
A-OKVQA: A Benchmark for Visual Question Answering using World Knowledge

Dustin Schwenk¹, Apoorv Khandelwal¹, Christopher Clark¹
Kenneth Marino², Roozbeh Mottaghi¹

¹ PRIOR @ Allen Institute for AI

² Carnegie Mellon University

Abstract

The Visual Question Answering (VQA) task aspires to provide a meaningful testbed for the development of AI models that can jointly reason over visual and natural language inputs. Despite a proliferation of VQA datasets, this goal is hindered by a set of common limitations. These include a reliance on relatively simplistic questions that are repetitive in both concepts and linguistic structure, little world knowledge needed outside of the paired image, and limited reasoning required to arrive at the correct answer. We introduce A-OKVQA, a crowdsourced dataset composed of a diverse set of about 25K questions requiring a broad base of commonsense and world knowledge to answer. In contrast to the existing knowledge-based VQA datasets, the questions generally cannot be answered by simply querying a knowledge base, and instead require some form of commonsense reasoning about the scene depicted in the image. We demonstrate the potential of this new dataset through a detailed analysis of its contents and baseline performance measurements over a variety of state-of-the-art vision–language models.

<http://a-okvqa.allenai.org/>

1 Introduction

The original conception of the Visual Question Answering (VQA) problem was as a Visual Turing Test [15]. Can we give a computer an image and expect it to answer any question we ask to fool us into thinking it is a human? To truly solve this Turing Test, the computer would need to mimic several human capabilities including: visual recognition in the wild, language understanding, basic reasoning capabilities and a background knowledge about the world. Since the VQA problem was formulated, many of these aspects have been studied. Early datasets mostly studied the perception and language understanding problem on natural image datasets [2, 34, 16]. Other datasets studied complex chains of reasoning about procedurally generated images [25]. More recently, datasets include questions which require factual [36, 56, 57] or commonsense knowledge [66].

But, to a large extent, VQA has been a victim of its own success. With the advent of large-scale pre-training of vision and language models [67, 62, 32, 33, 12, 43, 8] and other breakthroughs in multi-modal architectures, much of the low-hanging fruit in the field has been plucked and many of the benchmark datasets have seen saturated performance. Even performance on the newer knowledge-based datasets has been improved by such models [67]. So how can we continue developing yet more challenging datasets? What human capabilities are not yet expressed by current models?

We propose the following. First, continuing the direction of past work in knowledge-requiring VQA, we further expand the areas of knowledge required. Our dataset requires diverse forms of outside knowledge including explicit fact-based knowledge that is likely to be contained in knowledge bases,

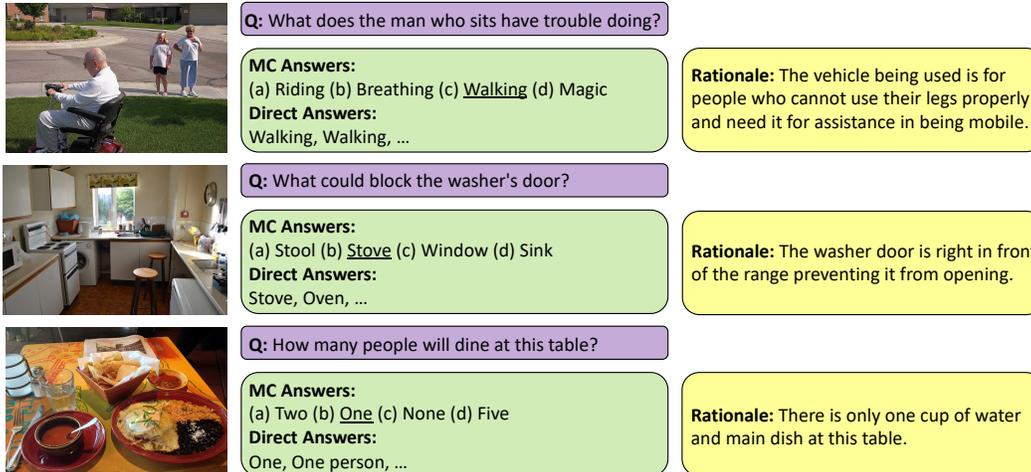


Figure 1: **A-OKVQA dataset**. The dataset includes questions that require reasoning using a variety of knowledge types such as commonsense, world knowledge and visual knowledge. We provide Multiple-Choice (MC) as well as Direct Answer evaluation settings. There is a rationale associated to each question in the train set providing the explanation/knowledge for answering the question.

commonsense knowledge about human social behavior, knowledge about the physics of the world, and visual knowledge. Not only do we expand the variety of knowledge our agent needs, but we also increase the complexity of reasoning systems needed to answer questions. We need models to recognize the image, understand the question, recall relevant knowledge, and use reasoning to arrive at an answer. For instance, in the first question shown in Figure 1, the model should reason that people use that type of cart to avoid walking. Therefore, the old man should have trouble walking. In general, our dataset requires additional types of world knowledge compared to our previous work OK-VQA [36]. Hence, we call it Augmented OK-VQA (A-OKVQA).

A-OKVQA is composed of about 25K questions paired with both multiple choice (MC) answer options and ten free-form answers to allow for direct answer (DA) evaluation. The MC component of the dataset bypasses many difficulties inherent in direct answer evaluation and allows for a simple, clean accuracy score. This is particularly helpful given the greater variety in answers in A-OKVQA questions. At the same time, we believe direct answer evaluation is important to encourage models with more real-world applicability. In addition to the questions and answers, we provide *rationales* for each question. This is to allow for models to use this extra annotation to train reasoning or knowledge retrieval methods or to build more explainable VQA models. The rationales also validate that both reasoning and knowledge are required for answering questions in the dataset.

In this work, our contributions are: (i) A new benchmark VQA dataset requiring diverse sources of outside knowledge and reasoning; (ii) A detailed analysis of the dataset that highlights its diversity and difficulty; (iii) An evaluation of a variety of recent baseline approaches in the context of the challenging questions in A-OKVQA; (iv) An extensive analysis of the results leading to interesting findings (e.g., how well models perform when answers are in the tail of the distribution, and also the complementarity of the studied models). We will release this dataset publicly.

2 Related Work

Visual Question Answering. Visual Question Answering (VQA) has been a common and popular form of vision and language reasoning. Many datasets on this task have been proposed [34, 2, 13, 65, 47, 69, 55, 27] but most of these do not require much outside knowledge or reasoning, often focusing on recognition tasks such as classification, attribute detection and counting.

Knowledge-based VQA datasets. Several previous works have studied the problem of knowledge-based VQA. The earliest explicitly knowledge-based VQA datasets were KB-VQA [56] and FVQA [57]. While these benchmarks did specifically require knowledge for questions, the knowledge required for these benchmarks is completely “closed”. FVQA [57] is annotated by selecting a triplet

from a fixed knowledge graph. This forces the questions to require knowledge, but because the question is written based on this knowledge, these questions are fairly trivial once the knowledge is known and do not require much reasoning. In addition, the knowledge required is explicitly closed to the knowledge graphs used to generate the dataset, so these datasets can only test knowledge retrieval on those specific graphs. KVQA [48] is based on images in Wikipedia articles. Because of the source of the images, these questions tend to mostly test recognizing specific named entities (e.g., Barack Obama) and then retrieving Wikipedia knowledge about that entity rather than commonsense knowledge.

Most similar to our work is OK-VQA [36]. This dataset was an improvement over prior work in terms of scale, and the quality of questions and images. It also has the property that the required knowledge was not “closed” or explicitly drawn from a particular source, and could be called “open”-domain knowledge. While this is an improvement over the previous works, it still suffers from problems which we address in this work. The knowledge required, while “open” is still biased towards simple lookup knowledge (e.g., what is the capital of this country?) and most questions do not require much reasoning. In contrast, our dataset is explicitly drawn to rely on more common-sense knowledge and to require more reasoning to solve. In addition, our dataset includes “rationale” annotations, which allow knowledge-based VQA systems to more densely annotate their knowledge acquisition and reasoning capabilities. S3VQA [23] analyzes OK-VQA and creates a new dataset which includes questions that require detecting an object in the image, replacing the question with the word for that object and then querying the web to find the answer. Like OK-VQA, it even more explicitly has the problem of questions usually requiring a single retrieval rather than much commonsense knowledge or reasoning.

