# Which visual questions are difficult to answer? Analysis with Entropy of Answer Distributions

Kento Terao, Toru Tamaki, Bisser Raytchev, Kazufumi Kaneda Hiroshima University, Japan

Shun'ichi Satoh National Institute of Informatics, Japan

## **Abstract**

We propose a novel approach to identify the difficulty of visual questions for Visual Question Answering (VQA) without direct supervision or annotations to the difficulty. Prior works have considered the diversity of ground-truth answers of human annotators. In contrast, we analyze the difficulty of visual questions based on the behavior of multiple different VQA models. We propose to cluster the entropy values of the predicted answer distributions obtained by three different models: a baseline method that takes as input images and questions, and two variants that take as input images only and questions only. We use a simple kmeans to cluster the visual questions of the VQA v2 validation set. Then we use state-of-the-art methods to determine the accuracy and the entropy of the answer distributions for each cluster. A benefit of the proposed method is that no annotation of the difficulty is required, because the accuracy of each cluster reflects the difficulty of visual questions that belong to it. Our approach can identify clusters of difficult visual questions that are not answered correctly by state-of-the-art methods. Detailed analysis on the VQA v2 dataset reveals that 1) all methods show poor performances on the most difficult cluster (about 10% accuracy), 2) as the cluster difficulty increases, the answers predicted by the different methods begin to differ, and 3) the values of cluster entropy are highly correlated with the cluster accuracy. We show that our approach has the advantage of being able to assess the difficulty of visual questions without ground-truth (i.e., the test set of VQA v2) by assigning them to one of the clusters. We expect that this can stimulate the development of novel directions of research and new algorithms.

Clustering results are available online<sup>1</sup>, in which we show lists of pairs of questions and clusters for both of the validation and test sets of the VQA v.2 dataset.

## 1. Introduction

Visual Question Answering (VQA) is one of the most challenging tasks in computer vision [40, 3]: given a pair of question text and image (a visual question), a system is asked to answer the question. It has been attracting a lot of attention in recent years because it has a large potential to impact many applications such as smart support for the visually impaired [15], providing instructions to autonomous robots [8], and for intelligent interaction between humans and machines [9]. Towards these goals, many methods and datasets have been proposed.

The VQA task is particularly challenging due to the diversity of annotations. Unlike common tasks, such as classification, where precise ground truth labels are provided by the annotators, a visual question may have multiple different answers annotated by different crowd workers, as shown in Figure 1. In VQA v2 [12] and VizWiz [6], which are commonly used in this task, each visual question was annotated by 10 crowd workers, and almost half of the visual questions in these datasets have multiple answers [14, 5], as shown in Table 1 for VQA v2. The metric for performance evaluation commonly used for these dataset has therefore the following form [3]:

$$accuracy = \min\left(\frac{\text{\# humans that provided that answer}}{3}, 1\right),$$

in other words, an answer is correct in 100% if at least three annotated answers match that answer.

The disagreement of crowd workers in ground truth annotations has been an annoying issue for researchers dealing with tasks which involve crowdsource annotations [7, 35, 28]. Recently some works on VQA have tackled this issue. Gurari *et al.* [14] analyzed the number of unique answers annotated by crowd workers and proposed a model that predicts when crowdsourcing answers (dis)agree by using binary classifiers. Bhattacharya *et al.* [5] categorized reasons why answers of crowd workers differ, and found which co-occurring reasons arise frequently.

These works have revealed why multiple answers may

<sup>1</sup>https://github.com/tttamaki/vqd



Q: How many of the men are wearing glasses? A: {'2'x5, '3'x4, '1'x1} Entropy: 0.943348392



Q: What color is the hat? A: {'red'x9, 'orange'x1} Entropy: 0.325082973



Q: How many people can fit in the 2 buses? A: {100'x4, 'many'x2, '80'x1, '40'x1, 'lot'x1, '200'x1} Entropy: 1.609437912



Q: How many birds? A: {'2'x10} Entropy: 0

Figure 1. Examples of visual questions and corresponding 10 answers of VQA v2 datasets, and corresponding entropy values.

Table 1. Numbers of unique answers per visual question of the validation set of VQA v2. The bottom row shows averages of unique answers.

Weis.								
	#Ans	Yes/No	Number	Other	All			
	1	41561	9775	18892	70228			
	2	33164	6701	18505	58370			
	3	5069	3754	15238	24061			
	4	621	2110	12509	15240			
	5	103	1528	10661	12292			
	6	23	1239	9186	10448			
	7	0	1062	7666	8728			
	8	0	952	6169	7121			
	9	0	726	4528	5254			
	10	0	287	2325	2612			
	total	80541	28134	105679	214354			
	ave	1.57±0.46	2.93±1.59	$4.04 \pm 1.75$	$2.97{\pm}1.60$			

arise and when they disagree, however this is not enough to find out how multiple answers make the visual question difficult for VQA models. Malinowski et al. [25] reported that the disagreement harms the performance of the VQA model, therefore the diversity of answers should be an important clue. However, formulating the (dis)agreement as binary (single or multiple answers) drops the information of the extent how diverse multiple answers are. For example, suppose two different answers are given to a visual question. This may mean that "five people gave one answer and the other five gave the other answer," or, that "one gave one answer and the rest 9 gave the other." In the latter case, the answer given by the first annotator may be noisy, hence not suitable for taking into account. To remove such noisy answers, prior work [14, 5] employed a minimum number of answer agreement. If the agreement threshold is set to m=2 (at least two annotators are needed for each answer to be valid), then the answer given by the single annotator is ignored. However setting a threshold is ad-hoc and different threshold may lead to different results when other datasets annotated by more (other than 10) workers would be available.

In this paper, we propose to use the entropy of answer distribution, instead of answer (dis)agreement. Let A is the set of answers, and the entropy H(A) is defined by

$$H(A) = -\sum_{a \in A} P(a) \ln P(a). \tag{1}$$

In general, entropy is large when the distribution is broad, and small when it has a narrow peak. This is a simple but useful indicator of the diversity of answers in ground truth annotations. Yang *et al.* [41] used the entropy as a metric of diversity for the task of predicting the answer distribution of ground truth annotations. This is beneficial for investigating how diverse human annotations are, and evaluating how difficult visual questions are for humans.

In contrast, we use the entropy values of answer predictions produced by different VQA models to evaluate the difficulty of visual questions for the models. Entropy values are available at no additional cost because it is common to predict an answer distribution by using softmax for computing the cross entropy loss. To the best of our knowledge, this is the first work to use entropy for analysing the difficulty of visual questions.

The use of the entropy of answer distribution enables us to analyse visual questions in a novel aspect. Prior works have reported overall performance as well as performances on three subsets of VQA v2 [12]; Yes/No (answers are yes or no for questions such as "Is it ..." and "Does she ..."), Numbers (answers are counts, numbers, or numeric, "How many ..."), and Others (other answers, "What is ..."). These three types have different difficulties (i.e., Yes/No type is easier, Other type is harder), and performances of each type are useful to highlight how models behave to different types of visual questions. In fact, usually the first two words carry the information of the entire question [14], and previous work [1] uses this fact to switch the internal model to adopt suitable components to each type. This categorization of question types is useful, however not enough to find which visual questions are difficult. If we can evaluate the difficulty of visual questions, this could push forward the development of better VQA models.

