concept-aware representation for visual question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*, pages 489–498.
- Liangke Gui, Borui Wang, Qiuyuan Huang, Alex Hauptmann, Yonatan Bisk, and Jianfeng Gao. 2021. Kat: A knowledge augmented transformer for vision-and-language. arXiv preprint arXiv:2112.08614.
- Dalu Guo, Chang Xu, and Dacheng Tao. 2021. Bilinear graph networks for visual question answering. *IEEE Transactions on Neural Networks and Learning Systems*.

- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. Realm: Retrieval-augmented language model pre-training. In *Proceedings of the 37th International Conference on Machine Learning*, ICML'20. JMLR.org.
- 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.
- Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. 2021. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, pages 4904–4916. PMLR.
- Huaizu Jiang, Ishan Misra, Marcus Rohrbach, Erik Learned-Miller, and Xinlei Chen. 2020. In defense of grid features for visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10267–10276.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547.
- 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)*, pages 6769–6781, Online. Association for Computational Linguistics.
- Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466.
- Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. 2019. Latent retrieval for weakly supervised open domain question answering. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6086–6096, Florence, Italy. Association for Computational Linguistics.
- 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.
- Guohao Li, Xin Wang, and Wenwu Zhu. 2020a. Boosting visual question answering with context-aware knowledge aggregation. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 1227–1235.
- Wei Li, Can Gao, Guocheng Niu, Xinyan Xiao, Hao Liu, Jiachen Liu, Hua Wu, and Haifeng Wang. 2021. UNIMO: Towards unified-modal understanding and generation via cross-modal contrastive learning. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 2592–2607, Online. Association for Computational Linguistics.
- Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. 2020b. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *European Conference on Computer Vision*, pages 121–137. Springer.
- Bill Yuchen Lin, Xinyue Chen, Jamin Chen, and Xiang Ren. 2019. KagNet: Knowledge-aware graph networks for commonsense reasoning. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2829–2839, Hong Kong, China. Association for Computational Linguistics.
- Man Luo, Yankai Zeng, Pratyay Banerjee, and Chitta Baral. 2021. Weakly-supervised visual-retriever-reader for knowledge-based question answering. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6417–6431, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Shangwen Lv, Daya Guo, Jingjing Xu, Duyu Tang, Nan Duan, Ming Gong, Linjun Shou, Daxin Jiang, Guihong Cao, and Songlin Hu. 2020. Graph-based reasoning over heterogeneous external knowledge for commonsense question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8449–8456.
- Kenneth Marino, Xinlei Chen, Devi Parikh, Abhinav Gupta, and Marcus Rohrbach. 2021. Krisp: Integrating implicit and symbolic knowledge for opendomain knowledge-based vqa. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14111–14121.
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/CVF Con*ference on Computer Vision and Pattern Recognition, pages 3195–3204.

- Medhini Narasimhan, Svetlana Lazebnik, and Alexander Schwing. 2018a. Out of the box: Reasoning with graph convolution nets for factual visual question answering. *Advances in neural information processing systems*, 31.
- Medhini Narasimhan, Svetlana Lazebnik, and Alexander Schwing. 2018b. Out of the box: Reasoning with graph convolution nets for factual visual question answering. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.
- Chen Qu, Hamed Zamani, Liu Yang, W Bruce Croft, and Erik Learned-Miller. 2021. Passage retrieval for outside-knowledge visual question answering. In Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, pages 1753–1757.
- 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.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.
- Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. How much knowledge can you pack into the parameters of a language model? In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5418–5426, Online. Association for Computational Linguistics.
- Amir Saffari, Armin Oliya, Priyanka Sen, and Tom Ayoola. 2021. End-to-end entity resolution and question answering using differentiable knowledge graphs. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4193–4200, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. 2019. Towards vqa models that can read. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8317–8326.
- Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Thirty-first AAAI conference on artificial intelligence*.
- Hao Tan and Mohit Bansal. 2019. LXMERT: Learning cross-modality encoder representations from transformers. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages

