

# Breaking Down Questions for Outside-Knowledge Visual Question Answering

Anonymous ACL submission

## Abstract

There is a recent trend towards Knowledge-Based VQA (KB-VQA) where different aspects of the question requires different source of knowledge including the image’s visual content and external knowledge such as common-sense concepts and factual information. To address this issue, we propose a novel approach that passes knowledge from various sources between different pieces of semantic content in the question. Questions are first segmented into several chunks, and each segment is used to generate queries to retrieve knowledge from ConceptNet and Wikipedia. Then, a graph neural network, taking advantage of the question’s syntactic structure, integrates the knowledge for different segments to jointly predict the answer. Our experiments on the OK-VQA dataset show that our approach achieves new state-of-the-art results.

## 1 Introduction

Over the past few years, Visual Question Answering (VQA) has emerged as a challenging task where a machine learning system needs to recognize and analyze key visual content within the image and predict an answer to a natural language question. Most recent systems (Yu et al., 2019; Lu et al., 2019; Tan and Bansal, 2019; Li et al., 2019; Zhou et al., 2020; Chen et al., 2020; Lu et al., 2020) utilize multi-modal transformers to jointly encode the entire question and the visual content, achieving a strong performance on various VQA benchmarks (Antol et al., 2015; Hudson and Manning, 2019; Singh et al., 2019).

There is a recent trend towards knowledge-based VQA (KB-VQA) (Wang et al.; Marino et al., 2019) where the information in the image alone is not sufficient for answering the visual questions. These questions cover a wide range of real-world topics, and therefore, require VQA systems to incorporate various types of external knowledge beyond the image content. For example, encyclopedia

articles provide factual statements, and common-sense knowledge bases offer everyday concepts and their relations. Both knowledge sources have been proven effective and are widely used in previous work (Wang et al.; Marino et al., 2019; Zhu et al., 2020; Li et al., 2020b; Marino et al., 2021; Wu et al., 2021).

While general VQA systems consider two modalities (*i.e.* question and image), the information across more modalities has to be properly utilized by KB-VQA systems to accommodate different types of knowledge input. This key difference introduces significant challenges to achieving good KB-VQA performance. First, knowledge representations can vary significantly across different knowledge sources, including factual sentences (Wu et al., 2021; Marino et al., 2019), knowledge triples (Wang et al.), concepts (Gardères et al., 2020) and images (Wu et al., 2021). More importantly, a system needs to understand which knowledge should be used for different semantic segments of the question. As shown in Fig. 1, KB-VQA systems need to link the segment “the vegetable that garnishes this dish” to the carrot on the plate and then query knowledge bases to find out which “human body part” particularly benefits from the nutrients in carrots.

Simply encoding the entire question for either retrieving or filtering the knowledge, as most KB-VQA systems (Wang et al.; Marino et al., 2019; Zhu et al., 2020; Li et al., 2020b; Marino et al., 2021; Wu et al., 2021) do, can cause confusion since different parts of the question focus on different aspects that can be either outside or inside the image. As depicted in Fig. 1, searching for “human body part” and “other surfaces” within the image may cause VQA systems to focus on irrelevant aspects of the image. To address this issue, we introduce a break-down VQA approach that segments visual questions into several semantic chunks, assuming that each chunk focuses on a sin-

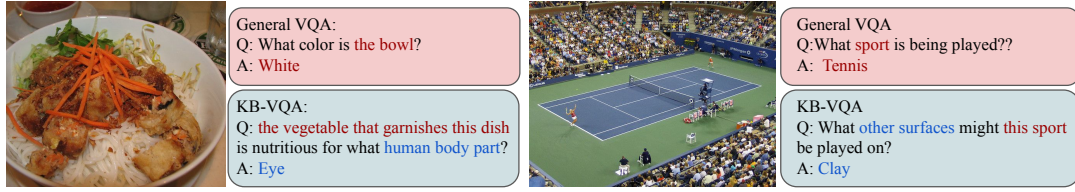


Figure 1: Examples of general and knowledge-based (KB) visual questions. The question and answer segments that focus on visual content within the image are highlighted in red, and the segments that requires external knowledge are highlighted in blue.

083 gle aspect. Those segments serve as semantic units  
 084 and are used to gather knowledge from various  
 085 sources. Finally, using a dependency parser (Hon-  
 086 nibal and Montani, 2017), a Graph Convolutional  
 087 Network (GCN) (Veličković et al., 2018) is con-  
 088 structed which assembles the retrieved knowledge  
 089 to predict the answer.

090 We evaluate our framework, break-down VQA,  
 091 on the OK-VQA dataset (Marino et al., 2019), the  
 092 largest KB-VQA dataset to date and our approach  
 093 achieves state-of-the-art results. This demonstrates  
 094 that breaking down questions and understanding  
 095 the role of each segment is especially important in  
 096 answering knowledge-based visual questions. In  
 097 summary, our main contributions are: (1) a novel  
 098 graph-based system that allows different seman-  
 099 tic segments in the questions to access different  
 100 sources of knowledge; and (2) we show that our  
 101 system, by integrating multiple sources of knowl-  
 102 edge, achieves state-of-the-art performance on OK-  
 103 VQA.

