

# Visual Question Answering for Blind People

Senyu Tong\*      Tianyi Lin\*      Yutian Zhao\*  
{senyut, tianyi12, yutianzh}@andrew.cmu.edu

## 1 Task Definition

In this project, we will use VizWiz dataset (Gurari et al., 2018) to address the visual question answering for blind people. The high level task of our project is to predict an accurate answer to a visual question given an image and question about it.

### 1.1 Subtasks

We will break down the task into the following three specific subtasks. Given the questions are collected from blind users of a mobile phone application, most of them are conversational and may not be answerable. Therefore, our first subtask would be to identify whether a question is valid and answerable according to the given image. The second task would be to output the correct answer if the question is answerable. Since every question in the VizWiz dataset has 10 answers, it is also important to determine how to select or combine these answers during training, which will be our third subtask.

### 1.2 Input-output Representation

Limited by computing resources, we would not try to modify any pre-training mechanism. Therefore, we would seek to incorporate task-specific knowledge to augment current large pre-trained model to boost performance. Specifically, according to our data analysis, we would like to address image framing issue particularly for images with text detected. Further discussion will be put into later sections.

Following (Tan and Bansal, 2019) and (Chen et al., 2020)’s work, we would like to use both joint representation and coordinated representations. Our tasks would be to generate two embeddings separately for images and texts, and a cross-modality output for their joint representation.

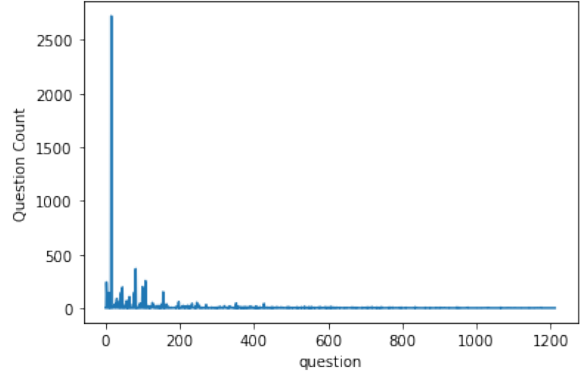


Figure 1: Question Frequency Distribution

For image representation, we would follow (Anderson et al., 2018). The input image embedding would be the features (positions and Region of Interest) of detected objects in it.

For language representation, the input would be word-level sentence embedding tokenized by WordPiece as in BERT (Devlin et al., 2019).

## 2 Data Analysis

### 2.1 Analysis of Questions

We examine the questions in both sentence and word level. We first calculate the distribution of the questions and observe that the percentage of questions that appear more than once is 54.4%. "What is this?" is the most common question that occupies more than 13% of all questions while the second most common question "What color is this?" is only about 1.7%. The frequency distribution is shown in Figure 1. We also analyze question diversity by computing statistics on sentence lengths. The mean question length is 6.76 words and the 25th and 75th percentile lengths are four and eight words respectively. The longest question has 62 words and the shortest one has only two words.

We also analyze the words in each sentence. We

---

\*Everyone Contributed Equally – Alphabetical order



Figure 2: Answer popularity word cloud excluding "unanswerable" and "unsuitable image" answers

consider the words that appear only once in all questions as rare words, which are about 2.64% of all words. The percentage of questions having at least one rare word is 12.09%. The most common first word of each question is "What", which is consistent with the observation of most common questions in the sentence-level analysis. We also observe that the questions in this dataset often begin with a rare first word. The percentage of questions starting with a first word that occurs for less than 5% and 10% of all questions is 32.68% and 38.24% respectively.

## 2.2 Analysis of Answers

We first analyze the percentage of answerable questions. Since VizWiz (Gurari et al., 2018) images are collected by blind people, a large portion of images have low quality and hence are not answerable. Based on our analysis, only 73.04% of visual questions are tagged as answerable.

We also analyze the diversity of answers by calculating statistics on answer length. The mean and median answer lengths are 1.66 and 1.00 words respectively, and the max answer length is 21 words. The result indicates that the visual questions in dataset tend to have short answers. We also generate a word cloud computed on all answers in the dataset excluding "unanswerable" and "unsuitable image", which is shown in Figure 2.