Another related line of work is Visual Commonsense Reasoning (VCR) [66] and VisualCOMET [39]. VCR is also a VQA dataset, but is collected from movie scenes and is quite focused on humans and their intentions (e.g. “why is [PERSON2] doing this”), whereas our dataset considers questions and knowledge about a variety of objects. Similarly, VisualCOMET tests commonsense language and vision models on a movie dataset, but its expected output is a scene graph for the image (e.g., “After, [PERSON] is likely to”). Additionally, the Ads Dataset [22] is a dataset requiring knowledge about the topic and sentiments of the ads. Other datasets have considered knowledge-based question answering for a sitcom [14] and by using web queries [9].

Explanation / Reasoning VQA. Visual reasoning on its own has been studied in several VQA datasets. In CLEVR [25], the image and question are automatically generated from templates and explicitly require models to go through multiple steps of reasoning to correctly answer. This dataset and similar datasets which rely on simulated images suffer from lack of visual realism and lack of richness in the images and questions and are thus prone to be overfit to with methods achieving nearly 100% accuracy [64]. Our dataset requires reasoning on real images and free-form language. Other works [38, 28] have collected or extracted justifications on the VQAv2 [16] dataset. However, VQAv2 mostly focuses on questions about object attributes, counting and activities, which do not require reasoning on outside knowledge.

Knowledge / Commonsense in NLP. Question answering with knowledge and commonsense is also a well-studied problem in natural language processing. This takes the form of knowledge base completion [6], knowledge-based question answering [46] to open-domain question answering [10]. [4, 63, 5] address question answering from specific knowledge sources. This includes open-domain question answering [10, 58, 61, 60, 52] and question answering from Wikipedia SQu-AD [46, 45]. Much work has also been done in commonsense question answering as in CommonsenseQA [53], where there is no direct source of knowledge but the agent must have general “commonsense” to answer the question.

3 A-OKVQA Collection

Image source. The first requirement of an image source for this knowledge-based VQA task is that it has an abundance of visually rich and interesting images. Images containing a small number of objects are typically quite challenging to write questions requiring outside knowledge to answer. We used images from the 2017 partitioning of the COCO dataset [29] in the creation of A-OKVQA because: (1) it has many images cluttered with multiple objects and entity types, (2) it is an established dataset with many associated models already in existence. To ensure suitable images for annotation, we do

some additional filtering to remove uninteresting images: For the training and validation sets, we define images with more than three objects as “interesting” and select those for question writing. For the test set, which lacks object annotation, we train a ResNet-50 classifier to distinguish “interesting” images based on this criteria, achieving an accuracy of 78% on the validation set. After multiple rounds of filtering (described below), we obtain 23.7K unique images.

Question collection & filtering. The questions in A-OKVQA were written and refined over several rounds of annotation by 437 crowd-workers on the Amazon Mechanical Turk platform and refined through several manual and automated filtering steps to increase overall quality. As a first quality assurance measure, workers completed a qualification task to demonstrate their ability to write questions that met our criteria, namely that questions require: (1) looking at the image to answer, (2) some commonsense or specialized knowledge, (3) some thinking beyond merely recognizing an object, and (4) not be too similar to previous questions.

To help ensure the last point, we clustered images by CLIP [42] visual features and batched similar images together so that the same worker wrote questions sequentially for related images (e.g., a worker might write questions for several images showing baseball games in one task) to cut down on repetitive questions. As an added measure to encourage question diversity, we maintained a database of questions written and required users to check a new question against these by displaying the five previous questions most similar in terms of their RoBERTa [31] embeddings. We used Pythia [51] pre-trained on VQAv2 as a first automated check for questions we considered trivial, removing any question for which the model predicted the correct answer choice. Next, questions were screened by three other workers and only included if the majority agreed that it met our criteria for inclusion. In all, 37,687 questions, or 60% of post-qualification questions were excluded from the dataset by this process. After questions and their multiple choice options were complete, nine additional free-form answers were collected for each question by a separate pool of workers.

Rationales. After questions and multiple-choice answer options were collected and validated, we initiated a task to collect rationales. Workers were given a question and answer options and asked to explain in one to two simple sentences why a particular answer was correct, including any necessary facts or knowledge about the world not contained in the images. Workers were given examples and went through a qualification process to assure high-quality output. For each question, we collected three rationales.

4 Dataset Statistics

Question/Answer/Rationale statistics. After all rounds of annotation, the A-OKVQA dataset contains 24,903 Question/Answer/Rationale triplets, split into 17.1K/1.1K/6.7K for train, validation and test. These preserve the COCO 2017 train/val/test splits. The average length of the questions, answers, and rationales, and the number of their unique words are shown in Table 1.

In Figure 2a we show the distribution of answer options in our dataset. What we see is a fairly typical long-tail distribution of labels, as is seen in many open-labeled image tasks [68]. A few answers occur quite often in the dataset, but overall, most answers in the dataset fall into the long tail of the distribution.

We are also interested to know the amount of overlap in the answer set between the train set and val and test sets. We find that in the val set, the ground-truth answer for 87.6% of questions appears in the train set while in the test set 82.7% do. This shows that there is indeed some reasonable similarity between the sets, but also that a significant portion of the held out sets require an answer that the model will not have seen during training, requiring the model to be able to generate out-of-distribution answers or generate answers based on some knowledge outside of the dataset.

Comparison with other datasets. In Table 1 we show dataset properties and statistics for A-OKVQA compared to related datasets. We see that compared to the more knowledge-focused natural image datasets such as OK-VQA, we have between 2-10x more questions while VCR (focused on images of people in movies) has about 10x as ours. This is unsurprising because we intentionally filter similar questions, making our questions more diverse (see Table 2), but difficult to collect at scale. Our dataset has annotations for both multiple choice and direct answer evaluation. Our dataset also has rationales, unlike OK-VQA, S3VQA and KB-VQA. FVQA has knowledge tuples as rationales rather than full sentences. Most similar to our rationales is VCR. Unlike all of these, we collect 3 rationales

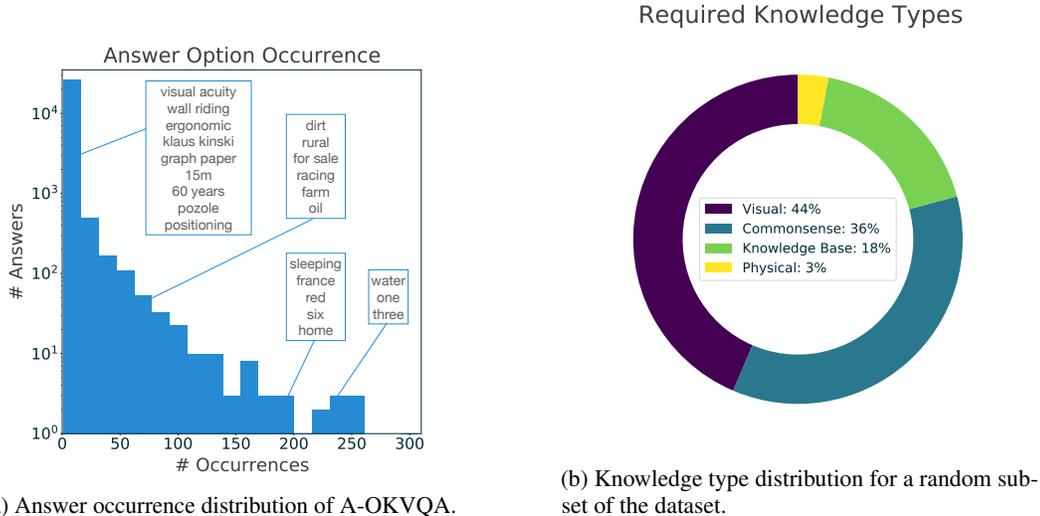


Figure 2: Dataset statistics.

	Q	I	Rationale	Knowledge type	Ans type	Avg. length (Q/A/R)	unique words (Q/A/R)
KB-VQA [56]	2,402	700	✗	fixed KB	DA	6.8/2.0/NA	530/1,296/NA
FVQA [57]	5,826	2,190	✓	fixed KB	DA	9.5/1.2/NA	3,010/1,287/NA
OK-VQA [36]	14,055	14,031	✗	factoid	DA	8.1/1.3/NA	5,703/11,125/NA
S3VQA [23]	7,515	7,515	✗	factoid	DA	12.7/2.8/NA	7,515/8,301/NA
VCR [66]	290k	99,904	✓	people actions	MC	8.7/7.7/16.8	11,254/18,861/28,751
A-OKVQA	24,903	23,692	✓	common/world	DA/MC	8.8/1.3/11.0	7,248/17,683/20,629

Table 1: Comparison of various knowledge-based VQA datasets. Data based on publicly reported numbers and/or our analysis of publicly available annotations (therefore some answer statistics may exclude test sets). Answer statistics for A-OKVQA based on the direct answer set. Q: question, I: image, A: answer, R: rationale, DA: Direct Answer, MC: Multiple Choice. NA indicates lack of rationales or, for FVQA are KB triplets, so we do not compare the lengths.

rather than just 1 as the rationales are more knowledge-based and have more possible variation within the same question. Our average question lengths are long compared to many of these datasets except for S3VQA which has the longest and VCR which is about on par. Ours also contains the most unique words except for VCR, likely because that dataset has more questions.