## 104 2 Related Work

### 105 2.1 Visual Question Answering

106 Visual Question Answering (VQA) has witnessed  
 107 significant progress with the introduction of multi-  
 108 modal transformers (Yu et al., 2019; Zhou et al.,  
 109 2020; Lu et al., 2020, 2019; Tan and Bansal, 2019;  
 110 Liu et al., 2019; Li et al., 2019, 2020a; Chen et al.,  
 111 2020). These transformers are pretrained on aux-  
 112 iliary tasks, including VQA, referring-expression  
 113 interpretation, image captioning, *etc.*, using vari-  
 114 ous multi-modal datasets (Sharma et al., 2018; An-  
 115 tol et al., 2015; Hudson and Manning, 2019; Suhr  
 116 et al., 2017; Yu et al., 2016; Young et al., 2014).  
 117 Cross attention modules are built over the textual  
 118 and visual modalities to learn a joint representation  
 119 for the entire question and the detected objects.

### 120 2.2 Knowledge-Based Visual Question 121 Answering

122 While VQA involves visual questions whose an-  
 123 swers can be directly found within the image, there  
 124 is a recent trend toward Knowledge-Based Visual  
 125 Question Answering (KB-VQA) that requires VQA  
 126 systems to incorporate knowledge from various ex-  
 127 ternal sources.

128 Recent high-performing KB-VQA systems are  
 129 mainly learning-based following general VQA sys-  
 130 tems, and incorporate additional modules to re-  
 131 trieve external knowledge. One (Narasimhan and  
 132 Schwing, 2018) learns to retrieve facts from a  
 133 knowledge base. Another (Narasimhan et al., 2018)  
 134 utilizes a GCN (Tompson et al., 2014) over the fact  
 135 graph where each node is a representation of an  
 136 image-question-entity triplet. A third (Li et al.,  
 137 2020b) introduces a knowledge-graph augmenta-  
 138 tion model to retrieve context-aware knowledge  
 139 sub-graphs, and then learns to aggregate the useful  
 140 visual and question relevant knowledge.

141 Although the knowledge is obtained from a wide  
 142 range of sources and encoded in different formats,  
 143 these previous systems simply learn to mine rel-  
 144 evant facts based on the entire question, which,  
 145 as mentioned above, could cause confusion. In  
 146 contrast, we present an approach that breaks the  
 147 question down into several segments and then uses  
 148 each of these segments to gather the appropriate  
 149 knowledge, which is then integrated to answer the  
 150 overall question.

151 The most similar work to ours is KRISP (Marino  
 152 et al., 2021), which combines knowledge from  
 153 both implicit question-image embeddings and ex-  
 154 plicit symbolic information from knowledge bases.  
 155 While KRISP aims to build a shared knowledge  
 156 graph for all KB questions, we build a question  
 157 graph for each specific visual question to address  
 158 two main issues with a shared knowledge graph: (1)  
 159 The knowledge is not question-specific and may

mislead the answer predictor since some knowledge may not be appropriate for a specific question; (2) The size of the knowledge graph is not scalable as it is hard to cover and process all the required knowledge for all questions.

### 2.3 Breaking Down Visual Questions.

Previous work has explored both rule-based (Andreas et al., 2016; Wolfson et al., 2020) and learning-based (Hu et al., 2017, 2018; Mao et al., 2019; Wolfson et al., 2020) approaches to break down visual questions. Rule-based approaches typically define a set of decomposition rules and a full decomposition is obtained by recursively applying these rules until no rule is matched. In particular, one method (Andreas et al., 2016) parses the questions and breaks it into a sequence of programs to execute. Another (Wolfson et al., 2020) breaks the question into several steps, each of which is encoded as a natural language expression. Learning-based approaches either learn to recursively rank some predefined modules to synthesize the entire network layout for solving a visual question (Hu et al., 2017, 2018) or directly learn to generate the steps using a seq2seq method (Mao et al., 2019; Wolfson et al., 2020). These approaches work especially well for datasets that represent queries as programs, including CLVER (Johnson et al., 2017) and GQA (Hudson and Manning, 2019).

### 2.4 Graph Convolutional Networks

Graph Convolutional Networks (GCNs) (Kipf and Welling, 2017) generalize Convolutional Networks (CNN) to accommodate graph-structured input. Various types of graph input for VQA have been explored including scene graphs generated by an object and relation detector (Ren et al., 2015; Yang et al., 2018), and knowledge graphs retrieved from a wide range of sources, such as DB-Pedia (Auer et al., 2007), ConceptNet (Liu and Singh, 2004), VisualGenome (Krishna et al., 2017) and hasPart KB (Bhakthavatsalam et al., 2020). Most KB-VQA systems (Ramnath and Hasegawa-Johnson, 2021; Narasimhan et al., 2018; Li et al., 2020b; Marino et al., 2021) build their GCNs on top of these knowledge graphs and extract relevant evidence using the entire question representation. Here, we explore an approach that constructs a reasoning graph from the question, where each node is a semantic segment of the question. Our graph utilizes the syntactic structure of the questions to better integrate the question segments that utilize both the visual content in the

image and relevant external knowledge.

## 3 Approach

We present the break-down VQA approach, a three-step framework. First, it segments visual questions into semantic chunks. Next, each segment, serving as a semantic unit, is used to retrieve knowledge from different external sources. Finally, a Graph Neural Network (GCN) integrates this retrieved knowledge to predict an answer. Fig. 2 illustrates the approach.

We instantiate our approach on top of the high-performing ViLBERT-multi-task as a base system (Lu et al., 2020) that provides a set of answer candidates  $A = \{a_1, \dots, a_n\}$  for each question-image pair. We also extract the product of its pooled features for the textual and visual BERT output,  $\mathbf{z}$ , as a joint representation of the question and the image.