Lastly, we analyze the answer confidence level and agreement level. There is a total of 20523 questions in the dataset. Each question has 10 answers, and each answer has a confidence level ranging from "yes", "maybe" to "no". Among all questions, there are 18.82% questions have all 10 answers with a confidence level of "yes", 57.02% questions have at least 8 answers with a confidence level of "yes" and 94.85% questions have at least half of the answers with a confidence level of "yes".

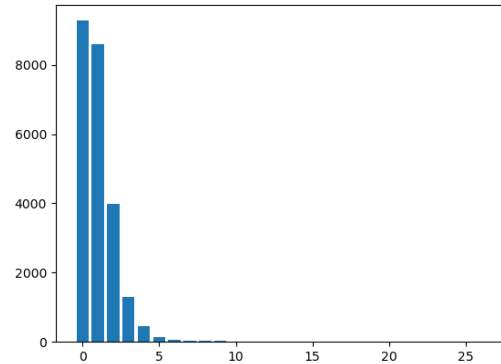


Figure 3: Number of bounding boxes

The dataset also has a pretty high human agreement level. Although we use a very strict agreement measure (exact string matching), we still observe that 40.54% of visual questions have more than 5 people agreed on the most popular answer, and 73.20% of visual questions have more than 3 people agreed on the most popular answer.

### 2.3 Analysis of Images

We run pre-trained YOLO (You Only Look Once) v3 model (Redmon et al., 2016) as a naive, coarse-grained, first-step baseline to analyze the training images. The weights of the model have been obtained by training on COCO dataset (Lin et al., 2015), and we hereby adapt the setting, trying to detect objects from 80 classes.

We count the number of bounding boxes on each image after we run objectness score thresholding and Non-Maximum Suppression to avoid overlapping. We observe that, we get 24,480 bounding boxes in total, and most images have either 0 (9,280 out of 23,953) or 1 (8,609 images) bounding box. The average number of bounding boxes per image is 1.02, and the third-quartile count is 2.0. The bar graph is shown in Figure 3.

After investigating the generated bounding boxes with labels, we observe that YOLO detects "Person" and "Bottle" objects most often. We also find that the coarse-grained model classify human hands as "Person" in most cases, and the 80 classes in COCO dataset are obviously insufficient.

From the experiment of running YOLO we conclude that, though most images may only contain one major object, it is difficult to successfully detect it. Ultrafine-grained semantic labels might be needed. Further study could be conducted when

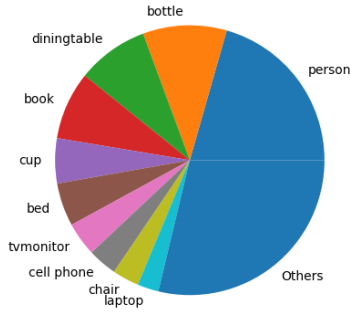


Figure 4: Predicted Classes Distribution

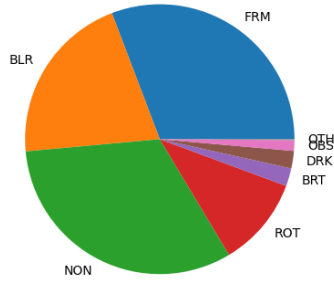


Figure 5: Images with Flaws Distribution

we incorporate the labeled image captions to our current predictions, to check where most errors are from. And we will adopt more advanced method such as Faster R-CNN (Changpinyo et al., 2019) to further explore the image data.

Apart from object detection bounding boxes analysis, we also analyze image quality. We observe that, a large number of pictures are covered by hands and many pictures are blurred. We then only focus on images that are tagged "answerable," and we find among 19,873 images that score 1 or below for category "recognizable," 52% of them have the "frame" flaw. Figure 5 shows the pie chart of flaws.

### 3 Related Work and Background

#### 3.1 Related Tasks

Over the years, there has been rapid progress in bridging vision and language in the research community. Various models have been developed and applied to a wide range of vision-and-language

tasks. We'll briefly describe the problem and training objective of five well-established tasks and include relevant papers.

**Visual Question Answering (VQA)** has received a lot of attention in recent years. A natural language question regarding an image is presented and the model will output the correct answer to it. Various models are proposed for this specific task (Agrawal et al., 2016; Tan and Bansal, 2019; Ben-younes et al., 2017; Fukui et al., 2016).