Our goal is to present a novel way of analysing visual questions by clustering the entropy values obtained from different models. Images and questions convey different information [13, 4], hence models that take images only or question only are often used as baselines [3, 5, 12]. Datasets often have the language bias [12], and then questions only may be enough to answer reasonably. However the use of the image information should help to answer correctly. Our key idea is that the entropy values of three models (that use image only (I), question only (Q), and both (Q+I)) are useful to characterize each visual question.

The contributions of this work can be summarized as follows.

- Instead of using the entropy of ground truth annotations, we use the entropy of the predicted answer distribution for the first time to analyse how diverse predicted answers are. We show that entropy values of different models are useful to characterize visual questions.
- We propose an entropy clustering approach to categorize the difficulty levels of visual questions. After training three different models (I, Q, and Q+I), predicting answer distributions and computing entropy values, the visual questions are clustered. This is simple yet useful, and enables us to find which visual questions are most difficult to answer.
- We discuss the performances of several state-of-the-art methods. Our key insight is that the difficulty of visual question clusters are common to all methods, and tackling the difficult clusters may lead to the development of a next generation of VQA methods.

## 2. Related work

The task of VQA has attracted a lot of attention in recent years. Challenges have been conducted since 2016, and many datasets have been proposed. In addition to the normal VQA task, related tasks have emerged, such as EmbodiedQA [8], TextVQA [33], and VQA requiring external knowledge [38, 26, 34]. Still the basic framework of VQA is active and challenging, and some tasks include VQA as an important component, such as visual question generation [27, 22], visual dialog [9, 16], and image captions [30].

VQA datasets have two types of answers. For multiple-choice [12, 45, 42], several candidate answers are shown to annotators for each question. For open-ended [3, 12, 6, 18, 29], annotators are asked to answer in free text, hence answers tend to differ for many reasons [5]. Currently two major datasets, VQA [3, 12] and VizWiz [6], suffer from this issue because visual questions in these datasets were answered by 10 crowd workers, while other datasets [29, 45, 19, 21, 38, 18, 42, 11] have one answer per visual question.

This disagreement between annotators has recently been investigated in several works. Bhattacharya *et al.* [5] proposed 9 reasons why and when answers differ: low-quality image (LQI), answer not present (IVE), invalid (INV), difficult (DFF), ambiguous (AMB), subjective (SBJ), synonyms (SYN), granular (GRN), and spam (SMP). The first six reasons come from both/either question and/or image, and the last three reasons are due to issues inherent to answers. They found that ambiguity occurs the most, and co-occurs with synonyms (same but different wordings) and granular (same but different concept levels). This work gives us quite an important insight about visual questions, however only for those that have multiple different answers annotated. Gurari *et al.* [14] investigated the number of unique

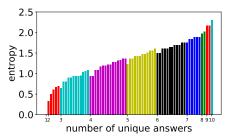


Figure 2. Entropy values of all possible combinations of unique number of answers.

answers annotated by crowd workers, but didn't consider how answers differ if disagreed. Instead they use a threshold of agreement to show how many annotators answered the same.

Our approach is to use the entropy of answer distributions of both ground truth and prediction. This is a novel aspect, and complementary to the prior works [14, 5]. Entropy takes into account by a single number the fraction of multiple answers as well as the distribution of answers. It therefore provides another modality to analyse visual questions at a fine-grained level. Figure 2 shows how entropy values change for the same number of unique answers. The leftmost bar's value is zero because there is only a single answer (*i.e.* all answers agree), and the rightmost bar represents the case when all 10 answers are different. In between, entropy values are sorted inside the same number of unique answers. This shows that entropy is finer than the number of unique answers.

We should note that this approach is different from uncertainty of prediction. Teney *et al.* [36] proposed a model using soft scores because scores may indicate uncertainty in ground truth annotations, and minimizing the loss between ground truth and prediction answer distribution. This approach is useful, yet it doesn't show the nature of visual questions.

Our approach is closely related to hard example mining [39, 31] and hardness / failure prediction [37]. Hard example mining approaches determine which examples are difficult to train during training, while hardness prediction jointly trains the task classifier and an auxiliary hardness prediction network. Compared to these works, our approach differs in the following two aspects. First, the VQA task is multi-modal and assessing the difficulty of visual questions has not been considered before. Second, our approach is off-line and can determine the difficulty without ground-truth, *i.e.*, before actually trying to answer the visual questions in the test set.

# 3. Clustering visual questions with entropy

# 3.1. Clustering method

To perform clustering, we hypothesize that "easy visual questions lead to low entropy while difficult visual questions to high entropy." This has been reported for the entropy of ground truth annotations by Malinowski *et al.* [25]. Here we extend this concept to the entropy of answer distributions produced by VQA models. This is reasonable because for easy visual questions VQA systems can predict answer distributions in which the correct answer category has large probability while other categories are low. In contrast, difficult visual questions makes VQA systems generate broad answer distributions because many answer candidates may be equally plausible. Entropy can capture the diversity of predicted answer distributions, and also that of ground truth annotations in the same manner.

We prepare three different models that use as input image only (I), question only (Q), and both question and image (Q+I). In this case, we expect the following three levels of difficulty of visual questions:

- Level 1: Reasonably answered by using question only.
- Level 2: Difficult to answer with question only but good with images.
- Level 3: Difficult even if both image and question are provided.

For a certain visual question, it is of level 1 if the answer distribution of the Q model has low entropy. It is of level 2 if the Q model is high entropy and the Q+I model is low entropy. If both the Q and Q+I models have high entropy, then the visual question is of level 3. This concept is realised by the following procedure. 1) Train the I, Q, and Q+I models on the training set with image only, questions only, and both images and questions, respectively. 2) Evaluate the validation set by using the three models and compute answer distributions and entropy values of each of visual questions. 3) Perform clustering on the validation set with entropy values. Clustering features are the entropy values of the three models.

## 3.2. Datasets and setting

We use VQA v2 [12]. It consists of training, validation, and test sets. To train models, we use the training set (82,783 images, 443,757 questions, and 4,437,570 answers). We use the validation set (40,504 images, 214,354 questions, and 2,143,540 answers) for clustering and analysis.

We choose Pythia v0.1 [17, 32] as a base model, and modify it so that it takes questions only (Q model), or images only (I model). To do so, we simply set either image

Table 2. Accuracy of models on the validation set of VQA v2.