### 3.1 Breaking Down Visual Questions

Given a visual question  $q$  that consists of  $l$  tokens  $(q_1, \dots, q_l)$  where a token is either a word or a WordPiece produced by a tokenizer (Vaswani et al., 2017), and its question segmentation is a set of token chunks  $X = (x^1, \dots, x^m)$  where each  $x^i$  consists of a sub-sequence of  $q$ , *i.e.*  $x^i = (q_1^i, \dots, q_{l_i}^i)$ .

Since different parts of a knowledge-based visual question need to focus on different sources of knowledge, encoding the entire question as a whole leads to inefficiency in both knowledge retrieval and answer prediction. To address this issue, our approach breaks down the question into segments where each segment contains only one semantic unit that can either be grounded in the image, or linked to external knowledge bases. Note that, we do not restrict each segment to only retrieve from a single knowledge source but let the VQA model choose the right source.

To this end, we first extract nouns, noun chunks, and verbs in the question as knowledge segments. For example, ‘other surfaces’, ‘this sport’, and ‘play’ are extracted for the second example in Fig. 1. Specifically, we utilize the ‘en\_core\_web\_sm’ spaCy parser (Honnibal and Montani, 2017) to dependency parse the question and POS-tag each word. Then, we extract the noun chunks (*i.e.* flattened phrases with a noun head in the parse tree) and lemmatized verbs. We also group any tokens between those extracted knowledge segments as additional segments to ensure completeness.

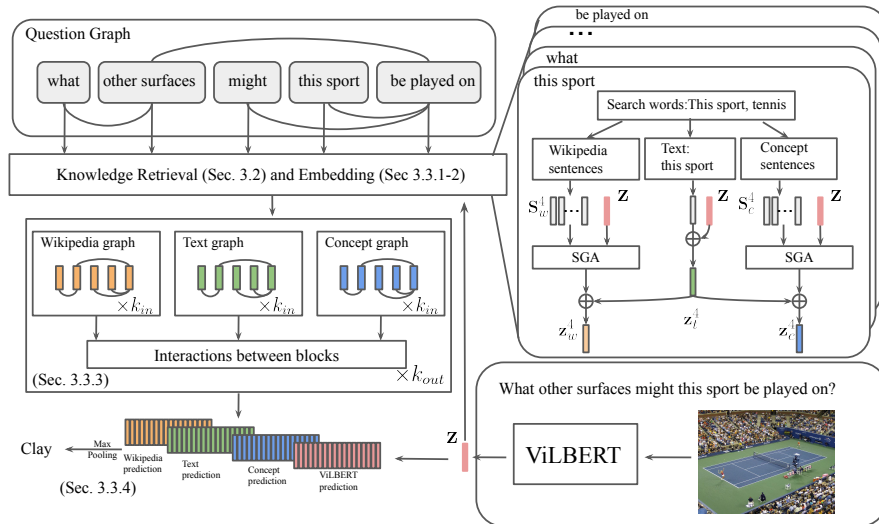


Figure 2: Model overview of the BreakDown VQA approach. The question is segmented into semantic chunks (left top). These chunks are used to retrieve external information from Wikipedia and ConceptNet. Each retrieved piece of knowledge is then encoded as a vector (right top), and fed to a graph neural net (left middle) to predict an answer for each knowledge source. The individual results are then max-pooled to get the final prediction (left bottom).

### 3.2 Knowledge Retrieval

We retrieve knowledge from Wikipedia and ConceptNet for each extracted segment inspired by the answer-guided knowledge retrieval (Wu et al., 2021) that ensure the relevancy of the external knowledge. Note that although our knowledge retrieval process is similar to (Wu et al., 2021), our main contribution lies in allowing different aspects of a question to access knowledge from different sources while (Wu et al., 2021) mine important knowledge using the entire question.

**Search Word Extraction.** We first remove the stop words in the segment and regard the remaining tokens as the search words. Then, we enrich these search words with object annotations, including linking the segment to objects in the image, text extracted using OCR, and brand detection following (Wu et al., 2021).<sup>1</sup> In particular, a pretrained ViLBERT-multi-task model (Lu et al., 2020) is used as the object linker. This system can generate linking scores indicating the confidence of linking phrases to detected objects. The linking is approved when its score is over 0.5. With the linked objects, a Google API is used to recognize words in text regions using OCR and company brands. In addition, we also detect common attributes of these objects using a Faster-RCNN (Ren et al., 2015) on a Detectron platform pretrained on Visual-Genome

<sup>1</sup>See section “S1: Answer-Agnostic Search Word Extraction”

data.<sup>2</sup> This process results in a set of search words for each segment. For example, the search-word set for the segment ‘the vegetable’ for the first example in Fig. 1 is {‘vegetable’, ‘carrot’, ‘red vegetable’}.

**Knowledge Retrieval.** We use two knowledge sources to extract information about the question segments in  $X$ , *i.e.* relevant textual facts and commonsense concepts as in (Wu et al., 2021). In contrast to (Wu et al., 2021), we retrieve knowledge independently for each segment instead of for the entire question. This ensures that the retrieved knowledge provides information about the given segment, and allows the VQA system to determine whether a particular piece of external knowledge about this segment is important.

*Retrieving from Wikipedia.* For each segment  $x_i$ , we query its search words and collect all sentences from the retrieved Wikipedia articles. We use answer-guided knowledge retrieval (Wu et al., 2021) to filter out irrelevant sentences. Specifically, we first keep the Wikipedia sentences that contain both at least one of the search words and one of the top 5 answer candidates predicted by ViLBERT-multi-task. Then, the remaining sentences are ranked according to the highest precision BERT-scores (Zhang et al., 2020) between the sentence and the question converted to a declarative statement (Demszky et al., 2018) and the top-5 an-

<sup>2</sup>We were careful to remove the OK-VQA test images from the training data for the Faster-RCNN system.