**Visual commonsense reasoning (VCR)** is similar to VQA, but besides predicting the answer, the model is also expected to select the correct answer justification among multiple choices (Lu et al., 2019; Li et al., 2019; Tan and Bansal, 2019).

**Visual grounding (VG)** aims to localize an image region or object given a natural language question. Common approaches to this task including bounding box proposals and reranking a set of image region candidates (Lu et al., 2019; Li et al., 2019; Fukui et al., 2016; Yu et al., 2018).

**Image captioning** recognizes the context of an image and adds a descriptive sentence to it. Models targeting this task are mainly attention-based neural networks that could extract deep features (Anderson et al., 2018; Donahue et al., 2016; Xu et al., 2016).

**Image-text matching** aims to measure the visual-semantic similarity between a text and an image. It has been widely applied to other applications such as image search for a given query (Kim et al., 2018; Nam et al., 2017).

Besides these five tasks, there are also many other popular vision-and-language transfer tasks such as caption-based image retrieval (Lu et al., 2019) and emotion recognition and sentiment analysis (Delbrouck et al., 2020).

#### 3.2 Related Techniques

In this section, we briefly review the techniques used in previous related work with a focus on the feature representation, attention mechanisms and fusion methods.

##### 3.2.1 Feature representation

Generating feature representation has been a crucial step in transforming raw data to meaningful features. There has been substantial past works in developing single modality models separately for vision and language representations. For image representations, common backbone models include Faster R-CNN (Ren et al., 2016), R2+1D-(152)

(Tran et al., 2015), VGG (Simonyan and Zisserman, 2015), ResNet (He et al., 2015) and fishNet (Sun et al., 2019). In terms of language representations, BERT (Devlin et al., 2019), a transformer language model pretrained on the BookCorpus (Zhu et al., 2015) and English Wikipedia is a popular choice nowadays.

Language embeddings and image embeddings are then combined through summation or multiplication (Lu et al., 2019; Delbrouck et al., 2020) to form a joint representation for fine tuning at a later stage. Another common approach to generate joint representation involves early fusion of image and text. Models are pretrained with multi-modal objectives such as masked language modeling with image and sentence-image prediction (Lu et al., 2019; Li et al., 2019; Tan and Bansal, 2019) to directly generate joint representations.

### 3.2.2 Attention Mechanism

Attention mechanism has been proven effective in many tasks including VQA. Previous attention approaches commonly used in VQA or related tasks can be classified into following four categories:

- Self attention that aggregate information inside each modality by query-key-value attention mechanism (Vaswani et al., 2017).
- Question-guided visual attention (Xu et al., 2016; Xu and Saenko, 2016) that compute attention on image region.
- Co-attention that jointly reason about both question and visual attention to interact across the two modalities. It uses image representation to guide the question attention and question representation to guide image attention (Lu et al., 2017, 2019; Yu et al., 2019). There also exist many variants of co-attention such as bilinear attention (Kim et al., 2018) and dense co-attention network (DCN) (Nguyen and Okatani, 2018) that considers interactions between every pair of question words and image regions.
- The intra- & inter-modal attention (DFAF) (Peng et al., 2019) that consider both inter-modality attention and intra-modality attention, where attention for intra-modality feature aggregation is dynamically modulated by the other modality using the pooled features.

### 3.2.3 Fusion Methods

The common multi-modal fusion approach is that visual and language features are extracted from the image and question independently as the first step, and then they are fused to compute the final results.

In previous studies, many works employed simple fusion methods that use element-wise product or summation of the visual and language features as final fused feature, and fed it to fully connected layers to predict results. (Lu et al., 2019; Delbrouck et al., 2020)

More complex fusion techniques have also been explored. For instance, Multimodal Compact Bilinear pooling (MCB) (Fukui et al., 2016) uses bilinear pooling method to compute the outer product between visual and language vectors as their fusion. In contrast to element-wise product approach, MCB allows a multiplicative interaction between all elements of both vectors. MUTAN fusion model (Ben-younes et al., 2017) is able to represent full bilinear interactions between visual and language modalities using Tucker decomposition of the correlation tensors, while maintaining the size of the model tractable. Multi-modality Latent Interaction Network (MLIN) (Gao et al., 2019) first encodes question and visual features into latent summarization vectors and then it go through the process of interaction, propagation and aggregation to achieve multi-modality reasoning.