- 5100–5111, Hong Kong, China. Association for Computational Linguistics.
- Niket Tandon, Gerard De Melo, and Gerhard Weikum. 2017. Webchild 2.0: Fine-grained commonsense knowledge distillation. In *Proceedings of ACL 2017*, *System Demonstrations*, pages 115–120.
- Peng Wang, Qi Wu, Chunhua Shen, Anthony Dick, and Anton Van Den Hengel. 2017. Fvqa: Fact-based visual question answering. *IEEE transactions on pattern analysis and machine intelligence*, 40(10):2413–2427.
- Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. 2022. SimVLM: Simple visual language model pretraining with weak supervision. In *International Conference on Learning Representations*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 38–45, Online. Association for Computational Linguistics.
- Jialin Wu, Jiasen Lu, Ashish Sabharwal, and Roozbeh Mottaghi. 2022. Multi-modal answer validation for knowledge-based vqa. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 2712–2721.
- Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu, Yumao Lu, Zicheng Liu, and Lijuan Wang. 2022. An empirical study of gpt-3 for few-shot knowledge-based vqa. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 3081–3089.
- Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. 2019. Deep modular co-attention networks for visual question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6281–6290.
- Zhou Yu, Jun Yu, Chenchao Xiang, Jianping Fan, and Dacheng Tao. 2018. Beyond bilinear: Generalized multimodal factorized high-order pooling for visual question answering. *IEEE transactions on neural networks and learning systems*, 29(12):5947–5959.
- Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. 2021. Vinvl: Revisiting visual representations in vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5579–5588.

Zihao Zhu, Jing Yu, Yujing Wang, Yajing Sun, Yue Hu, and Qi Wu. 2020. Mucko: Multi-layer cross-modal knowledge reasoning for fact-based visual question answering. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, pages 1097–1103. International Joint Conferences on Artificial Intelligence Organization. Main track.

#### A Training Details and Artifacts

Adam (Kingma and Ba, 2015) was used in this paper. In DPR pre-training, the retriever was trained for 6 epochs with a constant learning rate  $10^{-5}$ . In RA-VQA training (including RAG\*), the learning rates are  $10^{-5}$  for the retriever, and  $6 \times 10^{-5}$  for the answer generator, linearly decaying to 0 after 10 epochs. In the training of RA-VQA-NoDPR and TRiG\*, the initial learning rate is  $6 \times 10^{-5}$ . Empirically, the checkpoints at epoch 6 were used in testing. All experiments were run on Nvidia A-100 GPU clusters. With  $K_{\rm train} = 5$ , the RA-VQA training takes around 5 hours (10 epochs) while testing takes 5 minutes. The time cost increases as  $K_{\rm train}$  increases, approximately linearly.

Pre-trained model parameters (e.g. T5-large and BERT-base) are provided by huggingface (Wolf et al., 2020) accompanied by Python libraries (under Apache License 2.0). FAISS (Johnson et al., 2019) is under MIT License.

#### **B** Supplementary Tables and Figures

Limited by space available, we present supplementary tables and figures in this section, offering more findings to readers.

#### **B.1** Full Version of Table 1

We provide the full version of Table 1 in Table 8. Some more discussions about customized generation target in joint training is provided.

As noted, RA-VQA improves retrieval with the feedback of model predictions, and in turn the improved retrieval leads to improved answer generation by training towards  $s_k^*$ , a customized generation target for each retrieved document  $z_k$ . We remove this interaction from RA-VQA models by enforcing  $s_k^* = s^*$  (the most popular human response), independent of the retrieved  $z_k$ . The ablated models are denoted with a \*-NoCT suffix.