Model	Overall	Yes/No	Number	Other	
I	24.65±41.42	64.21±44.51	$0.27\pm3.16$	0.99±7.48	
Q	44.83±46.76	$68.48{\pm}43.00$	$32.05{\pm}43.45$	$30.21 \pm 42.91$	
Q+I	67.47±43.35	$84.52 {\pm} 33.01$	$47.55 \!\pm\! 46.25$	<b>59.78</b> ±45.00	
BUTD	63.79±44.61	81.20±35.84	$43.90 \pm 45.93$	55.81±45.78	
MFB	65.14±44.21	$83.11 \pm 34.31$	$45.32 {\pm} 46.16$	$56.72 \pm 45.60$	
MFH	66.23±43.85	$84.12 \pm 33.45$	$46.71 \!\pm\! 46.27$	$57.79 \pm 45.40$	
BAN-4	65.87±43.90	$83.57 \pm 33.88$	$47.23 \pm 46.17$	$57.34 \pm 45.43$	
BAN-8	66.00±43.87	$83.48 \pm 33.95$	$47.20 \!\pm\! 46.17$	$57.69 \pm 45.40$	
MCAN-small	67.20±43.42	$84.91 \pm 32.68$	$49.35 \!\pm\! 46.23$	$58.46 \pm 45.18$	
MCAN-large	67.47±43.33	$85.33 \pm 32.24$	$48.96{\pm}46.23$	$58.78 \pm 45.13$	
Pythia v0.3	65.91±44.42	$84.30 \pm 33.56$	$44.90{\pm}46.47$	$57.49 \pm 46.07$	

Table 3. Entropy of models on the validation set of VQA v2.

Model	Overall	Yes/No	Number	Other	
I	4.19±0.42	$4.16\pm0.42$	$4.19\pm0.43$	$4.21\pm0.41$	
Q	1.80±1.32	$0.59 {\pm} 0.22$	$2.28{\pm}0.94$	$2.60 \pm 1.22$	
Q+I	$0.84\pm1.06$	$0.20 \pm 0.27$	$1.39 \pm 1.18$	$1.19 \pm 1.15$	
BUTD	1.24±1.33	$0.32 \pm 0.29$	$1.86\pm1.25$	1.77±1.45	
MFB	1.76±1.86	$0.42 {\pm} 0.31$	$2.07 \pm 1.77$	$2.71\pm1.95$	
MFH	1.63±1.77	$0.40 {\pm} 0.31$	$2.00 \pm 1.76$	$2.46{\pm}1.89$	
BAN-4	$0.99\pm1.20$	$0.21 \pm 0.27$	$1.60 \pm 1.25$	$1.43 \pm 1.31$	
BAN-8	$0.95\pm1.17$	$0.20 {\pm} 0.26$	$1.53 \pm 1.23$	$1.36 \pm 1.27$	
MCAN-small	1.21±1.71	$0.17 \pm 0.27$	$1.66 \pm 1.82$	$1.89 \pm 1.91$	
MCAN-large	1.15±1.64	$0.16 {\pm} 0.26$	$1.63 \pm 1.76$	$1.78 \pm 1.84$	
Pythia v0.3	$0.59\pm0.82$	$0.13 \pm 0.22$	$0.86{\pm}0.93$	$0.89 {\pm} 0.87$	
GT	0.67±0.68	0.25±0.29	$0.66 \pm 0.68$	$0.99\pm0.71$	

features or question features to zero vectors. With no modification, it is Q+I model (*i.e.* Pythia v0.1). As in prior works [2, 23, 36, 43], 3129 answers<sup>2</sup> in the training set that occur at least 8 times are chosen as candidates, which results in a multi-class problem predicting answer distributions of 3129 dimension.

To compare the performance with state-of-the-art methods, we use BUTD [2], MFH [44], BAN [20] (including BAN-4 and BAN-8), MCAN [43] (including small and large), and Pythia v0.3 [32, 33].

First we show the performance of each model in Table 2. As expected, the I model performs worst because there is no clue of questions in the image. In contrast, Q model performs reasonably better, particularly for Yes/No type. Average performances of models (excluding I and Q) are about 84%, 47%, and 58% for types of Yes/No, Number, and Other, respectively.

Next in Table 3 we show the entropy values of the predicted answer distributions by different models for each of the three types, as well as ground truth annotations.<sup>3</sup> Average entropy values of models (excluding I and Q) for each type are 0.25, 1.62, 1.72, respectively. Yes/No type has smaller entropy than the others because answer distributions tend to gather around only two candidates ("Yes" and "No").

<sup>&</sup>lt;sup>2</sup>Other common choices are 3000 [44, 10] and 1000 [24]. Even when different numbers are used, our entropy clustering approach works and we expect our findings to hold.

 $<sup>^3</sup>$ Entropy ranges from 0 (single answer) to 2.303 (10 different answers) for ground truth answers, and from 0 (1 for a single entry, otherwise 0) to 8.048 (uniform values of 1/3129) for model predictions.

Table 4. Clustering results for the validation set of VQA v2. Each column corresponds to a different cluster and colors indicate cluster types (level 1 in gray, level 2 in yellow, and level 3 in red).