## 4 Baselines

### 4.1 Model Selection

We would run the following 2 models as baselines.

- LSTM+CNN (Kazemi and Elqursh, 2017) uses CNN to extract image features and uses LSTM for question embedding. Multiple attention distributions are computed over the spatial dimensions of the image features. Then the concatenation of image feature glimpses and the state of the LSTM is fed to two fully connected layers.
- LXMERT (Tan and Bansal, 2019) constructs separate representations for images and texts, use (Anderson et al., 2018) for images. Have two encoders each for one modality with self-attention, and then a cross-modality encoder with cross-attention.

Since their pre-trained models are all available and we are interested in how their different attention / representation / fusion methods would effect



Method	Accuracy				Overall		
	Binary (yes/no)	Number	Nonanswerable	Color	total	acc	edit dist
LSTM+CNN	66.42%	30%	83.02%	60%	4319	<b>48.97%</b>	9.05
LXMERT	77.50%	35.68%	69.54%	64.82%	4317	<b>49.63%</b>	8.80

Table 1: Performance results for models in different question categories.

the performance. Their overall results on VizWiz dataset are listed in table 1. All metrics follow VizWiz official setting, i.e.

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

We are particularly interested in each model’s performance with flawed images, and we mainly focus on binary (yes or no) questions and questions related to numbers. From VizWiz challenge leader board we observe that current SoTA model achieves 56.33% accuracy, but most models have relatively weaker performance in answering questions related to “numbers,” with only 27.1% accuracy. We first aim to explore this discrepancy.

According to our data analysis, many input images are flawed. Blurring and framing issues are the most common types of flaw. Therefore, we expect to produce model results for each sub-category of image quality issues.

Aiming to boost performance with flawed images, we would test different fusing methods that inject de-noising technique, such as applying multi-stage progressive image restoration (Zamir et al., 2021) to our baseline models. Additionally, We would evaluate alternative representation / encoding / fusion methods following VizWiz official evaluation metric.

## 4.2 Baseline Implementation Details

### 4.2.1 LSTM+CNN

We ran a variant of model described in “Show, Ask, Attend, and Answer” paper (Kazemi and Elqursh, 2017) referencing github<sup>1</sup>. Visual features are extracted using a pretrained ResNet-152 model on ImageNet. Questions are embedded and encoded with LSTM. Image features and question embedding are then combined to compute multiple attention distributions over image features. The attended image features and the questions are concatenated and fed into two fully connected layers. We trained model for 40 epochs with batch size 128 and learning rate 0.001.

<sup>1</sup><https://github.com/DenisDsh/VizWiz-VQA-PyTorch>

### 4.2.2 LXMERT

We ran LXMERT referencing github<sup>2</sup> to run Vizwiz Dataset. We used the image representation extracted using Faster-RCNN and sentence representation extracted using Bert. Following the original paper (Tan and Bansal, 2019), we finetuned the model for 20 epochs with batch size 32. The validation accuracy calculated using the formula on Vizwiz website is 49.63%. We also calculated the average edit distance between the output answer and all the answers provided by crowd workers as an intrinsic evaluation of the baseline model.

## 4.3 Error Analysis

We defined an output answer to be incorrect if it doesn’t equal to the most popular answer provided by the 10 crowd workers. Following the metric given on Vizwiz website, we also break down the questions into the following categories for more detailed error analysis:

1. Binary: if at least three reference answers are either “yes” or “no”.
2. Number: if at least three reference answers are numbers. We define an answer as number if it is either an integer, float or double without any other characters.
3. Unanswerable: if at least three reference answers are unanswerable.
4. Color: if at least three reference answers are color.

The overall accuracy for LSTM+CNN and LXMERT is 48.97% and 49.63% respectively. The average edit distances for all models are around 8-9. Please see table 1 for detail test results.

For incorrect outputs on each of the above question categories, we conducted error analysis in terms of their questions and images. Please see table 2 and 3 for test statistics of each question category.