As shown in Table 8, customizing generation targets for each retrieved  $z_k$  in training yields performance boost for both RA-VQA-FrDPR and RA-VQA, showing that this supervision signal is benefi-

Models	K	K = 5		K = 10 $K =$		= 20 $K = 50$		$K_{\text{test}} = 5$		
	P	R	P	R	P	R	P	R	EM	VQA Score
VRR (Luo et al., 2021)	- 51.00	80.40	-	88.55				97.11	-	42.54
RA-VQA-FrDPR RA-VQA	51.82 <b>57.39</b>	81.25 <b>82.84</b>	49.20 <b>54.83</b>	88.51 <b>89.00</b>	45.98 <b>51.48</b>	92.98 <b>93.62</b>		, 0., 0	55.77 <b>58.72</b>	51.22 <b>53.81</b>

Table 3: Comparing retrieval performance of VRR and our RA-VQA models. The same knowledge corpus (GS-full) was used. P: Pseudo Relevance Precision; P: Pseudo Relevance Recall; EM: Exact Match. P under K=5 refers to PRPrec@5. VRR was trained on  $K_{\text{train}}=100$ , while RA-VQAs were trained on  $K_{\text{train}}=5$ .

cial to overall system performance. We also notice that the improvement to RA-VQA (+1.14 VQA Score) is larger compared to RA-VQA-FrDPR (+0.56 VQA Score), showing that customizing the generation target brings more benefits when the retrieval is improved within our proposed RA-VQA joint training framework. This further confirms that the two components, retrieval and answer generation, complement each other bi-directionally.

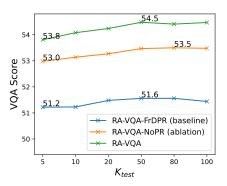
#### **B.2** Retrieval Performance of RA-VQA

In addition to Pseudo Relevance Recall (PRRecall) introduced in the paper, we further evaluate retrieval performance with **Pseudo Relevance Precision (PRPrec)**@**K**, which is calculated as the rate of pseudo positive documents in all the *K* documents retrieved for a question:

$$PRPrec@K = \frac{1}{K} \sum_{k=1}^{K} H(z_k, S) \qquad (15)$$

where  $H(\cdot)$  is the pseudo relevance function introduced in Sec. 3.2.

The success of our RA-VQA model can be further explained by Table 3. As expected, RA-VQA-FrDPR (pre-trained DPR) achieves similar retrieval performance as VRR (Luo et al., 2021) since they are both based on DPR and are trained with the same pseudo-relevance-based labels. Our proposed RA-VQA, with a substantial improvement in Recall over RA-VQA-FrDPR (82.84 PRRecall@5 vs 81.25 PRRecall@5), achieves significantly higher Precision (57.39 PRPrec@5 vs 51.82 PRPrec@5). This also yields substantial improvements to both EM (+3.05%) and VQA Score (+2.59%). This suggests that training the retriever jointly presents more potentially relevant documents to answer generation, improving the quality of the top-ranked documents.



(a) VQA Score vs  $K_{\text{test}}$ 

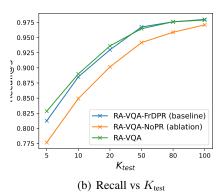


Figure 5: Comparison of model performance as more documents are retrieved in testing. These models are all trained with  $K_{\rm train}=5$ . In RA-VQA full joint training (green), combining model predictions with pseudo relevance labels yields higher PRRecall at low  $K_{\rm test}$ , showing that full joint training improves retrieval; RA-VQA-NoPR (orange), which uses only model predictions in training, achieves a higher VQA Score with lower Pseudo Relevance Recall compared to the RA-VQA-FrDPR with frozen DPR in training (blue), which suggests that Pseudo Relevance is only an approximate measurement of actual relevance.

# **B.3** Effects of Retrieving More Documents in Test

Fig. 5 presents the change of VQA Score and PRRecall as additional documents are retrieved in test (increasing  $K_{\text{test}}$ ).

PRRecall is improved dramatically as  $K_{\text{test}}$  increases from 5 to 50, after which only marginal improvement is observed. Similarly, the VQA Score of these models is improved as more documents are presented in test, and the performance peaks at  $K_{\text{test}} \sim 50$ . This suggests that including more additional documents in test is more likely to include the truly relevant document to help answer the question yet along with more distracting and misleading documents.

RA-VQA-NoPR (introduced in Sec. 4.1), which uses only model predictions in training to adjust document scores without pseudo relevance labels, yields a significantly lower PRRecall curve (orange curve in Fig. 5(b)) than RA-VQA-FrDPR (blue curve) across all  $K_{\text{test}}$ s, but achieves much higher VQA performance (Fig. 5(a)). This further confirms that Pseudo Relevance Labels are a weak signal and a high PRRecall does not necessarily guarantee to gather truly relevant knowledge in retrieval.