314 swer candidates. We keep the top-80 sentences in  
 315 total for each visual question and regard the other  
 316 sentences as irrelevant.

317 *Retrieving from ConceptNet.* Commonsense con-  
 318 cepts provide structured knowledge that is usually  
 319 not covered in factual Wikipedia sentences. Similar  
 320 to Wikipedia-article retrieval, we query the search  
 321 words for each segment and collect the retrieved  
 322 concepts. First, we keep all the concept triples  
 323 whose subjects and objects contain the search word  
 324 and one of the answer candidates from  $A$ . Then,  
 325 we convert other concept triples to sentences and  
 326 rank them according to the highest precision BERT-  
 327 scores between the sentence and the statements  
 328 from the question and answers. We also keep the  
 329 top-80 sentences in total for each visual question  
 330 and regard the other concepts as irrelevant.

331 **Matching Textual Knowledge.** For each query,  
 332 the sentences from Wikipedia and the concepts  
 333 from ConceptNet with a mean recall greater than  
 334 0.6 are matched to the search words. Mean recall  
 335 is defined as the average cosine similarity between  
 336 the GloVe embedding of the words in the search  
 337 word and their most similar word in the sentence  
 338 or in the concept. To ensure knowledge relevance,  
 339 we remove sentences that are matched to only a  
 340 single search word. We keep the top  $k_w$  sentences  
 341  $S_w^i = \{s_{w,1}^i, \dots, s_{w,k_w}^i\}$  according to the mean recall  
 342 as the textual facts for segment  $x_i$  from its search  
 343 word set, where  $s_{w,j}^i$  denotes the  $j$ -th Wikipedia  
 344 sentence for the  $i$ -th segment. Similarly, for con-  
 345 cepts, we keep the top  $k_c$  concept sentences  $S_c^i =$   
 346  $\{s_{c,1}^i, \dots, s_{c,k_c}^i\}$  for segment  $x_i$ , where  $s_{c,j}^i$  denotes  
 347 the  $j$ -th concept sentence for the  $i$ -th segment.

### 348 3.3 VQA model

349 This section describes the final VQA system that  
 350 incorporates the retrieved knowledge for each of  
 351 the semantic segments. We first generate features  
 352 for each knowledge sentence from Wikipedia and  
 353 ConceptNet. Then a representation of each source  
 354 for each segment is computed using these sentence  
 355 features. Finally, a GCN is employed that utilizes  
 356 the syntactic structure of the visual question and  
 357 produces joint features for predicting the answer.

#### 358 3.3.1 Knowledge Sentence Embedding

359 We use a word embedding matrix initialized by  
 360 GloVe vectors (Pennington et al., 2014) to com-  
 361 pute a word vector for each token in the knowl-  
 362 edge sentence. Then, a single layer LSTM with  
 363 a hidden states of 768 is built on top of the word

364 embeddings and the features for the last token are  
 365 extracted. This process produces a 768-d feature  
 366 vector for each sentence from both Wikipedia  $S_w^i$   
 367 and ConceptNet  $S_c^i$ , resulting in knowledge fea-  
 368 ture matrices  $\mathbf{S}_w^i \in R^{k_w \times 768}$  and  $\mathbf{S}_c^i \in R^{k_c \times 768}$  for  
 369 segment  $i$ , respectively.

#### 370 3.3.2 Segment Embedding

371 We produce an embedding for each segment by  
 372 integrating three representations, a content repre-  
 373 sentation of the text of the segment in the question  
 374 and two representations of relevant external knowl-  
 375 edge (Wikipedia + ConceptNet).

376 **Content Embedding.** To preserve all of the  
 377 information in the question, we employ the text  
 378 of each segment as input to the VQA model. We  
 379 use the GloVe embedding approach to encode seg-  
 380 ments. Similar to the knowledge sentence embed-  
 381 ding, an LSTM is used to sequentially encode the  
 382 GloVe vectors and the hidden state of the last token  
 383 is extracted as the content representation  $\mathbf{s}_t^i$ . The  
 384 final content embedding of segment  $i$  is computed  
 385 as the element-wise summation of  $\mathbf{s}_t^i$  and the pro-  
 386 jection of  $\mathbf{z}$ , *i.e.*  $\mathbf{z}_t^i = \mathbf{s}_t^i + \text{fc}(\mathbf{z})$ , where  $\text{fc}$  denotes  
 387 a fully connected layer.

388 **Knowledge Embedding.** As shown in Eqs. 1  
 389 and 2, we embed the knowledge matrices  $\mathbf{S}_w^i$  and  $\mathbf{S}_c^i$   
 390 for segment  $x_i$  into vector representations  $\mathbf{z}_w^i$  and  
 391  $\mathbf{z}_c^i$  that contain the question-relevant information  
 392 from the external knowledge source sentences  $S_w^i$   
 393 and  $S_c^i$ . In particular, we utilize a Self- and Guided-  
 394 Attention (SGA) module (Yu et al., 2019) where  
 395 the question and image representation  $\mathbf{z}$  from ViL-  
 396 BERT is used as a query, and the knowledge ma-  
 397 trices serve as keys and values. The SGA modules  
 398 provide a trainable method for mining question-  
 399 relevant knowledge from the retrieved materials  
 400 in contrast to the rule-based method used in the  
 401 knowledge retrieval process. In order to prevent  
 402 the case where the retrieved knowledge is empty,  
 403 we add the content embedding to the knowledge  
 404 embedding for each source.