Knowledge types. The most significant factor differentiating our dataset is the kind of knowledge required. Datasets such as FVQA have fixed knowledge bases that are used to write the questions, and so the knowledge required can be found in e.g. ConceptNet [30] directly. OK-VQA and S3VQA focus on more factoid knowledge (e.g., years of invention or countries of origin). In S3VQA in particular, researchers found that these datasets take the form of finding an entity in the image and/or question and searching and retrieving knowledge about that particular entity (see [23]). VCR is overwhelmingly images of people interacting in television shows and movies and requires images to have people in them. Thus, the required knowledge is very focused on commonsense about human behavior and intentions. In our dataset, we require broader areas of knowledge including the factoid knowledge likely to be contained in knowledge bases (as in FVQA, KBVQA, OKVQA and S3VQA), and commonsense knowledge (like VCR but broader than just about people).

To analyze the knowledge required in A-OKVQA more quantitatively, we annotated a randomly sampled subset of 1,000 questions in the A-OKVQA test. In this experiment, we ask the annotators to label what kind of knowledge type was required to answer the questions. The choices were: (1) **Commonsense** knowledge about human social behavior (e.g., that many donuts being made in a cart implies they are for sale rather than for personal consumption), (2) **Visual knowledge** (e.g., muted color pallets are associated with the 1950s), (3) **Knowledge bases** (e.g., hot dogs were invented in Austria), (4) **Physical knowledge** about the world that humans learn from their everyday experiences

Dataset	Mean Q distance	Mean rationale distance
FVQA [57]	0.6199	✗
VCR [66]	0.7095	0.8017
KB-VQA [56]	0.7192	✗
S3VQA [23]	0.8050	✗
VQAv2 [16]	0.8405	0.8228
OK-VQA [36]	0.8428	✗
A-OKVQA	0.8564	0.8779

Table 2: **Question and Rationale Diversity.** Mean pairwise cosine distances in a sentence transformer space for various datasets. ✗ indicates lack of rationale. We choose one rationale per question on A-OKVQA to make the comparison to other datasets with only one rationale. Rationales for VQAv2 come from the VQA-X dataset [38].

(e.g., shaded areas have a lower temperature than other areas). The distribution is shown in Figure 2b. Most of our questions focus around commonsense and visual knowledge. It should be noted that sometimes there is no clear distinction between these two categories, and a question can belong to either category.

Question diversity. To analyze the diversity of A-OKVQA compared to other datasets, as a proxy, we use the average pairwise cosine distance between the questions in the dataset. We run our questions through a sentence transformer¹ and compute the cosine distance between all pairs in the dataset. We then take the average of these. We see from Table 2 that our dataset has the most diversity on this metric. In particular, we see a large difference compared to VCR which has many similar questions such as “What is going to happen next?” and questions relating to what specific people in the scene are doing and why. We also compare the diversity of rationales to VCR and VQAv2 (using rationales from VQA-X [38] rationales). We also find that our rationales are much more diverse than in these datasets. Qualitatively, we also find that our dataset tends to have much more varied questions because it is taken from the more visually diverse COCO dataset (a quality shared by OK-VQA and VQAv2 which do almost as well on this metric) and requires more diverse kinds of knowledge.

Finally, we use the same mean pairwise distance to look in particular at how different our questions are from OK-VQA which is the most similar prior work to ours. To do this we compare the minimum pairwise distance between every question in the OK-VQA training set to every question in the OK-VQA test set and our test set. We find that the average minimum distance from OK-VQA train to test is **0.256** compared to **0.311** between OK-VQA train and our test set². This shows that there is in fact a significant difference between our question set and OK-VQA in this feature space.

5 Experiments

Next, we benchmark the A-OKVQA dataset and compare the performance of different models. We consider three classes of methods: (1) **large-scale pre-trained models** such as CLIP [42] and GPT-3 [7], (2) **models that generate and use rationales**, and (3) **specialized models** that are designed for knowledge-based VQA (KRISP [35]) or tested for VQA (e.g., ViLBERT [32]).

5.1 Evaluation

In the *multiple choice (MC)* setting, a model chooses its answer from one of four options and we compute accuracy as the evaluation metric. In the *direct answer (DA)* setting, a model can generate any text as its answer and we use the standard VQA evaluation from [2].

5.2 Large-scale Pre-trained Models

We compare three types of large-scale pre-trained models (discriminative, contrastive, and generative) in Table 3. We also test these models in different input settings (where questions, images, or both are provided).

¹Specifically multi-qa-MiniLM-L6-cos-v1 [19] to avoid overlap with RoBERTa.

²To make this comparison even, we chose a random subset of our test set to be the same size as OK-VQA test set so that the minimum is over the same number of possible choices in both cases.

Method	Multiple Choice		Direct Answer	
	Val	Test	Val	Test
(a) Random	26.70	25.36	0.03	0.06
(b) Random (weighted)	29.49	30.87	0.15	0.10
(c) Most Common	30.70	30.33	1.75	1.26
Question				
(d) BERT [12] (classifier)	32.93	33.54	9.52	8.41
(e) CLIP [42] (classifier)	32.74	33.54	13.10	10.24
(f) CLIP [42] (zero-shot)	30.42	30.58	0.44	0.57
(g) CLIP [42] (contrastive)	37.40	38.58	5.56	3.83
(h) GPT-3 [7]	35.07	35.21	12.98	11.49
Image				
(i) ResNet [17] (classifier)	28.19	28.81	2.68	2.30
(j) CLIP [42] (classifier)	33.21	32.56	5.15	4.38
(k) CLIP [42] (zero-shot)	56.28	53.94	2.24	2.29
(l) CLIP [42] (contrastive)	52.56	50.09	2.33	2.45
Question & Image				
(m) CLIP (classifier)	40.84	38.30	18.95	14.27
(n) CLIP (zero-shot)	48.19	45.72	1.08	0.71
(o) CLIP (contrastive)	53.77	51.01	10.36	7.10
(p) ClipCap [37]	56.93	51.43	30.89	25.90

Table 3: **Large-scale pre-trained models.** We also compare with no input heuristics (rows *a-c*) with choices (for MC) or vocabulary answers (for DA). *Random* is a uniform sampling. *Random (weighted)* uses weighted sampling proportional to correct answer frequencies in train. *Most Common* selects the most frequent answer in train.

We compute BERT [12, 20] and CLIP ViT-B/32 text encoder representations for questions. We also compute ResNet-50 [17] and CLIP ViT-B/32 features for images. These are provided as inputs to the appropriate discriminative and contrastive models. We provide questions as tokens and CLIP RN50x4 image representations as inputs to the generative models. We generate a vocabulary from a subset of training set answers and choices to use across all appropriate models. We describe this vocabulary further in Appx. B.

Discriminative models. We train a multi-label linear classifier (i.e. MLP with one hidden layer and sigmoid activation function) on top of BERT (row *d*), ResNet (row *i*), and CLIP (rows *e/j/m*) representations to score answers from the vocabulary. When questions and images are both provided, we first concatenate their representations. For the DA setting, we predict the top scoring vocabulary answer. For the MC setting, we instead predict the nearest neighbor³ choice to the top scoring vocabulary answer.

Contrastive models. We also evaluate models which match input questions and/or images with answers using their CLIP encodings. First, we evaluate the zero-shot setting (rows *f/k/n*). If both questions and images are provided as inputs, we first add their representations. We select the answer whose encoding has the greatest cosine similarity to our input representation. We select from vocabulary answers in DA and the given choices in MC.

We also train a single-layer MLP on top of our input representations (rows *g/l/o*). If both questions and images are provided, we first concatenate their representations. Our MLP produces a 512-d embedding and we train this with a CLIP-style contrastive loss between embeddings and their corresponding answers. We describe this loss further in Appx. B. We repeat the evaluation from the zero-shot setting, using these learned embeddings.

Generative models. We also evaluate models (GPT-3 [8] and ClipCap [37]) that generate answers directly as text. For both models, we predict the generated text for DA and the generated text’s nearest neighbor choice for MC.