## C Evaluation of Computational Cost

Model	GPU memory	Wall time	Batch Size
RAVQA	70G	180 min	2 × 16
RAVQA-FrDPR TRiG*	67G 70G	160 min 220 min	$4 \times 8$ $16 \times 2$

Table 4: Comparing computational cost of models. GPU memory: total GPU memory occupied in runtime (approximately); Wall time: the required training time for 8 epochs; Batch Size: per\_device\_batch\_size× gradient\_accumulation\_steps. The statistics was collected using one Nvidia A100 GPU (80G).

As shown in Fig. 4, comparing with the DPR baseline, the wall time is increased by only 20 mins for RAVQA joint training. Our method does not significantly increase the computation needs. In comparing with TRiG, the total training time is also reduced; this is because TRiG concatenates all *K* hidden states for decoding (Gao et al., 2022), which is computationally expensive.

Model	Wall Time (8 epochs)
RAVQA $K_{\text{train}} = 5$	3 hrs
RAVQA $K_{\text{train}} = 10$	7 hrs
RAVQA $K_{\text{train}} = 15$	12 hrs

Table 5: Wall time of the RAVQA model with an increasing  $K_{\text{train}}$ . More time is required for training 8 epochs as  $K_{\text{train}}$  increases.

As K increases, the DPR-based system takes significantly more time to train. This is a real issue for other DPR-based systems that use very high  $K_{\text{train}}$  (e.g. TRiG, KAT), but not the case for our framework: we verified in Sec.4.4.4 that K=5 achieves almost the same results as  $K\geq 10$ , which is a desirable feature that can significantly reduce the required computation cost for achieving the best performance.

## D Random Guess with OK-VQA Questions

Question	Prediction
What type of bird is this?	hawk
What time of day is it?	afternoon
What kind of dog is this?	chihuaha
What kind of bird is this?	hawk
What sport is this?	horse race
What breed of horse is that?	clydesdale
What century is this?	19th
What kind of birds are these?	pigeon
What city is this?	new york
What is the weather like?	rainy
What do these animals eat?	grass
What activity is this?	skateboard
How long do these animals live?	20 years
What type of train is this?	passenger
What is this used for?	travel
What is this room used for?	sleep
What kind of bird is that?	hawk
What food does the animal eat?	cat food
What type of dog is this?	chihuaha
What food do these animals eat?	cat food
What place is this?	switzerland
What breed of cat is this?	calico
What season is this?	winter

Table 6: Example of random guesses with only question input. Random guess achieved a good VQA Score by matching to the answers by chance. But the OK-VQA questions are still not directly answerable without access to the associated images.

From the feature ablation study we found that our RA-VQA-NoDPR achieved  $\sim\!28$  VQA Score relying on only questions. This is due to the fact that  $\sim\!75\%$  of answers to training questions appear in the answers to test questions. As shown in Table 6, for each distinct question, the model learned to generate the same answer without access to the associated images. These random guesses can match to the answers of some test questions by chance, leading to a good VQA Score. By inspection we report that most of the successful cases are random guesses, and these questions are still not directly answerable without reading the associated images.

#### E Case Study

We present a case study in Fig. 6 to compare RA-VQA-FrDPR and our proposed RA-VQA framework. Conclusions are provided to each case in the figure.

# F Generalizing RAVQA to Other Datasets

We are also interested in whether this approach is also generalisable to other similar VQA tasks that may benefit from improved passage retrieval.

We implement our framework on another knowledge-based VQA task, Fact-based VQA (FVQA) (Wang et al., 2017). This dataset contains commonsense factoid VQA questions, such as "Question: which object in the image can cut you? Answer: the knife". In contrast to OKVQA where no knowledge base is provided, FVQA grounds each question-answer pair with a fact (a triplet from several 'common sense' knowledge bases, including ConceptNet (Speer et al., 2017), Webchild (Tandon et al., 2017), and DBpedia (Auer et al., 2007)). A triplet contains a head node, a relation, and a tail node (e.g. [Car] /r/HasA [4 wheels]). To cope with passage retrieval, these knowledge triplets are flattened into surface texts (e.g. "[car] has [4 wheels]") such that DPR can be directly applied to retrieve them. We replace pseudo relevance with groundtruth relevance since relevant triplets for answering questions are given.