$$345 \mathbf{z}_w^i = \text{SGA}(\mathbf{z}, \mathbf{S}_w^i) + \mathbf{z}_t^i \quad (1) \quad 405$$

$$346 \mathbf{z}_c^i = \text{SGA}(\mathbf{z}, \mathbf{S}_c^i) + \mathbf{z}_t^i \quad (2) \quad 406$$

#### 407 3.3.3 Graph Neural Networks

408 **Building the Graph Structure.** We treat the seg-  
 409 ments' embeddings  $\{\mathbf{z}_k^i\}$  as nodes, where  $i$  denotes  
 410 the segment's index and  $k$  indexes the knowledge  
 411 source, and establish an edge between each pair if  
 412 there is a direct connection between tokens from

the two segments in the dependency parse tree. Given the parse tree  $\mathcal{E}_q$  of question  $q$ , which establishes edges between tokens in  $q$ , the edges of the segments  $\mathcal{E}$  are defined in Eq. 3:

$$\mathcal{E} = \{(\mathbf{z}_k^i, \mathbf{z}_k^j) \mid \exists(q^m \in x^i, q^n \in x^j)(q^m, q^n) \in \mathcal{E}_q\} \quad (3)$$

This produce a graph structure  $\mathcal{G}^k = (\{\mathbf{z}_k^i\}, \mathcal{E})$  for each modality  $k$ .

**Graph Neural Networks Architectures.** The networks consists of  $k_{out}$  blocks, where each block contains  $k_{in}$  graph layers. The node features within each block interact with other nodes’ features from the same modality, determining its importance to solve the visual question. The knowledge from different external modalities is fused outside the blocks to build connections to other types of knowledge.

Graph Neural Networks within Blocks. We formalize the input to the graph neural networks as  $\mathbf{H}_{i,0}^k$  where  $k$  denotes the source of the question segments’ features, and  $i$  is the index of the block. For layer  $l$  within block  $i$ , we use a graph layer that operates a non-linear function  $F(\mathbf{H}_{i,l}^k, \mathcal{G}^k)$ , producing the input to the next graph layer  $\mathbf{H}_{i,l+1}^k$ , *i.e.*  $\mathbf{H}_{i,l+1}^k = F(\mathbf{H}_{i,l}^k, \mathcal{G}^k)$ . The input  $\mathbf{H}_{i,0}^k$  for block  $i$  is the output of the previous block  $\mathbf{H}_{i-1,l}^k$  after interactions between modalities described below except for the first block that receives the segments’ features as inputs, *i.e.*  $\{\mathbf{z}_k\}$ .

Interactions between Modalities outside the Blocks. To give the graph neural networks access to the all types of external knowledge, we fused features from different modalities outside the blocks. The fused features serve as the inputs to the next blocks of graph neural nets. In particular, the input  $\mathbf{H}_{i+1,0}^k$  to the  $i + 1$  block is the concatenation of the segments representation  $\{\mathbf{z}_k\}$  and the summation of the output of the previous block from all modalities.

### 3.3.4 Answer Prediction

We build answer prediction heads for each knowledge source that compute a probability distribution over all answer candidates. The knowledge features from the last block, *i.e.*  $\mathbf{H}_{k_{out},k_{in}}^k$  are averaged and fed to the answer prediction head that consists of two consecutive fully-connected layers with ReLU activation. Then, we take the maximum value of these predictions for each answer candidate as the final answer predictions.

## 4 Implementation and Training Details

**Implementation.** Our break-down VQA approach is implemented on top of ViLBERT-multi-task (Lu et al., 2019), which utilizes a Mask-RCNN head (He et al., 2017) in conjunction with a ResNet-152 base network (He et al., 2016) as the object detection module. Convolutional features for at most 100 objects are extracted for each image as the visual features, *i.e.* a 2,048 dimensional vector for each object.

Since the OK-VQA test dataset contains COCO images from the validation set that are used to train the officially released ViLBERT model, we retrain the system from scratch using clean datasets where we remove all of the OK-VQA test images from the Visual Genome, MSCOCO, and GQA datasets. We used the default configuration when training the object detection module, pretraining on Conceptual Captions, and finally finetuning on the 12 visual-and-language tasks used in (Lu et al., 2020). We utilize a BERT tokenizer (Devlin et al., 2019) to tokenize the question and use the first 23 tokens of the question. We encode the top 5 Wikipedia sentences and top 10 ConceptNet concepts for each knowledge segment, *i.e.*  $k_w = 5$  and  $k_c = 10$ . The number of hidden units in the SGA modules in the knowledge embedding modules is set to 768. We use 4 attention heads in the SGA modules. The graph neural networks contain 2 blocks and 4 layers within each block. A SAGE (Hamilton et al., 2017) layer with transformed root node features is used as the graph layer. The Pytorch Geometric toolbox (Fey and Lenssen, 2019) is used for the GCN implementation.

**Training.** For training, we optimize the answer predictions for each knowledge source using the standard VQA loss, together with the VQA loss on the final predictions. We train the system for 75 epochs using a learning rate of  $2e-5$  for the ViLBERT parameters and  $5e-5$  for the additional parameters introduced in the BreakDown VQA system. We freeze the first 10 layers of the ViLBERT base network.

## 5 Experiments

This section evaluates our BreakDown VQA approach on the OK-VQA dataset (Marino et al., 2019). We first briefly describe the dataset, and then present results comparing to current state-of-the-art systems.