cluster 0 1 2 3 4 5 6 7 8 9								0			
	cluster	0	1 47 + 0.24			4	5	6	7	8	9
base model entropy	1	3.77±0.27	4.47±0.24	4.22±0.39	4.09±0.40	4.22±0.40	4.19±0.40	4.23±0.40	4.27±0.39	4.23±0.41	4.42±0.41
ba ntr	Q	0.60±0.24	0.61±0.23	2.69±0.32	1.78±0.29	4.09±0.47	1.48±0.41	2.61±0.36	4.01±0.49	2.68±0.49	4.33±0.62
- o	Q+I	0.20±0.26	0.21±0.27	0.24±0.27	0.25±0.28	0.63±0.46	1.45±0.45	1.46±0.35	2.25±0.43	2.77±0.47	3.79±0.56
	BUTD	0.38±0.45	0.42±0.48	0.85±0.98	0.77±0.83	1.73±1.35	1.79±0.92	2.03±0.97	3.13±1.06	3.01±0.90	$3.98\pm0.91$
Ħ	MFB	0.55±0.67	$0.57\pm0.68$	1.41±1.41	$1.15\pm1.16$	$2.77 \pm 1.80$	$2.39\pm1.43$	$2.85{\pm}1.48$	$4.44\pm1.36$	$4.03\pm1.36$	$5.47 \pm 1.27$
state-of-the-art entropy	MFH	$0.49\pm0.56$	$0.51 \pm 0.57$	1.23±1.27	$1.00\pm1.03$	$2.47 \pm 1.75$	$2.19\pm1.30$	$2.60 \pm 1.37$	$4.21\pm1.42$	$3.83\pm1.33$	$5.37 \pm 1.33$
e-of-the- entropy	BAN-4	$0.25\pm0.37$	$0.27\pm0.39$	$0.61 \pm 0.81$	$0.53 \pm 0.67$	$1.32 \pm 1.21$	$1.46 \pm 0.83$	$1.68\pm0.89$	$2.67 \pm 1.07$	$2.63\pm0.89$	$3.61\pm0.99$
ent	BAN-8	$0.23\pm0.36$	$0.26\pm0.37$	$0.57\pm0.77$	$0.50\pm0.64$	$1.21 \pm 1.16$	$1.40 \pm 0.82$	$1.60 \pm 0.87$	$2.55\pm1.07$	$2.54\pm0.90$	$3.48\pm1.01$
ţ	MCAN-small	0.23±0.45	$0.25\pm0.49$	$0.70\pm1.16$	$0.54\pm0.89$	$1.74 \pm 1.78$	$1.66 \pm 1.3$	$2.0\pm1.45$	$3.66 \pm 1.67$	$3.35{\pm}1.46$	$4.95 \pm 1.52$
×.	MCAN-large	$0.21\pm0.42$	$0.23\pm0.45$	$0.64\pm1.07$	$0.51\pm0.84$	$1.62\pm1.71$	$1.59\pm1.23$	$1.9\pm1.37$	$3.5\pm1.65$	$3.21\pm1.43$	$4.82\pm1.53$
	Pythia v0.3	$0.14\pm0.28$	$0.16\pm0.29$	$0.28\pm0.50$	$0.25\pm0.44$	$0.69 \pm 0.78$	$0.89 \pm 0.68$	$0.99 \pm 0.69$	$1.59 \pm 0.85$	$1.76\pm0.80$	$2.20\pm0.90$
test set entropy	I	$3.77\pm0.26$	$4.47\pm0.24$	4.21±0.39	$4.08\pm0.40$	$4.22\pm0.41$	$4.19\pm0.41$	$4.23\pm0.40$	$4.28\pm0.40$	$4.24\pm0.41$	4.43±0.41
test set entropy	Q	$0.60\pm0.23$	$0.61\pm0.23$	$2.70\pm0.32$	$1.78\pm0.29$	$4.09\pm0.47$	$1.48\pm0.41$	$2.62\pm0.37$	$4.02\pm0.49$	$2.68\pm0.50$	$4.33\pm0.61$
5 E	Q+I	$0.18\pm0.25$	$0.20 \pm 0.26$	$0.24\pm0.27$	$0.26\pm0.29$	$0.63 \pm 0.45$	$1.45\pm0.45$	$1.46 \pm 0.35$	$2.25 \pm 0.44$	$2.78\pm0.49$	$3.78\pm0.56$
. <del>.</del> .	I	53.13±47.09	54.81±46.93	0.67±6.88	1.61±11.58	0.75±6.89	2.33±13.20	$0.69 \pm 5.87$	1.05±7.18	$0.90\pm6.28$	1.13±7.00
base model acc.	Q	69.70±42.54	$66.73 \pm 43.73$	32.91±45.13	$46.91\pm47.63$	$16.32\pm35.09$	$34.99 \pm 42.62$	$24.13\pm37.8$	$8.71\pm23.97$	$14.54\pm29.74$	$5.50\pm18.76$
7 8 "	Q+I	86.15±31.32	$82.87 \pm 34.48$	84.53±32.9	83.08±34.53	<b>67.79</b> ±42.20	$47.10\pm44.24$	$45.73 \pm 43.64$	26.04±38.10	22.53±34.56	$9.32\pm25.05$
	BUTD	82.55±34.81	$78.85\pm37.62$	$78.32 \pm 38.36$	$77.39\pm39.08$	60.18±44.87	45.77±44.21	43.06±43.62	23.53±36.76	22.6±35.03	9.93±25.37
E	MFB	84.20±33.29	$80.84 \pm 36.24$	79.77±37.26	$78.73 \pm 38.16$	$60.93 \pm 44.76$	$47.09\pm44.38$	$43.56\pm43.56$	$24.25\pm37.11$	$23.16\pm35.32$	$10.39\pm25.85$
- s-	MFH	85.38±32.21	$81.89 \pm 35.41$	80.72±36.51	$80.06\pm37.09$	$62.84 \pm 44.23$	$47.97 \pm 44.36$	$44.82 \pm 43.69$	25.17±37.79	23.77±35.47	$10.97 \pm 26.82$
rac rt	BAN-4	84.89±32.65	81.35±35.75	80.11±36.95	$79.70\pm37.32$	$61.98 \pm 44.48$	$47.97 \pm 44.20$	$45.27 \pm 43.72$	24.67±37.24	23.84±35.47	$10.76\pm26.41$
te-of-the- accuracy	BAN-8	84.88±32.64	$81.25 \pm 35.85$	80.51±36.67	$79.67 \pm 37.37$	$62.79 \pm 44.23$	$48.33 \pm 44.24$	$45.65 \pm 43.74$	$25.07 \pm 37.62$	$24.00 \pm 35.59$	$10.62\pm26.21$
state-of-the-art accuracy	MCAN-small	86.06±31.49	$82.76 \pm 34.63$	81.31±35.97	$80.78 \pm 36.50$	$63.52 \pm 43.96$	<b>49.86</b> ±44.24	$46.67 \pm 43.72$	$26.10 \pm 38.02$	$25.46 \pm 36.33$	$11.56\pm27.42$
20	MCAN-large	86.44±31.10	83.18±34.25	81.58±35.79	$80.86 \pm 36.53$	$63.79 \pm 43.88$	$49.33 \pm 44.30$	<b>47.27</b> ±43.78	$26.35 \pm 38.05$	$25.50\pm36.27$	$11.58\pm27.30$
	Pythia v0.3	85.56±32.46	$82.29 \pm 35.46$	82.35±36.42	$80.19\pm37.90$	$65.21 \pm 45.60$	$48.54 \pm 46.50$	$47.13\pm46.57$	<b>27.31</b> ±42.05	<b>26.32</b> ±40.98	<b>11.70</b> ±30.28
	entropy	0.30±0.36	$0.29\pm0.36$	$0.60\pm0.60$	$0.50\pm0.54$	$0.99 \pm 0.66$	$0.97 \pm 0.65$	$1.13\pm0.65$	$1.37\pm0.67$	$1.45\pm0.65$	$1.34\pm0.70$
	ave # ans	1.72±0.98	$1.69\pm0.96$	$2.71\pm1.91$	$2.39 \pm 1.67$	$3.98{\pm}2.30$	$3.84{\pm}2.26$	$4.43\pm2.37$	$5.42\pm2.57$	$5.75\pm2.53$	$5.39 \pm 2.68$
GT statistics	total	42637	52600	20235	21643	10631	13516	19010	12608	12620	8854
atis	yes/no	35483	44426	22	262	6	288	23	13	13	5
St	number	1194	1471	2770	5778	253	4489	4969	1255	3790	2165
5	other	5960	6703	17443	15603	10372	8739	14018	11340	8817	6684
	# agree	20762	26338	6770	8528	1488	1912	1988	954	699	789
	# disagree	21875	26262	13465	13115	9143	11604	17022	11654	11921	8065
:	LQI	0.05±0.04	0.05±0.04	$0.02\pm0.03$	0.03±0.04	$0.01\pm0.03$	0.03±0.04	0.03±0.04	$0.02\pm0.04$	$0.03\pm0.04$	$0.06\pm0.08$
_	IVE	$0.48\pm0.21$	$0.48\pm0.21$	$0.11\pm0.14$	$0.18\pm0.19$	$0.10\pm0.10$	$0.23 \pm 0.20$	$0.19\pm0.19$	$0.15\pm0.15$	$0.24\pm0.20$	$0.21\pm0.16$
Щe	INV	$0.26\pm0.14$	$0.25\pm0.14$	$0.02\pm0.03$	$0.03\pm0.05$	$0.01\pm0.02$	$0.04\pm0.06$	$0.03\pm0.04$	$0.02\pm0.02$	$0.03\pm0.04$	$0.03\pm0.03$
÷	DFF	$0.09\pm0.07$	$0.08\pm0.07$	$0.08\pm0.10$	$0.11\pm0.11$	$0.07\pm0.07$	$0.15\pm0.12$	$0.13\pm0.12$	$0.10\pm0.09$	$0.17\pm0.14$	$0.12\pm0.08$
s tc	AMB	0.75±0.13	$0.76\pm0.12$	$0.95 \pm 0.05$	$0.93 \pm 0.07$	$0.96\pm0.04$	$0.91 \pm 0.08$	$0.93 \pm 0.06$	$0.94 \pm 0.06$	$0.92 \pm 0.07$	$0.91 \pm 0.08$
uo.	SBJ	$0.32\pm0.23$	$0.30\pm0.23$	$0.13\pm0.09$	$0.14\pm0.11$	$0.12 \pm 0.08$	$0.13\pm0.12$	$0.12\pm0.09$	$0.11\pm0.09$	$0.11\pm0.09$	$0.10\pm0.09$
reasons to differ	SYN	$0.25\pm0.27$	$0.25 \pm 0.26$	$0.81 \pm 0.18$	$0.72\pm0.25$	$0.87 \pm 0.12$	$0.65 \pm 0.27$	$0.72\pm0.24$	$0.80 \pm 0.18$	$0.68 \pm 0.24$	$0.72\pm0.21$
-	GRN	$0.35\pm0.22$	$0.35 \pm 0.22$	$0.79\pm0.15$	$0.71\pm0.20$	$0.82 \pm 0.11$	$0.66 \pm 0.21$	$0.71\pm0.19$	$0.76\pm0.16$	$0.67\pm0.19$	$0.69\pm0.17$
	SPM	$0.03\pm0.02$	$0.02 \pm 0.01$	$0.01\pm0.01$	$0.01 \pm 0.01$	$0.01 \pm 0.01$	$0.01 \pm 0.01$	$0.01 \pm 0.01$	$0.01 \pm 0.01$	$0.01 \pm 0.01$	$0.02 \pm 0.02$
	OTH	$0.01\pm0.01$	$0.01 \pm 0.01$	$0.00 \pm 0.00$	$0.00 \pm 0.00$	$0.00 \pm 0.00$	$0.00 \pm 0.00$	$0.00 \pm 0.00$	$0.00 \pm 0.00$	$0.00 \pm 0.00$	$0.01 \pm 0.01$