### 4.3.1 Question Analysis

**Binary questions** constitute 6.34% of the total dataset, and LSTM+CNN and LXMERT’s accu-

<sup>2</sup><https://github.com/airsplay/lxmert>

Method	Binary (yes/no)			
	yes/no swapped	unanswerable	ques begins with "can you .."	ques length stats
LSTM+CNN	31.0%	50.3%	19.3%	max: 157 min: 10 ave: 51
LXMERT	<b>52.9%</b>	<b>46.2%</b>	27.9%	max: 204 min: 10 ave: 54
Method	Number			
	% of non-numbers	unanswerable	ques begins with "can you .."	ques length stats
LSTM+CNN	76.7%	60.0%	10.0%	max: 129 min: 9 ave: 44
LXMERT	93.6%	50.0%	9.4%	max: 249 min: 8 ave: 43
Method	Nonanswerable			
	% of binary	% of numbers	ques begins with "can you .."	ques length stats
LSTM+CNN	9.6%	1.3%	7.6%	max: 249 min: 8 ave: 42
LXMERT	10.1%	0.8%	10.2%	max: 249 min: 8 ave: 42
Method	Color			
	% of non-color	% of unanswerable	ques begins with "can you .."	ques length stats
LSTM+CNN	55.0%	31.9%	0.6%	max: 118 min: 11 ave: 72
LXMERT	58.7%	24.0%	0.7%	max: 118 min: 11 ave: 77

Table 2: Error analysis of models’ incorrect outputs on different question categories.

racy is about 66 and 77% respectively. We first analyzed the number of workers that agreed on the incorrect output answer, since 10 workers themselves often give different reference answers. Figure 6 is the plot for the distribution of wrong answer votes for each model. We can see a similar pattern for all of the models, which the majority of incorrect answers completely different from all the reference answers provided by the crowd workers (0 votes mean that no worker give the same reference answer as the output answer).

However, it’s interesting to note that LXMERT is able to identify 52.88% of the binary questions but output opposite answers while LSTM+CNN only outputs "yes" or "no" for binary questions 31% among all. Moreover, LSTM+CNN outputs more "unanswerable" than LXMERT on binary questions and these statistics indicate that LXMERT performs better than LSTM+CNN on binary questions, possibly due to their deeper model architecture.

**Number questions** only constitute 1.63% of the total dataset, and the baseline models’ accuracy is about 30-36%. We also analyzed the number of workers that agreed on the incorrect output answer and similarly to binary questions, more than 70% of the output answer doesn’t exist in the workers reference answers. All three models tend to out-

put "unanswerable" to number questions, failing to recognize the numbers in the image.

**Unanswerable questions** constitute 25.6% of the total dataset, and the baseline models’ accuracy is above 83%.

We’re interested in the models’ ability to distinguish different kinds of questions and from the result in table 1, we can see that 10% of the incorrect answers for unanswerable questions are binary while only about 1% are number. LXMERT output "unanswerable" 63% of the time for binary questions and output binary answers for 10% of the unanswerable questions. We would like to identify the reasons behind such confusion between different question categories and improve upon it.

#### 4.3.2 Image Analysis

We performed image error analysis on vizwiz visual question analysis validation dataset for baseline LSTM+CNN model (Kazemi and Elqursh, 2017) and LXMERT model (Tan and Bansal, 2019). We joined the vizwiz image quality issues dataset (Chiu et al., 2020) with our predictions from baseline models to classify input images into 6 flaw categories. Vizwiz image quality issues dataset (Chiu et al., 2020) is annotated with number of votes, out of five crowd workers, for quality flaws. We classified an image as one of the flaw types if at least one