We prompt GPT-3⁴ (row *h*) with 10 random questions and answers from the training set, followed by a new question, and let GPT-3 generate an answer to that question, in a manner similar to [59]. We

³Cosine similarity between mean GloVe [40, 21] word embeddings.

⁴We use the second largest available GPT-3 model, Curie, as in [59].

Method	Multiple Choice		Direct Answer	
	Val	Test	Val	Test
(a) ClipCap → Cap. → GPT	42.51	43.61	16.59	15.79
(b) ClipCap → Ratl. → GPT	44.00	43.84	18.11	15.81
Oracles				
(c) GT Caption → GPT	45.40	—	16.39	—
(d) GT Rationale → GPT	56.74	56.75	24.02	20.75

Table 4: **Models using generated and GT rationales** as described in Sec. 5.3. We are unable to evaluate the GT Caption → GPT setting on the test set, as captions are not available in the COCO [11] test set.

provide GPT-3 with the prompt template “Question: ... Answer: [...]”, expecting it to complete the answer for each evaluation question.⁵

ClipCap [37] (row *p*) is an image captioning method that passes CLIP image features through a trained network to GPT-2 (as input tokens). We adapt this model by adding question tokens (and answer choices if applicable) to the prompt of GPT-2, generate answers instead of captions, and fine-tune on our data. We provide additional details, diagrams, and variations in Appx. B.

Results. Table 3 shows the results of our evaluation of these models. Rows *a-c* show the biases in our dataset, but that the direct answer setting is appropriately challenging. Question-only baselines (rows *d-h*) show poor performance in both MC and DA settings. However, it is interesting that GPT-3 performs similarly to the fine-tuned CLIP models (whichever is better per setting). The zero-shot CLIP model (row *f*) is least effective, indicating that training is necessary to repurpose CLIP text encodings for language-only tasks. Unsurprisingly, CLIP image features are very strong for zero-shot multiple choice matching (row *k*). However, they are not as strong as for the fine-tuned classifier (row *j*) in DA. ClipCap (row *p*) outperforms all other baselines in DA, because we use powerful image features and also fine-tune a strong language model for our task.

5.3 Rationale Generation

We are interested in whether we can improve GPT-3 prompting results by providing additional image- and question- specific context and report results for the following methods in Table 4. So, we fine-tune ClipCap (given images and questions, but not choices) as above, but for the task of generating rationales instead of answers. Our model scores **10.2** (val) / **9.58** (test) on SacreBLEU [41] and **0.271** (val) / **0.256** (test) on METEOR [3]. We can then prompt GPT-3 (as above) but also provide these generated rationales as “Context: ...”. This model is denoted by ‘ClipCap → Ratl. → GPT’. We provide additional details, diagrams, and examples of generated rationales in Appx. C. We repeat this experiment using captions (generated from only images) from the original ClipCap model: ‘ClipCap → Cap. → GPT’.

Results. We show results from these experiments in Table 4. Interestingly, prompting GPT-3 with ground-truth rationales (row *d*) is competitive with the best model in Sec. 5.2 (Table 3, row *p*) in MC and significantly outperforms the question-only GPT-3 method (Table 3, row *h*). When we prompt GPT-3 with ground-truth rationales (row *d*), we see higher performance than when we provide ground-truth captions (row *c*). This affirms that rationales contain useful information (i.e. specific to our questions and answers) in addition to captions. However, the additional performance of prompting GPT-3 using generated rationales (row *b*) over generated captions (row *a*) is not as significant. This indicates potential room for improvement in our approach for generating rationales.

5.4 Specialized Models

In this section, we evaluate some recent high-performing, open-source models trained on knowledge-based VQA or the traditional VQA. The models we consider are Pythia [24], ViBERT [32], LXMERT [54], KRISP [35], and GPV-2 [26]. As the first four models are part of MMF [50], it is easier to compare them fairly. KRISP is a high-performing model on OK-VQA [36]. It provides a suitable baseline as it addresses knowledge-based VQA. GPV-2 performs multiple vision and

⁵During MC, we also tried prompting GPT-3 with “Choices: ...”, but find that this actually hurts performance.

Method	Multiple-Choice		Direct Answer	
	Val	Test	Val	Test
(a) Pythia [24]	49.0	40.1	25.2	21.9
(b) ViLBERT [32] - OK-VQA	32.8	34.1	9.1	9.2
(c) ViLBERT [32] - VQA	47.7	42.1	17.7	12.0
(d) ViLBERT [32]	49.1	41.5	30.6	25.9
(e) LXMERT [54]	51.4	41.6	30.7	25.9
(f) KRISP [35]	51.9	42.2	33.7	27.1
(g) GPV-2 [26]	60.3	53.7	48.6	40.7
Oracles				
(h) GPV-2 [26] + Masked Ans.	65.1	58.3	52.7	43.9
(i) GPV-2 [26] + GT Ratl.	73.4	67.2	58.9	51.7

Table 5: **Specialized models results.** Baselines trained for VQA or knowledge-based VQA, and fine-tuned on A-OKVQA. The bottom two rows are not comparable with the others since they use ground-truth rationales at test time.

vision–language tasks and has learned a large number of concepts, so it can be a strong baseline for A-OKVQA. All of these models are fine-tuned on A-OKVQA to predict answers directly for DA evaluation. We adapt them to MC using the nearest choice method described above. See Appx. D for the details of each model.

Results. Unsurprisingly, these models, which are specialized for DA and some of which are specialized for knowledge-based VQA perform very well on the DA evaluation and quite well on MC. Of the models trained only on A-OKVQA KRISP does the best, likely because it is trained to directly use outside knowledge graphs. GPV-2, however, performs best of all, beating all other models (that do not use ground-truth rationales) in all settings, possibly because of the large number of concepts it has learned.

Transfer results. We train ViLBERT on VQAv2 and OK-VQA datasets (denoted by ‘ViLBERT-VQA’ and ‘ViLBERT-OK-VQA’ in Table 5) to evaluate whether the knowledge from those datasets is sufficient for A-OKVQA. The low performance shows the significant difference between these datasets.

Ground-truth Rationales. To evaluate how well the model performs if it is provided with high-quality rationales, we use ground-truth rationales at test. We show these results with GPV-2 (our best model). Ground-truth rationales are appended to questions as additional input text (‘GPV-2 + GT Ratl.’). For this experiment, we used only one of the rationales. Comparing rows *g* and *i* of Table 5 shows rationales are helpful. To evaluate how much of this improvement can be attributed to rationales and not the fact that sometimes rationales contain the answer, we replaced answers in the rationales with [answer] token. The performance drops (row *i* vs row *h*), however, it is still higher than the case that we do not use rationales (row *h* vs row *g*).

6 Analysis of Models

Next, we analyze the predictions that our baseline models make to see if we can learn more about A-OKVQA: what kinds of questions do different types of approaches do better / worse on? For these experiments, we choose some of the best performing models on Direct Answer: ViLBERT [32], LXMERT [54], KRISP [35], ClipCap [37] and GPV-2 [26]. We also use the ClipCap → Rationale → GPT model from Table 4, which will be referred to as ‘GR-GPT’ for Generated Rationales GPT.

Answer Frequency. First, we look at how answer frequency affects performance in Table 6. We first count the number of times any answer appears in the direct answers in the training set. We then divide these into bins and look at the direct DA test accuracy of our baselines for each of these frequency bins. We find that GPV-2, and to a lesser extent ClipCap and GR-GPT perform better on questions whose answers do not appear often in the training set (1-5 and 6-10 columns of Table 6). GPV-2 in particular (which is fine-tuned on several vision and language tasks) is able to predict these tail answers much better than other methods, especially the discriminative methods such as LXMERT.

Knowledge Types. Next, we use the subset of test that we collected knowledge types on (see Sec. 4) to look at the direct answer accuracy of these models for different types of knowledge. In Table 7, we

Model	1-5	6-10	11-20	21-50	51-100	101-200	201+
VilBERT [32]	0.00	0.00	3.68	10.97	19.95	26.53	35.91
LXMERT [54]	0.00	0.00	4.29	13.73	20.18	26.69	34.31
KRISP [35]	0.00	0.61	6.34	13.99	21.78	28.55	35.22
ClipCap [37]	4.71	4.24	9.10	17.90	25.93	29.44	33.99
GR-GPT	8.18	9.29	9.41	17.39	18.31	21.98	24.65
GPV-2 [26]	10.16	12.12	22.60	31.04	38.40	41.60	44.69

Table 6: **Results across different answer frequencies.** The questions are categorized based on the frequency of the GT answer in the training set. Columns show accuracy for answers that appear 1-5 times, 6-10 times, etc. If multiple direct choices, we default to most common one.