## 3.3. Clustering results

Now we show the clustering results in Table 4. We used k-means to cluster the 3-d vectors of 214,354 visual questions into k = 10 clusters.<sup>4</sup>

Each column of Table 4 shows the statistics for each cluster. Clusters are numbered in ascending order of the entropy for the Q+I model. The top rows with 'base model entropy' show the entropy values for the three base models.

To find three levels of visual questions, we divide the clusters by the following simple rule. For each cluster, if 'Q entropy' < 1 then it is level 1, else if 'Q+I entropy' > 2 then it is level 3, otherwise level 2. Column colors of Table 4 indicate levels; level 1 (clusters 0 and 1) are in gray, level 2 in yellow (2 to 6), and level 3 (7, 8, and 9) in red.

Below we describe other rows of Table 4.

**base model acc.** Accuracy values of the three base models. Accuracy of Q+I model tends to decrease as Q+I entropy increases, which we will discuss later.

**state-of-the-art entropy and accuracy** Entropy and accuracy values of 9 state-of-the-art methods.

**test set entropy** Entropy values of the test set of VQA v2. We assign test visual questions to one of these clusters (we will discuss this later).

**GT statistics** Statistics of ground truth annotations. Row 'entropy' shows entropy values of ground truth annotations. Row 'ave # ans' shows the average number of unique answers per visual question. These two rows show how ground truth answers differ in each cluster.

Row 'total' shows total numbers of visual questions. Rows 'yes/no', 'number', and 'other' shows numbers of each type in that cluster. Rows '# agree' and '# disagree' show numbers of visual questions for which 10 answers agree (all are the same) and disagree (all are not the same), as in [5].

<sup>&</sup>lt;sup>4</sup>Many factors (*e.g.* initialization and number of clusters, chosen algorithms) affect the clustering result, but we have seen that similar clustering results are obtained with different parameter settings in preliminary experiments. Here we use the simplest algorithm, and a reasonable number of clusters

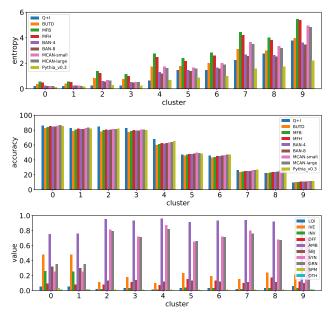


Figure 3. Values of (top) entropy, (middle) accuracy, and (bottom) reasons to differ for each cluster. Entropy values increase while accuracy decreases from cluster 0 (left) to 9 (right), while the predicted values of 9 reasons to differ are not well correlated to the order of clusters.

reasons to differ Average values obtained by reason classifiers [5] that output values from 0 (not that reason) to 1 (it is this reason) to each reason independently. We train classifiers on the subset of the VQA v2 training set provided by [5], then apply to VQA v2 validation set.

## 3.4. Discussion

Entropy suggests accuracy. We performed the clustering by using the entropy values of the three models based on Pythia v0.1 [17, 32]. Using a different base model may lead to different clustering results, however the values of entropy and accuracy of different state-of-the-art models exhibit similar trends; entropy values increase while accuracy decreases from cluster 0 to 9, as shown in Figure 3. This suggests that clusters with large (or small) entropy values have low (high) accuracy, as shown in Figure 4, and this tells us that entropy values are an important cue for predicting accuracy.

Entropy is different from reasons to differ and question types. Most frequent reasons to differ shown in [5] are AMB, SYN, and GRN, but Figure 3 shows that predicted values of those reasons are not well correlated to the order of clusters. For question types, Number and Other types looks not related to these clusters. Therefore our approach using entropy captures different aspects of visual questions.

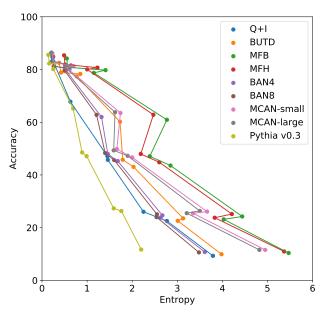


Figure 4. Scatter plot of entropy and accuracy of different models. Dots of each model are connected by lines in the order of cluster from top-left (cluster 0) to bottom-right (cluster 9).