The metrics used for assessing performance are Accuracy and Recall, with their standard deviations of 5 splits. Accuracy counts the portion of questions that are successfully answered, while Recall@K measures how likely the retrieved Kknowledge triplets contain the answer node. Since FVQA was designed for answer selection instead of open-ended answer generation, prior works used accuracy as "whether the answer node is successfully selected from all KG nodes". To enable fair evaluation with our open-ended framework, in calculating accuracy, a question is considered successfully answered if the answer node is the closest node to the generated answer string (shortest in Levenshtein distance). For example, the generated answer 'knives' is still a valid answer since the answer node '[knife]' can be matched with a shortest Levenshtein distance.

The significance of performance is guaranteed by reporting the average of 5 splits (as in the official FVQA evaluation). In total we trained 5 DPR models and  $5 \times 3$  models (RAVQA, RAVQA-FrDPR, and RAVQA-NoDPR) with the same hyperparameters. Each split has approximately half questions for training and the remaining for testing.

We compare with three systems in prior work: (1) FVQA (Wang et al., 2017): the baseline sys-

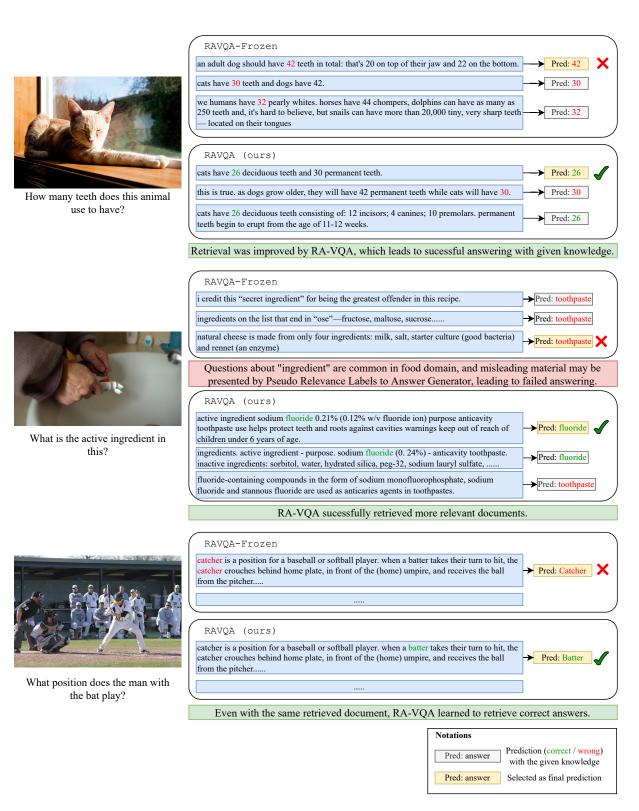


Figure 6: A case study comparing RA-VQA-FrDPR (baseline) and our RA-VQA that benefits from joint training of retrieval and answer generation.

tem provided in the official FVQA dataset paper.

- (2) GCN (Narasimhan et al., 2018b): a model that leverages graph convolutional networks (GCNs) to aggregate features from visual/language/fact modalities.
- (3) Mucko (Zhu et al., 2020): the current state-of-the-art system that uses GCNs to combine visual, fact, and semantic graphs.

Model	Accuracy (Std.)	Recall@5 (Std.)
Mucko	73.06 ( - )	- (0.77 (0.07)
<b>RAVQA</b> GCN	69.88 (0.13) 69.35 ( - )	68.77 (0.87)
RAVQA-FrDPR	68.81 (0.59)	64.54 (0.80)
RAVQA-NoDPR FVOA	67.93 (0.82) 58.76 (0.92)	-
гуул	36.70 (0.92)	-

Table 7: Model performance on the FVQA dataset (sorted by accuracy). Our proposed systems are in bold.