Model	Commonsense	Knowledge Base	Physical Knowledge	Visual Knowledge
VilBERT [32]	24.30	19.96	29.76	26.55
LXMERT [54]	25.51	16.01	27.38	27.23
KRISP [35]	26.63	20.72	39.29	26.09
ClipCap [37]	27.19	16.57	30.95	33.41
GR-GPT	21.42	12.99	17.86	24.79
GPV-2 [26]	39.76	25.24	44.05	41.19

Table 7: **Analysis of results based on knowledge type.**

see that while again GPV is the best overall and in every category, the results show some interesting distinctions. KRISP, which is specifically designed with access to explicit knowledge sources such as ConceptNet [30] performs better on “Knowledge Base” questions compared with other discriminative multi-modal transformer methods such as VilBERT and LXMERT as well compared to ClipCap which has an overall higher performance. It also performs better on “Physical Knowledge” which also tends to overlap with the knowledge sources it has.

Prediction overlap/difference. Finally, we look at some statistics on a question by question level in the A-OKVQA test set. Specifically we look at the overlap in which methods answered which questions correctly⁶. We use the same models as in Tables 6 & 7.

First, we find that only **5.85%** of questions in test were answered correctly by all models and **30.96%** of questions had no model predict a correct answer for. Considering the worst performing model of these gets **15.81%** DA accuracy and the best gets **40.7%**, it implies that there is actually a large variation between these models beyond some just being generally better than others and thus getting “hard” questions right and keeping performance on “easy” questions.

In Table 8, we show the difference between the questions each model gets right on A-OKVQA test. Each row shows the percentage of that method’s correctly answered questions that were not correctly answered by the comparison model in each column. If we look at the row for the lowest performing model (GR-GPT) for the column for the best performing model (GPV-2), we still see that **29.2%** of GR-GPT’s correctly answered questions are answered wrong by GPV-2!

Finally, to further illustrate the point that different models have very different mistake patterns, we take the prediction of all of these models except for GPV-2 for each question and take the majority vote between these. This majority vote combination gets an accuracy of **29.5** compared to the best of these models which gets **27.1**. This does not work when GPV-2 is added (this majority model gets **35.60** which is lower than GPV-2’s **40.7**). We can also look at the Oracle combination accuracy. That is, from our six models, choose the answer with the highest ground-truth value and take that as the oracle combination answer. This DA accuracy is **56.87** versus the single best performance of **40.7**, again showing that even worse performing models get lots of questions right that the best model gets wrong.

Qualitative Analysis. We analyzed our models and extracted questions that all of the discussed models fail at. Figure 3 shows an example from each knowledge type. This qualitative example shows what type of reasoning is missing in our current top performing models.

⁶For ease of analysis we count a binary yes/no of whether a model answered correctly if it answered any possible answer in the direct answer set.

Model	ViBERT	LXMERT	KRISP	ClipCap	GR-GPT	GPV-2
ViBERT [32]	0.00	29.00	27.19	43.72	59.72	26.33
LXMERT [54]	28.07	0.00	26.57	44.39	59.73	27.44
KRISP [35]	30.44	30.76	0.00	44.18	60.29	27.43
ClipCap [37]	48.72	49.98	46.76	0.00	55.94	26.64
GR-GPT	50.27	50.91	48.67	40.30	0.00	29.20
GPV-2 [26]	51.09	52.46	49.57	46.56	61.94	0.00

Table 8: **Pairwise difference between correctly answered questions.** For row i and column j of this table the value is percentage of questions answered correctly by model i that j did not answer correctly.

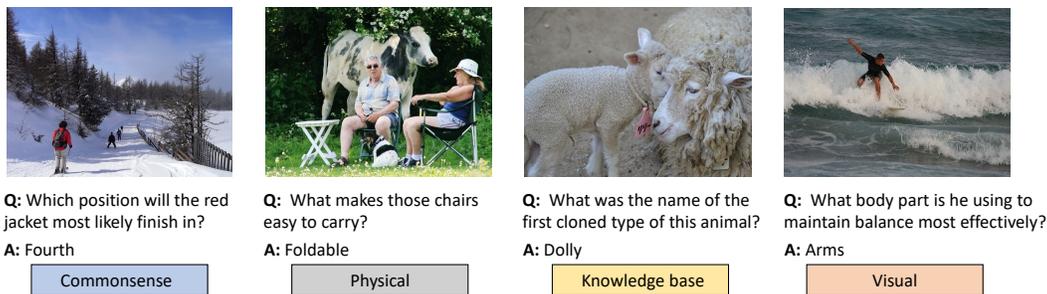


Figure 3: **Example questions that all discussed models fail at.**

All of these analyses together provide several interesting findings. First, aside from being generally difficult, the A-OKVQA dataset shows a surprising lack of overlap in the specific questions different models answer correctly. Second, we see that different methods handle rare answers very differently. Moreover, different methods perform differently based on the type of knowledge. All of this suggests that A-OKVQA provides many different kinds of challenging questions which bring out different strengths and weaknesses of methods.

7 Conclusion

Vision and language models have become progressively more powerful, however, evaluation of the reasoning capabilities of these models have not received adequate attention. To take a step in this direction, we propose a new knowledge-based VQA benchmark called A-OKVQA, which primarily includes questions that require reasoning using commonsense and world knowledge. We provide *rationales* for each question so models can learn the line of reasoning that leads to the answer. We evaluate a large set of recent, high performance baselines. While they show impressive performance on the proposed task, it is evident that they lack the reasoning capability and/or the knowledge required to answer the questions, and there is a large room for improvement. Through extensive analyses, we show different models have different weaknesses and strengths. To solve A-OKVQA and to move towards general multi-modal intelligence, we need to combine many types of capabilities from many different methods.

References

- [1] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *CVPR*, 2018. 19
- [2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. VQA: visual question answering. In *ICCV*, 2015. 1, 2, 6
- [3] Satyanjeev Banerjee and Alon Lavie. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, 2005. 8
- [4] Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. Semantic parsing on freebase from question-answer pairs. In *EMNLP*, 2013. 3

- [5] Antoine Bordes, Sumit Chopra, and Jason Weston. Question answering with subgraph embeddings. In *EMNLP*, 2014. 3
- [6] Antoine Bordes, Jason Weston, Ronan Collobert, and Yoshua Bengio. Learning structured embeddings of knowledge bases. In *AAAI*, 2011. 3
- [7] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, T. J. Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *NeurIPS*, 2020. 6, 7
- [8] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *NeurIPS*, 2020. 1, 7
- [9] Yingshan Chang, Mridu Baldevraj Narang, Hisami Suzuki, Guihong Cao, Jianfeng Gao, and Yonatan Bisk. WebQA: Multihop and multimodal qa. *arXiv*, 2021. 3
- [10] Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. Reading Wikipedia to answer open-domain questions. In *ACL*, 2017. 3
- [11] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO Captions: Data collection and evaluation server. *arXiv*, 2015. 8, 19
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 2019. 1, 7
- [13] Haoyuan Gao, Junhua Mao, Jie Zhou, Zhiheng Huang, Lei Wang, and Wei Xu. Are you talking to a machine? dataset and methods for multilingual image question. In *NeurIPS*, 2015. 2
- [14] Noa García, Mayu Otani, Chenhui Chu, and Yuta Nakashima. KnowIT VQA: Answering knowledge-based questions about videos. In *AAAI*, 2020. 3
- [15] Donald Geman, Stuart Geman, Neil Hallonquist, and Laurent Younes. Visual turing test for computer vision systems. *Proceedings of the National Academy of Sciences*, 2015. 1
- [16] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V in VQA matter: Elevating the role of image understanding in visual question answering. In *CVPR*, 2017. 1, 3, 6, 19
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 7
- [18] Drew A. Hudson and Christopher D. Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *CVPR*, 2019. 19
- [19] HuggingFace. <https://huggingface.co/sentence-transformers/multi-qa-MiniLM-L6-cos-v1>. 6
- [20] HuggingFace. <https://huggingface.co/sentence-transformers/nli-bert-base>. 7
- [21] HuggingFace. https://huggingface.co/sentence-transformers/average_word_embeddings_glove.6B.300d. 7
- [22] Zaeem Hussain, Mingda Zhang, Xiaozhong Zhang, Keren Ye, C. Thomas, Zuha Agha, Nathan Ong, and Adriana Kovashka. Automatic understanding of image and video advertisements. In *CVPR*, 2017. 3
- [23] Aman Jain, Mayank Kothiyari, Vishwajeet Kumar, Preethi Jyothi, Ganesh Ramakrishnan, and Soumen Chakrabarti. Select, substitute, search: A new benchmark for knowledge-augmented visual question answering. In *SIGIR*, 2021. 3, 5, 6
- [24] Yu Jiang, Vivek Natarajan, Xinlei Chen, Marcus Rohrbach, Dhruv Batra, and Devi Parikh. Pythia v0.1: the winning entry to the VQA challenge 2018. *arXiv*, 2018. 8, 9, 19
- [25] Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning. In *CVPR*, 2017. 1, 3
- [26] Amita Kamath, Christopher Clark, Tanmay Gupta, Eric Kolve, Derek Hoiem, and Aniruddha Kembhavi. Webly supervised concept expansion for general purpose vision models. *arXiv*, 2022. 8, 9, 10, 11, 19, 20
- [27] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. Visual Genome: connecting language and vision using crowdsourced dense image annotations. *IJCV*, 2017. 2, 19
- [28] Qing Li, Jianlong Fu, Dongfei Yu, Tao Mei, and Jiebo Luo. Tell-and-answer: Towards explainable visual question answering using attributes and captions. In *EMNLP*, 2018. 3
- [29] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014. 3
- [30] Hugo Liu and Push Singh. ConceptNet—a practical commonsense reasoning tool-kit. *BT technology journal*, 2004. 5, 10
- [31] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A robustly optimized bert pretraining approach. *arXiv*, 2019. 4