Cluster 0 is easy, cluster 9 is hard. Level 1 (clusters 0 and 1) are dominant, and covers 44% of the entire validation set, including 99% of Yes/No type. Low entropy values and few number of unique answers (row 'ave # ans') of these cluster can be explained by the fact that typical answers are either 'Yes' or 'No'. Accuracy of Yes/No type is expected to be about 85% (Table 2), and it is close to the accuracy for these clusters. In contrast, level 3 (clusters 7, 8, and 9) looks much more difficult to answer. In particular, accuracy values of cluster 9 are about 10% compared to over 80% of level 1. This is due to the fact that visual questions with disagreed answers gather in this level; GT entropy is about 1.3, with more than five unique answers. However, values of DFF, AMB, SYN and GRN of level 3 are not so different from level 2, which may suggest that the quality of visual questions is not the main reason for difficulty.

Difficulty of the test set can be predicted. This finding enables us to evaluate the difficulty of visual questions in the test set. To see this, we applied the same base models (that are already trained and used for clustering) to visual questions in the test set, and computed entropy values to assign each visual questions to one of the 9 clusters. Rows with 'test set entropy' in Table 4 show the average entropy values of those test set visual questions. Assuming that the validation and test sets are similar in nature, we now are able to evaluate and predict the difficulty of test-set visual questions without computing accuracy. This is the most inter-

esting result, and we have released a list<sup>5</sup> that shows which visual questions in the train / val / test sets belong to which cluster. This would be extremely useful when developing a new model incorporating the difficulty of visual questions, and also when evaluating performances for different difficulty levels (not for different question types).

Qualitative evaluation of cluster difficulty Figure 6 shows some examples of visual questions in each level (from cluster 0, 4, 8, and 9). Entropy values of different methods tend to be larger in cluster 9, and visual questions in cluster 9 seem to be more difficult than those in cluster 0. To answer easy questions like "Is the catcher wearing safety gear?" or "What is the player's position behind the batter?" in cluster 0, images are not necessary and the Q model can correctly answer with low entropy. The question in cluster 9 at the bottom looks pretty difficult for the models to answer because of the ambiguity of the question ("What is this item?") and of the image (containing the photos of vehicles on the page of the book) even when the human annotators agree on the single ground-truth annotation.

# 3.5. Disagreement of predictions of different models

For difficult visual questions the number of unique answers is large, *i.e.* annotators highly disagree, while for easy questions numbers are small and they agree (5.39 for cluster 9, 1.72 for cluster 0). Now the following question arises; how much do different models (dis)agree, *i.e.* do they produce the same answer or different answers?

To see this, we define the overlap of model predictions. We have 9 models (BUTD, MFB, MFH, BAN-4/8, MCAN-small/large, Pythia v0.3 and v0.1 (Q+I)), and we define the "overlap" of the answers to be 9 when all models predict the same answer. For example, if we have two different answers to a certain question, each answer produced (supported) by respectively four and five models, then the answer overlaps are four and five, and we call the larger one a *max overlap*. Therefore, larger max overlap indicates a higher degree of agreement among the models. Figure 5 shows histograms of visual questions with different number of unique answers. The legend shows the details of max overlap.

For clusters 0 and 1, almost visual questions have one or two unique answers, and the models highly agree (max overlap of 9 is dominant). This is expected because most visual questions in these clusters are of Yes/No type, and models tend to agree by predicting either of two answers. Apparently clusters 2, 3, and 4 look similar; dominant max overlap is 9. This means that all of 9 models predict the same answer to almost half of visual questions even when annotators disagree to five different answers. In contrast, models predict different answers to visual questions of clusters 6-9 even when annotators agree and there is a single

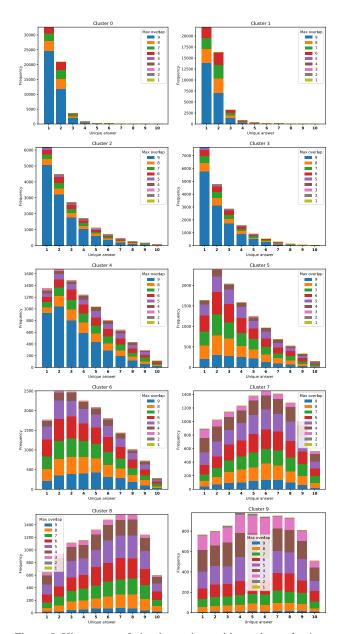


Figure 5. Histograms of visual questions with numbers of unique answers of ground truth annotations, and max overlap of predicted answers by 9 models.

ground truth answer (this is the case in the middle of cluster 9 column in Figure 6). Filling this gap may be a promising research direction for the next generation VQA models.

# 4. Conclusions

We have presented a novel way of evaluating the difficulty of visual questions of the VQA v2 dataset. Our approach is surprisingly simple, using three base models (I, Q, Q+I), predicting answer distributions, and computing entropy values to perform clustering with a simple k-

<sup>5</sup>https://github.com/tttamaki/vqd

Q: Is the catcher wearing safety gear? yes x 10

yes x 10

l: no (3.5743)
Q: yes (0.1587)
Q+l: yes (0.0000)
BUTD: yes (0.0208)
MFB: yes (0.0847)
MFH: yes (0.0500)
BAN8: yes (0.0001)
BAN8: yes (0.0001)
MCAN-5mall: yes (0.0021)
MCAN-large: yes (0.0011)
Pythia v0.3: yes (0.0000)
DFF: 0.0252
AMB: 0.8112
SYN: 0.1155
GRN: 0.3432



Q: How many couches? 2 x 10

I: yes (3.8459) Q: 2 (1.0936) Q+1: 2 (0.0778) BUTD: 2 (0.5404) MFB: 2 (0.7904) MFH: 1 (0.5361) BAN8: 1 (0.7193) MCAN-Small: 1 (0.6287) MCAN-large: 2 (0.7597) Pythia v 0.3: 2 (0.0000) DFF: 0.2610 AMB: 0.8199 SYN: 0.2606 GRN: 0.3463



Q: What is the player's position behind the batter? catcher x  $10\,$ 

I: no (3.6192)
Q: catcher (0.0060)
Q+I: catcher (0.0010)
BUTD: catcher (0.6489)
MFB: catcher (0.6489)
MFB: catcher (0.649)
BAN4: catcher (0.1424)
BAN4: catcher (0.1424)
MCAN-small: catcher (0.2366)
MCAN-large: catcher (0.0547)
Pythia v0.3: catcher (0.0169)
DFF: 0.0777
AMB: 0.9716
SYN: 0.9055
GRN: 0.8595



Q: When the man bought the sandwich, did he also buy a beverage? yes x 9 unknown x 1

I: no (4.2563)
Q: recently (3.9350)
Q+I: no (1.1395)
BUTD: yes (3.4916)
MFB: no (5.8325)
MFH: unknown (5.7551)
BAN4: no (1.5772)
BAN8: yes (3.2947)
MCAN-small: yes (2.7787)
MCAN-large: beer (4.5833)
Pythia v0.3: dinner (1.3259)
DFF: 0.0406
AMB: 0.8392
SYN: 0.1314
GRN: 0.3068