As shown in Table 7, RAVQA-NoDPR achieves an already strong result (67.93% accuracy) compared to early work in FVQA, showing that the extracted vision-to-language features are useful and text-based Transformers are able to learn to answer commonsense VQA questions well without accessing the provided knowledge graph (ConceptNet). The incorporation of DPR boosts the performance to 68.81% with 64.54% Recall@5, showing that retrieval works as expected and the retrieved knowledge triplets are exploited in answer generation. The joint training scheme improves the retrieval (64.54% to 68.77% Recall@5) as well as the overall performance (68.81% to 69.88% Accuracy). This demonstrates that our proposed joint training framework is generalizable to other KB-VQA tasks, though the passages used in retrieval are simply flattened surface texts of KG triplets.

In comparing with other systems in the FVQA benchmark, our best system ranks second without an explicit design for leveraging KG structures. This shows the power of open-ended answer generation with text-based Transformers. But we emphasise that better performance could be achieved through designing a more specialised retrieval component for the structured knowledge base used in this task.

To summarise, our system shows great generalizability in an external KB-VQA task that was constructed very differently. Therefore, the proposed framework can serve as a strong basis for future improvements.

Model	T5	GPT-3	$K_{ m train}$	$K_{ m test}$	Knowl. Src.	PRRecall	EM	VQA Score
ConceptBERT	×	×	-	-	С			33.66
KRISP	X	×	-	-	C + W			38.35
VRR	×	×	100	100	GS (trn/full)			39.22 / 45.08
MAVEx	×	×	-	-	W + C + GI			39.40
KAT-T5	$\checkmark$	×	40	40	W			44.25
TRiG	$\checkmark$	×	5	5	W		49.21	45.51
TRiG	$\checkmark$	×	100	100	W		53.59	49.35
TRiG-Ensemble	$\checkmark$	×	100	100	W		54.73	50.50
TRiG*	$\checkmark$	×	5	5	GS (trn/full)		52.79	45.11 / 48.32
RAG*	$\checkmark$	×	5	5	GS (trn/full)	82.34	52.52	44.90 / 48.22
RA-VQA (Ours)	$\checkmark$	×	5	5	GS (trn/full)	82.84	58.72	48.77 / 53.81
RA-VQA (Ours)	$\checkmark$	×	5	50	GS (trn/full)	96.55	59.41	49.24 / <b>54.48</b>
Ablation Study								
RA-VQA-FrDPR	<b>√</b>	×	5	5	GS (trn/full)	81.25	55.77	47.05 / 51.22
RA-VQA-NoPR	$\checkmark$	×	5	5	GS (full)	77.67	57.80	52.98
RA-VQA-FrDPR-NoCT	$\checkmark$	×	5	5	GS (full)	81.25	54.99	50.66
RA-VQA-NoCT	$\checkmark$	×	5	5	GS (full)	83.77	57.51	52.67
GPT-3-based Systems (>	GPT-3-based Systems (>175 Billion Parameters)							
PICa	×	<b>√</b>	-	-	GPT-3			48.00
KAT-Knowledge-T5	$\checkmark$	$\checkmark$	40	40	W + GPT-3			51.97
KAT-Ensemble	$\checkmark$	$\checkmark$	40	40	W + GPT-3			54.41

Table 8: Full table of RA-VQA vs Baseline Systems. Knowledge Sources:  $\underline{\mathbf{C}}$  onceptNet;  $\underline{\mathbf{W}}$  ikipedia;  $\underline{\mathbf{G}}$  oogle  $\underline{\mathbf{S}}$  earch;  $\underline{\mathbf{G}}$  oogle  $\underline{\mathbf{I}}$  mages;  $\underline{\mathbf{GPT-3}}$  closed book knowledge. 'GS (trn/full)' indicates if the Google Search data is constrained to contain only documents relevant to training questions. PRRecall, HSR, FSR, and EM are reported on 'GS (full)' systems as percentage (%). PRRecall is reported at the corresponding  $K_{\text{test}}$ .