- [32] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *NeurIPS*, 2019. 1, 6, 8, 9, 10, 11, 19
- [33] Jiasen Lu, Vedanuj Goswami, Marcus Rohrbach, Devi Parikh, and Stefan Lee. 12-in-1: Multi-task vision and language representation learning. In *CVPR*, 2020. 1
- [34] Mateusz Malinowski and Mario Fritz. A multi-world approach to question answering about real-world scenes based on uncertain input. In *NeurIPS*, 2014. 1, 2
- [35] Kenneth Marino, Xinlei Chen, Devi Parikh, Abhinav Kumar Gupta, and Marcus Rohrbach. KRISP: Integrating implicit and symbolic knowledge for open-domain knowledge-based VQA. In *CVPR*, 2021. 6, 8, 9, 10, 11, 19
- [36] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. OK-VQA: A visual question answering benchmark requiring external knowledge. In *CVPR*, 2019. 1, 2, 3, 5, 6, 8
- [37] Ron Mokady, Amir Hertz, and Amit H. Bermano. ClipCap: CLIP prefix for image captioning. *arXiv*, 2021. 7, 8, 9, 10, 11, 16
- [38] Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Anna Rohrbach, Bernt Schiele, Trevor Darrell, and Marcus Rohrbach. Multimodal explanations: Justifying decisions and pointing to the evidence. In *CVPR*, 2018. 3, 6
- [39] Jae Sung Park, Chandra Bhagavatula, Roozbeh Mottaghi, Ali Farhadi, and Yejin Choi. VisualCOMET: Reasoning about the dynamic context of a still image. In *ECCV*, 2020. 3
- [40] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. GloVe: Global vectors for word representation. In *EMNLP*, 2014. 7
- [41] Matt Post. A call for clarity in reporting BLEU scores. In *Conference on Machine Translation*, 2018. 8
- [42] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 4, 6, 7
- [43] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 2019. 1
- [44] Colin Raffel, Noam M. Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *JMLR*, 2020. 19
- [45] Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don't know: Unanswerable questions for SQuAD. In *ACL*, 2018. 3
- [46] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In *EMNLP*, 2016. 3
- [47] Mengye Ren, Jamie Kiros, and Richard S. Zemel. Exploring models and data for image question answering. In *NeurIPS*, 2015. 2
- [48] Sanket Shah, Anand Mishra, Naganand Yadati, and Partha Pratim Talukdar. KVQA: Knowledge-aware visual question answering. In *AAAI*, 2019. 3
- [49] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *ACL*, 2018. 19
- [50] Amanpreet Singh, Vedanuj Goswami, Vivek Natarajan, Yu Jiang, Xinlei Chen, Meet Shah, Marcus Rohrbach, Dhruv Batra, and Devi Parikh. MMF: A multimodal framework for vision and language research. <https://github.com/facebookresearch/mmf>, 2020. 8
- [51] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards VQA models that can read. In *CVPR*, 2019. 4
- [52] Haitian Sun, Bhuwan Dhingra, Manzil Zaheer, Kathryn Mazaitis, Ruslan Salakhutdinov, and William W. Cohen. Open domain question answering using early fusion of knowledge bases and text. In *EMNLP*, 2018. 3
- [53] Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In *NAACL-HLT*, 2019. 3
- [54] Hao Hao Tan and Mohit Bansal. LXMERT: Learning cross-modality encoder representations from transformers. In *EMNLP*, 2019. 8, 9, 10, 11, 19
- [55] Makarand Tapaswi, Yukun Zhu, Rainer Stiefelwagen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. MovieQA: understanding stories in movies through question-answering. In *CVPR*, 2016. 2
- [56] Peng Wang, Qi Wu, Chunhua Shen, Anthony R. Dick, and Anton van den Hengel. Explicit knowledge-based reasoning for visual question answering. In *IJCAI*, 2017. 1, 2, 5, 6
- [57] Peng Wang, Qi Wu, Chunhua Shen, Anton van den Hengel, and Anthony R. Dick. FVQA: fact-based visual question answering. *TPAMI*, 2017. 1, 2, 5, 6
- [58] Shuohang Wang, Mo Yu, Xiaoxiao Guo, Zhiguo Wang, Tim Klinger, Wei Zhang, Shiyu Chang, Gerald Tesauro, Bowen Zhou, and Jing Jiang. R3: Reinforced reader-ranker for open-domain question answering. In *AAAI*, 2018. 3
- [59] Peter West, Chandra Bhagavatula, Jack Hessel, Jena D Hwang, Liwei Jiang, Ronan Le Bras, Ximing Lu, Sean Welleck, and Yejin Choi. Symbolic knowledge distillation: from general language models to commonsense models. *arXiv preprint arXiv:2110.07178*, 2021. 7

- [60] Wei Yang, Yuqing Xie, Aileen Lin, Xingyu Li, Luchen Tan, Kun Xiong, Ming Li, and Jimmy Lin. End-to-end open-domain question answering with BERTserini. In *NAACL*, 2019. 3
- [61] Yi Yang, Wen-tau Yih, and Christopher Meek. WikiQA: A challenge dataset for open-domain question answering. In *EMNLP*, 2015. 3
- [62] Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu, Yumao Lu, Zicheng Liu, and Lijuan Wang. An empirical study of GPT-3 for few-shot knowledge-based VQA. *arXiv*, 2021. 1
- [63] Xuchen Yao and Benjamin Van Durme. Information extraction over structured data: Question answering with Freebase. In *ACL*, 2014. 3
- [64] Kexin Yi, Jiajun Wu, Chuang Gan, Antonio Torralba, Pushmeet Kohli, and Joshua B. Tenenbaum. Neural-symbolic VQA: Disentangling reasoning from vision and language understanding. In *NeurIPS*, 2018. 3
- [65] Licheng Yu, Eunbyung Park, Alexander C. Berg, and Tamara L. Berg. Visual madlibs: Fill in the blank description generation and question answering. In *ICCV*, 2015. 2
- [66] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *CVPR*, 2019. 1, 3, 5, 6
- [67] Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. VinVL: Revisiting visual representations in vision-language models. In *CVPR*, 2021. 1, 19
- [68] Xiangxin Zhu, Dragomir Anguelov, and Deva Ramanan. Capturing long-tail distributions of object subcategories. In *CVPR*, 2014. 4
- [69] Yuke Zhu, Oliver Groth, Michael S. Bernstein, and Li Fei-Fei. Visual7W: grounded question answering in images. In *CVPR*, 2016. 2, 19

A Additional details of dataset collection

A.1 Examples of rejected questions

With a focus on overall question quality, we removed around 60% of questions written for having any of several flaws. The vast majority of questions removed exhibited one or more four flaws: 1) Only required recognition of a common object, 2) only required counting a readily specified object, 3) did not require looking at the image to answer, 4) only asked about the color of a readily specified object. Examples of questions from each of these categories are shown in Fig. 4.