**OK-VQA dataset.** This is currently the largest

Method	Knowledge Resources	Performance
ViLBERT (Lu et al., 2019)	—	36.1
MMBERT (Marino et al., 2019)	—	37.1
ConceptBERT (Gardères et al., 2020)	—	33.7
KRISP (Marino et al., 2021)	Wikipedia + ConceptNet	37.8
MAVEx (Wu et al., 2021)	Wikipedia + ConceptNet + Google Images	38.7
Ours	Wikipedia + ConceptNet	39.1
Ours + MAVEx	Wikipedia + ConceptNet + Google Images	40.8
Ours + MAVEx (oracle)	Wikipedia + ConceptNet + Google Images	<b>42.5</b>

Table 1: Our approach outperforms current state-of-the-art approaches on the OK-VQA dataset. The middle column lists the external knowledge sources, if any, used by each VQA system.

511 knowledge-based VQA dataset. The questions are  
512 crowdsourced from human workers on Amazon  
513 Mechanical Turk instead of artificially synthesized  
514 from knowledge bases. Human judges are asked to  
515 ensure that outside knowledge beyond the image  
516 is required. Also, since it is not synthesized, there  
517 are no ground truth knowledge bases that can pro-  
518 vide a VQA system all of the necessary external  
519 knowledge. Therefore, systems have to retrieve  
520 knowledge from a variety of knowledge sources.  
521 The dataset contains 14,031 images and 14,055  
522 questions covering a variety of topics, including  
523 transportation, brands, material, sports, cooking,  
524 geography, plants, animals, science, weather, *etc.*

## 525 5.1 Main Results

526 We report results on version 1.1 of the OK-VQA  
527 dataset in Table 1, unlike the original version (*i.e.*  
528 version 1.0), answers are lemmatized to improve  
529 scoring. Our BreakDown VQA approach outper-  
530 forms all previous systems, achieving a new state-  
531 of-the-art accuracy score of 39.1%.

## 532 5.2 Ablation Study on Source Knowledge

533 This sections gives results when we ablate the ex-  
534 ternal knowledge sources. In particular, we manu-  
535 ally zero out the knowledge features  $\mathbf{z}_k$  to exclude  
536 the external information obtained from knowledge  
537 source  $k$  during training and test. We use 2 blocks  
538 and 4 layers within each block in the graph neural  
539 networks. As shown in Table 2, each knowledge  
540 source helps improve the overall performance, indi-  
541 cating the need to access to a variety of knowledge  
542 sources for solving the KB visual questions.

## 543 5.3 Ablation Studies on the Graph Model

544 Table 3 shows results on how different values for  
545 the hyper-parameters in the GCN influence VQA

Sources	Performance
Wikipedia	38.2
ConceptNet	38.5
Wikipedia + ConceptNet	39.1

Table 2: Ablation study of knowledge sources.

546 performance. It includes an extreme case using  
547 only one graph block (*i.e.*  $k_{out} = 1$ ), where the  
548 knowledge sources do not interact and predict the  
549 answer independently. We also tested two ablated  
550 models to test the contribution of the graph struc-  
551 ture that exploits the parse tree of the question. We  
552 simply build the answer prediction heads on top  
553 of the knowledge embedding of each source,  $\mathbf{z}_k$ ,  
554 where  $k$  is the knowledge source indicator. This  
555 baseline system achieves a score of 38.5, and a  
556 fully-connected graph achieves 38.7. That indi-  
557 cates that building the segments’ graph using the  
558 question’s syntactic structure helps the VQA sys-  
559 tem improve its use of the retrieved knowledge,  
560 improving the results.

$k_{out}$	$k_{in}$	Performance
1	4	38.6
2	4	39.1
2	6	38.7
3	4	38.8

Table 3: Ablation study using different GCN hyper-  
parameter values.  $k_{out}$  and  $k_{in}$  denote the number of  
blocks and the number of layers within each block.

## 561 5.4 Using BERT for Knowledge Embedding

562 We also tested a BERT-based knowledge embed-  
563 ding for encoding the retrieved sentences from the  
564 external knowledge sources. We used a pretrained  
565 BERT-base-uncased model (Devlin et al., 2019) to

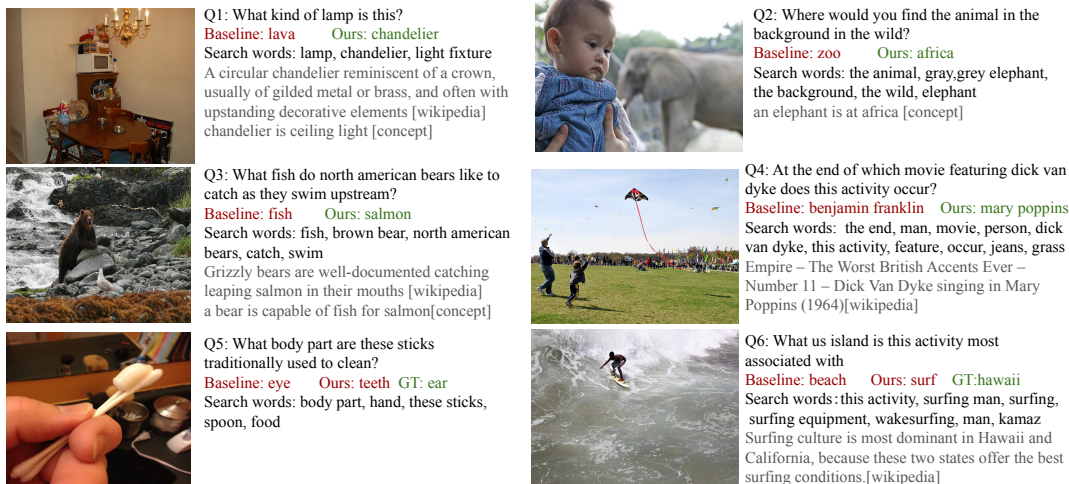


Figure 3: Qualitative results from our Break Down VQA and a ViLBERT baseline. Q1-Q4 show success cases and Q5 and Q6 illustrate a couple failure cases. Red and green denote wrong and right answers, respectively.