Q: What is the number of the bus?  $15 \times 10$ 

I: yes (4.2003) Q: 106 (4.0407) Q+I: 15 (1.0893) BUTD: 24 (3.6212) MFB: 1 (5.6941) MFH: 1 (5.8010) BAN4: 2 (3.2430) BAN8: 24 (2.7979) MCAN-small: 24 (5.3440) MCAN-large: 24 (4.5901) Pythia v0.3: 15 (1.0030) DFF: 0.2425 AMB: 0.7270 SYN: 0.2650 GRN: 0.3269



Q: What is the red object the Q: What is the r lady is holding? purse x 4 phone x 3 cell phone x 1 shawl x 1 bag x 1

I: no (4.6562)
Q: umbrella (4.6081)
Q+I: phone (1.3123)
BUTD: phone (0.9196)
MFB: phone (3.9944)
MFH: phone (3.3278)
BAN4: phone (2.8702)
BAN8: phone (2.8702)
BAN8: phone (2.5336)
MCAN-small: 0 (4.7481)
MCAN-large: phone (3.2502)
Pythia v0.3: cell phone (2.2534)
DFF: 0.0868
AMB: 0.9535 AMB: 0.9535 SYN: 0.8550 GRN: 0.8102



Q: What are the walls wooden? no x 8 mildew x 1 floor x 1

I: no (3.4617) I: no (3.4617)
Q: wooden (3.0521)
Q+I: tiles (3.4396)
BUTD: tile (1.9585)
MFB: tile (4.1105)
MFH: tile (2.7648)
BAN4: wall (4.2221)
BAN8: yes (4.3106)
MCAN-small: tile (4.1377)
MCAN-large: bathroom (5.6611)
Pythia v0.3: tile (1.0919)
DFF: 0.0519
AMB: 0.9777 AMB: 0.9777 SYN: 0.8639 GRN: 0.8408



Q: How many children are in this

I: no (4.3152) Q: 1 (1.8925) Q+1: 20 (3.0316) BUTD: 20 (3.4771) MFB: 15 (4.1004) MFH: 2 (4.5634) BAN4: 20 (3.3831) BAN8: 20 (3.1911) MCAN-small: 20 (4.1850) MCAN-large: 50 (3.6752) Pythia v0.3: 20 (2.5838) DFF: 0.3417 AMB: 0.8537 SYN: 0.4023 GRN: 0.4710



Q: Are these separate bananas or in a bunch? bunch x 7 both x 3

I: yes (4.6209)
Q: bananas (1.9826)
Q+I: ripe (2.4034)
BUTD: yes (1.8573)
MFB: both (3.2306)
MFH: bananas (3.5689)
BAN4: bananas (2.4397)
BAN8: bananas (1.5407)
MCAN-small: yes (5.5443)
MCAN-large: yes (4.7539)
Pythia v0.3: fruit (0.9098)
DFF: 0.0352
AMB: 0.9802 AMB: 0.9802 SYN: 0.8927 GRN: 0.8694



Q: What is the girl dancing? nowhere x 1 she is standing still x 1 can't say x 1

Cart Cay X 1

I: no (4.6622)
Q: kite (4.1933)
Q+I: nothing (4.6329)
BUTD: nothing (2.1515)
MFB: nothing (0.9136)
MFH: nothing (2.0973)
BAN8: nothing (2.0974)
BAN8: nothing (2.0974)
MCAN-Isrge: nothing (0.4467)
MCAN-Isrge: nothing (0.8140)
Pythia v0.3: nothing (0.3068)
DFF: 0.0425
AMB: 0.9386
SYN: 0.7515
GRN: 0.6765



Q: How many bananas are on Q: How many bananas are or display next to the oranges?  $17\times 2$   $19\times 2$   $19\times 2$   $15\times 1$   $40\times 1$   $16\times 1$   $50\times 1$   $30\times 1$   $20\times 1$ 

I: yes (4.3595) Q: 0 (3.4289) Q+I: 0 (3.6978) BUTD: 20 (3.8565) MFB: 6 (3.9011) MFH: 10 (3.9825) BAN4: 15 (3.8222) BAN8: 20 (3.7235) MCAN-small: 20 (3.7058) MCAN-large: 40 (4.3153) Pythia v0.3: 20 (3.9288) DFF: 0.4179 AMB: 0.8326 SYN: 0.3601 GRN: 0.4308



Q: What is this item? book x 10

I: no (4.0659)
C: vase (4.9113)
C+I: computer (4.1763)
BUTD: book (4.6128)
MFB: airplane (5.5921)
MFH: airplane (4.7641)
BAN8: plane (1.8373)
MCAN-small: book (2.6817)
MCAN-large: plane (1.3094)
Pythia v0.3: plane (1.3099)
DFF: 0.0372
AMB: 0.9748 AMB: 0.9748 SYN: 0.9105 GRN: 0.8587



cluster 0 cluster 4 cluster 8 cluster 9

Figure 6. Examples of visual questions in cluster 0, 4, 8, and 9 (from left to right). For each visual question, question text, answers, predicted answers and entropy values (in parenthesis) of each method are shown, followed by values of DFF, AMB, SYN and GRN.

means. Experimental results have shown that these clusters are strongly correlated with entropy and accuracy values of many models including state-of-the-art methods. By providing the correspondences between clusters and visual questions in the test set as the indicator of difficulty, our approach explores a novel aspect of evaluating performances of VQA models, suggesting a promising direction for future development of a next generation of VQA models.

# Acknowledgement

We would like to thank Yoshitaka Ushiku for his comments as a mentor of the PRMU mentorship program. This work was supported by JSPS KAKENHI grant number JP16H06540.

## References

- [1] Aishwarya Agrawal, Dhruv Batra, Devi Parikh, and Aniruddha Kembhavi. Don't just assume; look and answer: Overcoming priors for visual question answering. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2
- [2] 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 *The IEEE Conference on Com*puter Vision and Pattern Recognition (CVPR), June 2018. 4
- [3] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. VQA: Visual Question Answering. In *International Conference on Computer Vision (ICCV)*, 2015. 1, 2, 3
- [4] Abhishek Das Harsh Agrawal Larry Zitnick Devi Parikh Dhruv Batra. Human attention in visual question answering: Do humans and deep networks look at the same regions? In *International Conference on Machine Learning (ICML) Workshop on Visualization for Deep Learning*, 2016. 2
- [5] Nilavra Bhattacharya, Qing Li, and Danna Gurari. Why does a visual question have different answers? In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. 1, 2, 3, 5, 6
- [6] Jeffrey P. Bigham, Chandrika Jayant, Hanjie Ji, Greg Little, Andrew Miller, Robert C. Miller, Robin Miller, Aubrey Tatarowicz, Brandyn White, Samual White, and Tom Yeh. Vizwiz: Nearly real-time answers to visual questions. In Proceedings of the 23Nd Annual ACM Symposium on User Interface Software and Technology, UIST '10, pages 333–342, New York, NY, USA, 2010. ACM. 1, 3
- [7] Florian Daniel, Pavel Kucherbaev, Cinzia Cappiello, Boualem Benatallah, and Mohammad Allahbakhsh. Quality control in crowdsourcing: A survey of quality attributes, assessment techniques, and assurance actions. ACM Comput. Surv., 51(1):7:1–7:40, Jan. 2018. 1
- [8] Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra. Embodied Question Answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1, 3