 <p><i>What are the round fruits in the bowl</i></p> <p>A) orange B) grapefruit C) apple D) peach</p>	 <p><i>How many pieces of luggage are in the snow</i></p> <p>A) two B) three C) four D) five</p>	 <p><i>When were modern traffic lights invented</i></p> <p>A) 1920 B) 1930 C) 1940 D) 1950</p>	 <p><i>What color is the door in this image</i></p> <p>A) tan B) gray C) burgundy D) white</p>
Reason rejected: Question only requires recognition of a common object.	Reason rejected: Question only requires counting specified objects.	Reason rejected: Question doesn't require looking at the image to answer.	Reason rejected: Question only asks about the color of a clearly specified object.

Figure 4: Examples of questions rejected for not meeting our criteria.

A.2 Data collection interface

The data-collection interface used by crowdworkers to write questions is shown in Fig. 5. Detailed instructions along with examples of good and bad questions were provided. After writing a question, workers were required to press the “Check for similar question” button. This sent a request to a server which returned the five questions closest to those already written in our growing dataset. We asked workers to rewrite or rephrase questions that were too similar, but did not enforce a minimum distance cutoff. The set of questions queried were reset when collecting the val and test sets to allow a greater degree of overlap with the training set. After satisfied with their question, workers advanced to the next image. Each task workers performed included four images, nearby neighbors in a CLIP embedding space, which encouraged creative differences in questions written for similar images. Workers were only required to write two questions (out of four possible images) to allow them to skip images they didn’t feel they could write a suitable questions for. This cut down on unsuitable questions that they would have otherwise been forced to write in order to complete the task. After completing two questions, workers were allowed to submit their work and advance to the next image set.

The data-collection interface used by crowdworkers to write rationales is shown in Fig. 6. Detailed instructions along with examples of good rationales were provided. We first asked workers to confirm the correct answer or provide the answer they thought was correct. This allowed a check on the correctness of the original question, and questions with a disagreement were removed from the dataset. Workers then provided a 1-2 sentence explanation of why the answer was correct that included any external knowledge needed to arrive there.

B Additional Details for Large-scale Pre-trained Models

We produce the vocabulary for the experiments in Sec. 5.2 from the training set by selecting all correct choices, as well as all choices and direct answers that appear in at least three questions. This results in a vocabulary with 10,424 answers.

Instructions

In this task you will be given four images and asked to write a question for each. You will also need to provide the correct answer along with at least 3 other answer options that are incorrect. Here are the guidelines to follow when writing:

- The questions must require looking at the image to answer.
- Questions should require some knowledge outside the image to answer (either from your common sense knowledge about the world, or something that you search google or wikipedia for).
- Answer options should be 1-2 words at most.
- Questions should be unique, try not to write questions that are too similar to others.
- To check previously written questions, you can press the first check question button to see similar questions others have written. If they're too similar, you can modify your questions and check again.
- After finishing a question, press next to load a new image and start back at step 1. After the writing a question for the 4th image, you can press the green submit button to finish the HIT.
- Make the correct answer option (A), the first of your answer choices, so we know which it is.
- Make sure to avoid questions that only require recognising an object, with no outside knowledge or thinking needed

Examples:






<p>What is this type of accident called?</p> <p>Answers:</p> <p>A) jackknifing B) head-on-collision C) rear-ending D) pileup</p> <p>Clues:</p> <p>the truck has folded in on itself there are no other cars involved</p>	<p>Which car is in the greatest danger?</p> <p>Answers:</p> <p>A) silver car B) red car C) white car D) black car</p> <p>Clues:</p> <p>the silver car is spanning the track there is a train coming</p>	<p>What fuels this type of train?</p> <p>Answers:</p> <p>A) diesel B) gasoline C) coal D) electricity</p> <p>Clues:</p> <p>this is a modern train there are no wires or smokestack</p>	<p>Why is the man wearing an orange vest?</p> <p>Answers:</p> <p>A) visibility B) fashion C) camouflage D) dress code</p> <p>Clues:</p> <p>he is probably working workers wear orange vests to be visible to others</p>
--	---	--	---

Things to avoid:






<p>What is this plane doing?</p> <p>Reason:</p> <p>It's hard to tell whether it's landing or taking off</p>	<p>What airline is this the logo of?</p> <p>Reason:</p> <p>This question is too specific.</p>	<p>What type of plane is this?</p> <p>Reason:</p> <p>There are too many possible answers.</p>	<p>What is the white and red object in the center of the image?</p> <p>Reason:</p> <p>This only requires recognizing the object.</p>
--	--	--	---

Here is your image:



Step 1 - Write the question Create a multiple choice question, with 4 answer options labeled (A) (B) (C) and (D), for example:

Which car is in the greatest danger? (A) silver car (B) red car (C) white car (D) black car (make sure the correct answer is always option A)

Where is this video game being played? (A) store display (B) electronics expo (C) at home (D) at manufacturer

The question is valid. Please proceed.

Step 2 - Make sure your question is unique Click the button: [Check for similar questions](#)

question	similarity score
Who are playing video game	0.7602785
What video game console are they playing	0.72090036
What video game console are they playing	0.72090036
What video game system is the person playing on	0.71226996
What video game system is the person playing on	0.71226996

Back
1 / 4
Next

Submit

Figure 5: Instructions and interface used for question collection.

B.1 Discriminative models

We train all of our discriminative train models for 500 epochs with a learning rate of 0.01 and batch size of 128, except the model with ResNet input features, which is trained with a learning rate of 0.001.

B.2 Contrastive models

The CLIP zero-shot setting requires no training. In the trained setting, we train our linear layer for 500 epochs with a learning rate of 0.01 and batch size of 128. We further elaborate on our “CLIP-style contrastive loss” below and visualize it in Fig. 7.

Recall that we have passed CLIP representations (for questions and/or images) through a linear layer to produce a 512-d embedding (the same size as a CLIP text encoding). For a batch of embeddings E and the CLIP text encodings of their corresponding answers A , we produce a cosine similarity matrix between E and A (i.e. the purple matrix in Fig. 7, showing a batch size of 4). We apply softmax over each matrix row (producing embedding-answer matching probabilities per embedding over answers in A) and compute a cross-entropy loss to maximize the similarity between each embedding and its corresponding answer.

B.3 Generative models

We show our modified ClipCap model in Fig. 8. As in ClipCap [37], we provide CLIP image representations to a mapping network, which produces prefix tokens as input for GPT-2. We then tokenize our question and ground-truth answer (appended with an end-of-sequence string, $\langle \text{EOS} \rangle$) and also provide these tokens as input. The remaining input tokens (in black) are zero-padding. As mentioned in our paper, we also appended the (pre-tokenized) question string with “Choices: ...” during the MC setting.

Answer the question and give a brief explanation.

This task gives you a question written about an image along with four answer options. We'd like you to give the correct answer and then explain how you arrive at it.

- First, tell us which answer you think is correct by writing the letter of the correct answer- a, b, c, or d on the first line.
- Skip a line, and then write one or two sentences that explain why that is the correct answer.
- These explanations should be simple and to the point.
- They will generally include pointing out some detail in the image and stating some fact about the world.
- If possible, try not to use exactly the same word or phrase from the correct answer in your explanation.
- Answer option A will be correct most of the time, but not always, so please look at all options before deciding.
- Please see examples below to get an idea of what's expected

Show / Hide Examples



Why is the person staying on top of the jet? (A) performing (B) seeking help (C) being trapped (D) sightseeing

A

The person is wearing a harness attached to the plane. The style of the plane is generally only seen at airshows.



What is the last time for parking on Thursdays in this street?

(A) 9:30 pm (B) 2:00 pm (C) 7:00 am (D) 8:00 am

A

The sign says there's parking Monday-Saturday that ends at this time. Thursday falls in the range



Where will the person put the skateboard? (A) Sidewalk (B) air (C) grass (D) street

A

Skateboards don't roll well on grass. They boy is too young too skate in the street.



How is this automated kiosk powered? (A) Solar energy (B) gas (C) coal (D) manual cranking

A

Parking kiosks run on electricity. This unit has a solar panel on the top.



What is he focused at?

Answer options:

- (A) television
- (B) another person
- (C) street
- (D) window

Please give the correct answer and explain your choice:

A

He is holding a game controller which requires looking at a tv to play.

[Submit](#)

Figure 6: Instructions and interface used for rationale collection.

This model is trained autoregressively. I.e., O_i is generated conditionally, given $I_0 \cdots I_i$ (for input tokens I and output logits O), and supervised with a cross-entropy loss against the next sequence token I_{i+1} . In our case, we only compute this cross-entropy loss for outputs corresponding with the ground-truth answer tokens (including $\langle \text{EOS} \rangle$).