compute the features for each sentence. We extract the final layer representation for the “[CLS]” token as the sentence embedding to replace the GloVe embedding used in Sec. 3.3.1. Note that this BERT model is not finetuned for VQA. The BERT Embedding approach achieves a score of 38.9 compared to 39.1 using the GloVe embedding. Our hypothesis is that though BERT features may encode richer information, fine-tuning on the down-stream task is important for the final performance.

### 5.5 Combining with Answer Validation

(Wu et al., 2021) introduced an answer validation module (MAVEx) that reweights the answer confidence with a verification score obtained by examining the knowledge retrieved for each of the top answer candidates. MAVEx also uses retrieved images from Google as a third knowledge source to provide visual external knowledge. We combined our BreakDown VQA approach with a static MAVEx system that provides the weights of the top 5 answer candidates. As shown in Table 1, we achieve a score of 40.8 when combining the MAVEx weights using predicted answer candidates and 42.5 when using an oracle answer candidate set where a ground truth answer is manually inserted into the answer candidate set during validation. This shows that our approach can be effectively combined with other recent advances in KB-VQA to further improve the state-of-the-art.

### 5.6 Qualitative Results

We show some representative examples of our approach versus a ViLBERT baseline system in Fig.

3. Q1 shows an example where the answer is already in the search word list (*i.e.* chandelier), illustrating the effectiveness of enriching the segments parsed from the question with various types of annotations. Q2-Q4 show examples where our approach successfully retrieves relevant knowledge about specific segments which allows it to predict the correct answer. Q4 shows an example where Wikipedia knowledge is especially helpful and Q3 shows an example where both knowledge sources provide useful information.

We also show some common failure cases from our approach in Q5 and Q6. Q5 shows an example where object recognition fails since the cotton swabs are just annotated as sticks, making it hard to retrieve the relevant knowledge. Q6 shows a case where the retrieved knowledge seems helpful but the final prediction is wrong. It seems the VQA system failed to understand that the question is asking about a location rather than an activity.

## 6 Conclusion

We have introduced a novel approach to knowledge-based VQA that breaks down visual questions into multiple semantic segments which are used to drive the retrieval and utilization of relevant knowledge from multiple external sources. This approach achieves a new state of the art on the challenging OK-VQA benchmark. We find that segmenting questions is especially helpful for open-domain KB-VQA because different parts of the question require utilizing different types of information, such as linking to objects in the image and exploiting factual information or commonsense knowledge.



631  
632  
633  
634  
  
635  
636  
637  
638  
  
639  
640  
641  
642  
  
643  
644  
645  
646  
  
647  
648  
649  
650  
  
651  
652  
653  
654  
  
655  
656  
657  
658  
  
659  
660  
661  
662  
  
663  
664  
665  
666  
667  
  
668  
669  
670  
  
671  
672  
  
673  
674  
675  
  
676  
677  
678  
679  
  
680  
681  
682

## References

Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. 2016. Neural Module Networks. In *CVPR*.

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. VQA: Visual Question Answering. In *ICCV*.

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*. Springer.

Sumithra Bhakthavatsalam, Kyle Richardson, Niket Tandon, and Peter Clark. 2020. Do Dogs Have Whiskers? A New Knowledge Base of hasPart Relations. *arXiv preprint arXiv:2006.07510*.

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 *ECCV*. Springer.

Dorottya Demszky, Kelvin Guu, and Percy Liang. 2018. Transforming Question Answering Datasets into Natural Language Inference datasets. *arXiv preprint arXiv:1809.02922*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *ACL*.

Matthias Fey and Jan E. Lenssen. 2019. Fast Graph Representation Learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*.

François Gardères, Maryam Ziaeeafard, Baptiste Abeloos, and Freddy Lecue. 2020. ConceptBert: Concept-aware representation for visual question answering. In *Findings of the Association for Computational Linguistics: EMNLP 2020*.

William L Hamilton, Rex Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. *NeurIPS*.

Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask R-CNN. In *ICCV*.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *CVPR*.

Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural Language Understanding with Bloom Embeddings, Convolutional Neural Networks and Incremental Parsing. To appear.

Ronghang Hu, Jacob Andreas, Trevor Darrell, and Kate Saenko. 2018. Explainable Neural Computation via Stack Neural Module Networks. In *ECCV*.