- [9] Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José M.F. Moura, Devi Parikh, and Dhruv Batra. Visual Dialog. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 1, 3
- [10] Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus Rohrbach. Multimodal compact bilinear pooling for visual question answering and visual grounding. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 457–468, Austin, Texas, Nov. 2016. Association for Computational Linguistics. 4
- [11] 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 C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, Advances in Neural Information Processing Systems 28, pages 2296–2304. Curran Associates, Inc., 2015. 3
- [12] 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 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 1, 2, 3, 4
- [13] Yash Goyal, Akrit Mohapatra, Devi Parikh, and Dhruv Batra. Towards transparent ai systems: Interpreting visual question answering models. In *International Conference on Machine Learning (ICML) Workshop on Visualization for Deep Learning*, 2016. 2
- [14] Danna Gurari and Kristen Grauman. Crowdverge: Predicting if people will agree on the answer to a visual question. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, CHI '17, pages 3511–3522, New York, NY, USA, 2017. ACM. 1, 2, 3
- [15] Danna Gurari, Qing Li, Abigale J. Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P. Bigham. Vizwiz grand challenge: Answering visual questions from blind people. In *The IEEE Conference on Computer Vision* and Pattern Recognition (CVPR), June 2018.
- [16] Unnat Jain, Svetlana Lazebnik, and Alexander G. Schwing. Two can play this game: Visual dialog with discriminative question generation and answering. In *The IEEE Conference* on Computer Vision and Pattern Recognition (CVPR), June 2018. 3
- [17] Yu Jiang, Vivek Natarajan, Xinlei Chen, Marcus Rohrbach, Dhruv Batra, and Devi Parikh. Pythia v0.1: the winning entry to the VQA challenge 2018. *CoRR*, abs/1807.09956, 2018. 4, 6
- [18] 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 *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. 3
- [19] Kushal Kafle and Christopher Kanan. An analysis of visual question answering algorithms. In *ICCV*, 2017. 3
- [20] Jin-Hwa Kim, Jaehyun Jun, and Byoung-Tak Zhang. Bilinear attention networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett,

- editors, Advances in Neural Information Processing Systems 31, pages 1564–1574. Curran Associates, Inc., 2018. 4
- [21] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, Michael Bernstein, and Li Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. In *CoRR*, 2016.
- [22] Yikang Li, Nan Duan, Bolei Zhou, Xiao Chu, Wanli Ouyang, Xiaogang Wang, and Ming Zhou. Visual question generation as dual task of visual question answering. In *The IEEE Conference on Computer Vision and Pattern Recog*nition (CVPR), June 2018. 3
- [23] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *arXiv preprint arXiv:1908.02265*, 2019. 4
- [24] Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. Hierarchical question-image co-attention for visual question answering. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, Advances in Neural Information Processing Systems 29, pages 289–297. Curran Associates, Inc., 2016. 4
- [25] Mateusz Malinowski, Marcus Rohrbach, and Mario Fritz. Ask your neurons: A neural-based approach to answering questions about images. In *The IEEE International Confer*ence on Computer Vision (ICCV), December 2015. 2, 4
- [26] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *The IEEE Conference on Computer Vision and Pattern Recognition* (CVPR), June 2019. 3
- [27] Nasrin Mostafazadeh, Ishan Misra, Jacob Devlin, Margaret Mitchell, Xiaodong He, and Lucy Vanderwende. Generating natural questions about an image. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1802–1813, Berlin, Germany, Aug. 2016. Association for Computational Linguistics. 3
- [28] Jafar Muhammadi and Hamid Reza Rabiee. Crowd computing: a survey. *CoRR*, abs/1301.2774, 2013. 1
- [29] Mengye Ren, Ryan Kiros, and Richard Zemel. Exploring models and data for image question answering. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, Advances in Neural Information Processing Systems 28, pages 2953–2961. Curran Associates, Inc., 2015. 3
- [30] Tingke Shen, Amlan Kar, and Sanja Fidler. Learning to caption images through a lifetime by asking questions. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. 3
- [31] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors with online hard example mining. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. 3
- [32] Amanpreet Singh, Vivek Natarajan, Yu Jiang, Xinlei Chen, Meet Shah, Marcus Rohrbach, Dhruv Batra, and Devi Parikh. Pythia-a platform for vision & language research. In SysML Workshop, NeurIPS, volume 2018, 2018. 4, 6

- [33] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 3, 4
- [34] Ajeet Kumar Singh, Anand Mishra, Shashank Shekhar, and Anirban Chakraborty. From strings to things: Knowledgeenabled vqa model that can read and reason. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. 3
- [35] Guillermo Soberón, Lora Aroyo, Chris Welty, Oana Inel, Hui Lin, and Manfred Overmeen. Measuring crowd truth: Disagreement metrics combined with worker behavior filters. In Proceedings of the 1st International Conference on Crowdsourcing the Semantic Web Volume 1030, CrowdSem'13, pages 45–58, Aachen, Germany, Germany, 2013. CEURWS.org. 1
- [36] Damien Teney, Peter Anderson, Xiaodong He, and Anton van den Hengel. Tips and tricks for visual question answering: Learnings from the 2017 challenge. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 3, 4
- [37] Pei Wang and Nuno Vasconcelos. Towards realistic predictors. In *The European Conference on Computer Vision (ECCV)*, September 2018. 3
- [38] Peng Wang, Qi Wu, Chunhua Shen, Anthony R. Dick, and Anton van den Hengel. FVQA: fact-based visual question answering. *IEEE Trans. Pattern Anal. Mach. Intell.*, 40(10):2413–2427, 2018. 3
- [39] Xiaolong Wang and Abhinav Gupta. Unsupervised learning of visual representations using videos. In *The IEEE Inter*national Conference on Computer Vision (ICCV), December 2015. 3
- [40] Qi Wu, Damien Teney, Peng Wang, Chunhua Shen, Anthony Dick, and Anton van den Hengel. Visual question answering: A survey of methods and datasets. *Computer Vision and Image Understanding*, 163:21 – 40, 2017. Language in Vision.
- [41] Chun-Ju Yang, Kristen Grauman, and Danna Gurari. Visual question answer diversity. In *Proceedings of the Sixth AAAI Conference on Human Computation and Crowdsourcing, HCOMP 2018, Zürich, Switzerland, July 5-8, 2018*, pages 184–192, 2018. 2
- [42] Licheng Yu, Eunbyung Park, Alexander C. Berg, and Tamara L. Berg. Visual madlibs: Fill in the blank description generation and question answering. In *The IEEE Inter*national Conference on Computer Vision (ICCV), December 2015. 3
- [43] Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. Deep modular co-attention networks for visual question answering. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 4
- [44] Zhou Yu, Jun Yu, Chenchao Xiang, Jianping Fan, and Dacheng Tao. Beyond bilinear: Generalized multimodal factorized high-order pooling for visual question answering. *IEEE Transactions on Neural Networks and Learning Sys*tems, 29(12):5947–5959, 2018. 4

[45] Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. Visual7W: Grounded Question Answering in Images. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 3