At inference time, we prompt GPT-2 with our image prefix and question tokens. We have the model predict the most likely next token (i.e. generating a token in the answer) from the output logits. We append this token to the input and repeat this step, until the model predicts $\langle \text{EOS} \rangle$. We can use the tokenizer to decode these output tokens (excluding $\langle \text{EOS} \rangle$), producing our model’s textual answer prediction. Note that beam search is an alternative way to generate text from autoregressive language models, but we found that it led to worse results, likely because the answers we are trying to generate are short (e.g. 1-3 words).

We fine-tuned the models in our experiments (choosing the checkpoint with the best F1 validation score for generated answers over 10 epochs), using the settings and COCO pre-trained weights (for the MLP mapping network) made available by the ClipCap authors⁷. For the pre-trained MLP model, they used CLIP ViT-B/32 features, produced 10 image prefix tokens, and had also fine-tuned GPT-2 (for their image captioning task). We further fine-tuned the GPT-2 weights on our task.

C Additional Details for Rationale Generation

We generated rationales from ClipCap in a nearly identical manner to how we generated answers (see Sec. B.3 and Fig. 8 above). However, we replace the ground-truth answer string/tokens with a ground-truth rationale. And, we don’t provide “Choices: ...” in the ClipCap prompt for the MC setting. We also use beam search during generation, as it seems to perform better for these longer

⁷https://github.com/rmokady/CLIP_prefix_caption

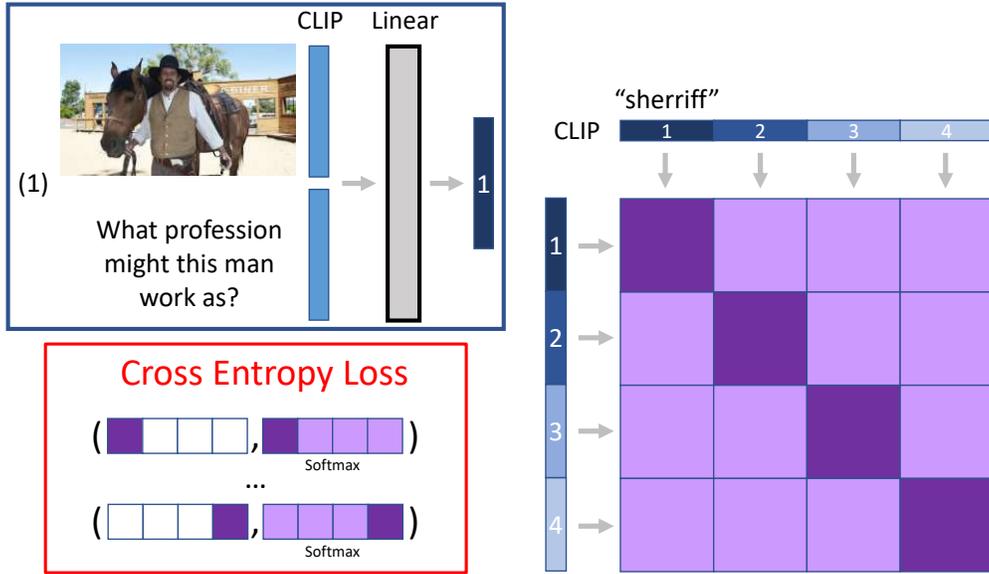


Figure 7: As described in Sec. B.2. CLIP-style contrastive loss between embeddings (of questions and images) and CLIP text encodings (of answers). Shown for a batch size of 4.

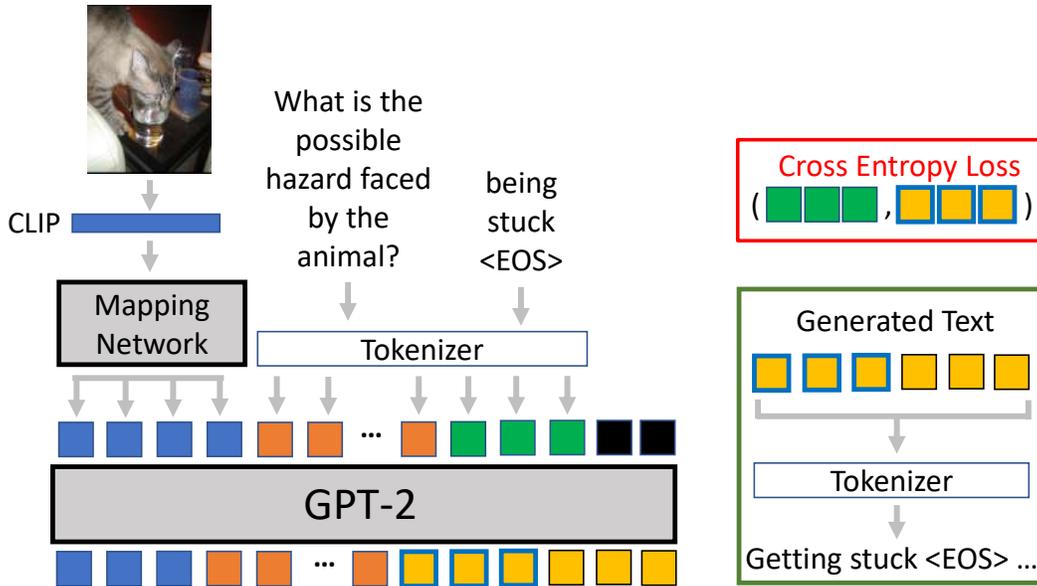


Figure 8: Diagram of modified ClipCap architecture for VQA tasks.

strings. We also use the MLP mapping network and continue to fine-tune GPT-2, as it demonstrates the best performance for this task. We again fine-tuned this model on our training data for 10 epochs and picked the checkpoints with best BLEU and METEOR validation scores.

We show some examples of generated rationales in Fig. 9.

D Additional Details for Specialized Models

For all of these models, we use the same training hyperparameters as the original implementation. For all of the discriminative methods in the paper we use a fixed vocabulary constructed from direct



Figure 9: Examples of rationales generated by our modified ClipCap method for examples in our validation set.

answers that appeared two or more times in the training set. This includes 2,133 bi-grams or unigrams, with 1,937 words.

Pythia [24] Pythia is a modification of [1] that introduces changes to the architecture and learning schedule and utilizes more training data. We fine-tune it on the A-OKVQA dataset. For fine-tuning, we replace the top classification layer with a randomly initialized layer for our set of answer vocabulary.

LXMERT [54] LXMERT is a Transformer-based vision and language model pre-trained using a large amount of image-sentence pairs for a set of pre-training tasks such as masked language modeling and object prediction. The model is pre-trained on VQAv2 [16], GQA [18], VG-QA [69], COCO captions [11], and Visual Genome captions [27]. We then fine-tune the model using the training set of A-OKVQA.

ViLBERT [32] ViLBERT is an extension of the BERT architecture to process vision and language modalities for learning a joint representation for them. ViLBERT has been pre-trained on proxy tasks, but it has been evaluated on VQA as a downstream task. ViLBERT is pre-trained using Conceptual Captions [49] and fine-tuned on A-OKVQA. To evaluate how well a model trained on VQAv2 or OK-VQA performs on A-OKVQA, we fine-tune ViLBERT after training them on those datasets. These models are referred to as ‘ViLBERT-VQA’ and ‘ViLBERT-OK-VQA’ in Table 5.

KRISP [35] KRISP is a method for knowledge-based VQA which combines multi-modal Transformers with graph neural networks methods on knowledge graphs. We use the same models and data and knowledge sources and pre-processing steps as in that work, but filter the knowledge graph based on A-OKVQA rather than OK-VQA (see Sec. 3.2 of [35]).

GPV-2 [26] GPV-2 [26] is a generative vision and language model built using the T5 [44] language model and VinVL [67] image features. It was pre-trained on Conceptual Captions [49] and then fine-tuned in a multi-task setting on image captioning, visual question answering, object localization, and classification, as well on web-search images for 10,000 visual concepts.

We fine-tune the fully-trained model on A-OKVQA by training it to generate the most common answer for each question. For direct answer evaluations, answers are then generated using beam search with 20 beams. For multiple choice, the answers are ranked by the log-probability score assigned to them by the model.

We perform two additional experiments with rationales with this model. First, ground-truth rationales are appended to the question as additional input text. Recall that we do not provide rationales at test

time. However, for this experiment we use them during test. We refer to this model as ‘GPV-2 + GT Ratl.’. Second, we use the same setting, but we replace every occurrence of the ground-truth answer in the rationale with the [answer] token. We refer to this model as ‘GPV-2 [26] + Masked Ans.’ in Table 5.