Ronghang Hu, Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Kate Saenko. 2017. Learning to Reason: End-to-End Module Networks for Visual Question Answering. In *ICCV*. 683  
684  
685  
686

Drew A Hudson and Christopher D Manning. 2019. GQA: A New Dataset for Compositional Question Answering over Real-World Images. *CVPR*. 687  
688  
689

Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and Ross Girshick. 2017. CLEVR: A Diagnostic Dataset for Compositional Language and Elementary Visual Reasoning. In *CVPR*. 690  
691  
692  
693  
694

Thomas N Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. *ICLR*. 695  
696  
697

Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. 2017. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations. *IJCV*. 698  
699  
700  
701  
702  
703

Gen Li, Nan Duan, Yuejian Fang, Ming Gong, and Daxin Jiang. 2020a. Unicoder-vl: A UNiversal encoder for Vision and Language by Cross-Modal Pre-training. In *AAAI*. 704  
705  
706  
707

Guohao Li, Xin Wang, and Wenwu Zhu. 2020b. Boosting Visual Question Answering with Context-aware Knowledge Aggregation. In *ACM Conference on Multimedia*. 708  
709  
710  
711

Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. Visualbert: A Simple and Performant Baseline for Vision and Language. *arXiv preprint arXiv:1908.03557*. 712  
713  
714  
715

Bei Liu, Zhicheng Huang, Zhaoyang Zeng, Zheyu Chen, and Jianlong Fu. 2019. Learning Rich Image Region Representation for Visual Question Answering. *arXiv preprint arXiv:1910.13077*. 716  
717  
718  
719

Hugo Liu and Push Singh. 2004. ConceptNet: a Practical Commonsense Reasoning Tool-kit. *BT technology journal*. 720  
721  
722

Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks. In *NeurIPS*. 723  
724  
725  
726

Jiasen Lu, Vedanuj Goswami, Marcus Rohrbach, Devi Parikh, and Stefan Lee. 2020. 12-in-1: Multi-Task Vision and Language Representation Learning. In *CVPR*. 727  
728  
729  
730

Jiayuan Mao, Chuang Gan, Pushmeet Kohli, Joshua B Tenenbaum, and Jiajun Wu. 2019. The Neuro-Symbolic Concept Learner: Interpreting Scenes, Words, and Sentences from Natural Supervision. *ICLR*. 731  
732  
733  
734  
735

736	Kenneth Marino, Xinlei Chen, Devi Parikh, Abhinav Gupta, and Marcus Rohrbach. 2021. KRISP: Integrating Implicit and Symbolic Knowledge for Open-Domain Knowledge-Based VQA. In <i>CVPR</i> .	Tomer Wolfson, Mor Geva, Ankit Gupta, Matt Gardner, Yoav Goldberg, Daniel Deutch, and Jonathan Berant. 2020. Break It Down: A Question Understanding Benchmark. <i>TACL</i> .	789
737			790
738			791
739			792
740	Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. OK-VQA: A Visual Question Answering Benchmark Requiring External Knowledge. In <i>CVPR</i> .	Jialin Wu, Jiasen Lu, Ashish Sabharwal, and Roozbeh Mottaghi. 2021. Multi-Modal Answer Validation for Knowledge-Based VQA. <i>arXiv preprint arXiv:2103.12248</i> .	793
741			794
742			795
743			796
744	Medhini Narasimhan, Svetlana Lazebnik, and Alexander Schwing. 2018. Out-of-The-Box: Reasoning with Graph Convolution Nets for Factual Visual Question Answering. In <i>NeurIPS</i> .	Jianwei Yang, Jiasen Lu, Stefan Lee, Dhruv Batra, and Devi Parikh. 2018. Graph R-CNN for Scene Graph Generation. In <i>ECCV</i> .	797
745			798
746			799
747			
748	Medhini Narasimhan and Alexander G Schwing. 2018. Straight to the Facts: Learning Knowledge Base Retrieval for Factual Visual Question Answering. In <i>ECCV</i> .	Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. From Image Descriptions to Visual Denotations: New Similarity Metrics for Semantic Inference over Event Descriptions. <i>TACL</i> .	800
749			801
750			802
751			803
752	Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. Glove: Global Vectors for Word Representation. In <i>EMNLP</i> .	Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. 2016. Modeling Context in Referring Expressions. In <i>ECCV</i> , pages 69–85. Springer.	804
753			805
754			806
755	Kiran Ramnath and Mark Hasegawa-Johnson. 2021. Seeing is Knowing! Fact-based Visual Question Answering using Knowledge Graph Embeddings. <i>AAAI</i> .	Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. 2019. Deep Modular Co-Attention Networks for Visual Question Answering. In <i>CVPR</i> .	808
756			809
757			810
758	Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster R-CNN: Towards Real-time Object Detection with Region Proposal Networks. In <i>NIPS</i> .	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2020. BERTScore: Evaluating Text Generation with BERT. <i>ICLR</i> .	811
759			812
760			813
761	Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual Captions: A Cleaned, Hypernymed, Image Alt-text Dataset For Automatic Image Captioning. In <i>ACL</i> .	Luowei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason J. Corso, and Jianfeng Gao. 2020. Unified Vision-Language Pre-Training for Image Captioning and VQA. In <i>AAAI</i> .	814
762			815
763			816
764			817
765	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 <i>CVPR</i> .	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 <i>IJCAI</i> .	818
766			819
767			820
768			821
769	Alane Suhr, Mike Lewis, James Yeh, and Yoav Artzi. 2017. A corpus of natural language for visual reasoning. In <i>ACL</i> .		
770			
771			
772	Hao Tan and Mohit Bansal. 2019. LXMERT: Learning Cross-Modality Encoder Representations from Transformers. In <i>EMNLP</i> .		
773			
774			
775	Jonathan J Tompson, Arjun Jain, Yann LeCun, and Christoph Bregler. 2014. Joint Training of a Convolutional Network and a Graphical Model for Human Pose Estimation. In <i>NeurIPS</i> .		
776			
777			
778			
779	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is All You Need. In <i>NeurIPS</i> .		
780			
781			
782			
783	Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2018. Graph Attention Networks. <i>ICLR</i> .		
784			
785			
786	Peng Wang, Qi Wu, Chunhua Shen, Anthony Dick, and Anton Van Den Hengel. Fvqa: Fact-Based Visual Question Answering. <i>PAMI</i> .		
787			
788			