Flaw	LSTM+CNN Accuracy					LXMERT Accuracy				
	binary	number	unans.	color	total	binary	number	unans.	color	total
BLR	68.52%	30.95%	82.21%	61.45%	51.46%	80.37%	34.52%	77.16%	65.95%	51.51%
FRM	67.42%	32.77%	82.49%	63.06%	50.67%	76.28%	37.85%	78.92%	66.15%	50.74%
DRK	70.18%	46.15%	80.47%	61.90%	54.49%	77.78%	48.72%	76.29%	70.83%	54.76%
BRT	72.99%	30.00%	86.76%	61.31%	53.03%	82.18%	33.33%	80.72%	67.26%	53.51%
OBS	60.16%	22.22%	82.50%	72.00%	57.18%	81.30%	11.11%	78.45%	71.33%	57.57%
ROT	67.67%	25.49%	81.70%	59.49%	50.42%	74.00%	39.22%	77.66%	61.54%	50.14%
NON	65.44%	26.44%	84.05%	59.03%	46.82%	79.47%	32.76%	80.09%	62.85%	47.83%
<b>Total</b>	<b>66.42%</b>	<b>30.05%</b>	<b>83.02%</b>	<b>60.10%</b>	<b>48.97%</b>	<b>77.49%</b>	<b>35.68%</b>	<b>78.94%</b>	<b>64.82%</b>	<b>49.63%</b>

Table 3: Performance analysis of different image quality issue types on LSTM+CNN and LXMERT model. The quality issue categories are provided by vizwiz image quality issues dataset: Blur(BLR), Framing(FRM), Dark(DRK), Bright(BRT), Obscured(OBS), Rotation(ROT) and No Flaws(NON). In each of the image flaw category, we calculate accuracy for 4 different answer types: binary(yes/no), number, unanswerable and color.

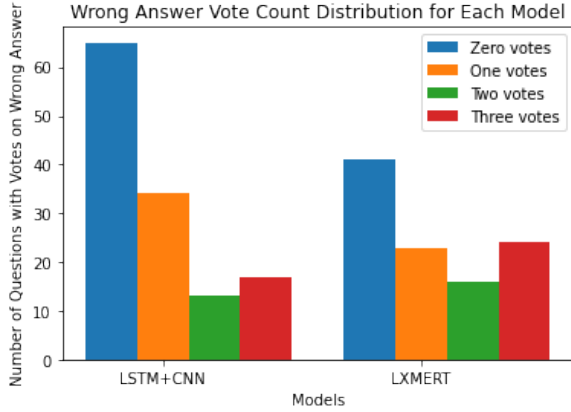


Figure 6: Wrong Answer Vote Count Distribution on Binary Questions

crowd worker votes for it. The flaw categories are defined as following: Blur(BLR), Framing(FRM), Dark(DRK), Bright(BRT), Obscured(OBS), Rotation(ROT) and No Flaws(NON). We further calculated accuracy for 4 different question types in each of the image flaw categories: binary(yes/no), number, unanswerable and color. See Table 3 for the result.

Based on our analysis, Framing(FRM), Blur(BLR), and Rotation(ROT) are the three most common types of flaws, consisting of 79.28%, 64.24%, and 30.63% of total images, see Figure 5. Among all 6 flaw types, Obscured(OBS) performs best with around 57% overall accuracy for both baseline models, while Rotation(ROT) and Framing(FRM) perform worst with around 50% accuracy for both models. We also broke down flaw types into different question types and investigated model performance on each

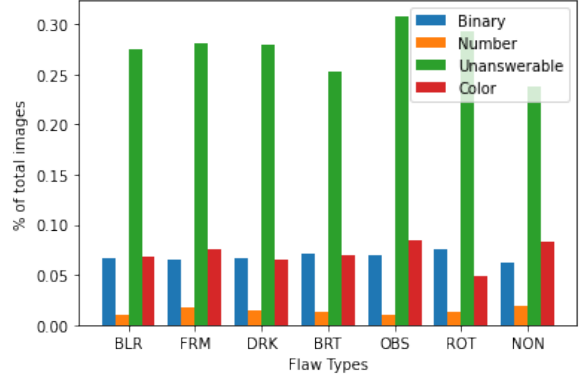


Figure 7: Image flaw type percentage based on different question types

question type. Although LXMERT achieved a better overall result than LSTM+CNN model, both models achieved similar results on flawed images. Furthermore, LXMERT performed worse on the unanswerable question type while performed better on all other question types than LSTM+CNN model, which indicates a future improvement for LXMERT on unanswerable question type.

It is also interesting to note that No Flaw(NON) image categories indeed perform worse than flaw image categories, yielding 46.82% and 47.83% overall accuracy for LSTM+CNN and LXMERT model, respectively. In order to explain this phenomenon, we further classified image flaw types into different question types in Figure 7. We noticed that flawed images have higher unanswerable question types than NON (no flaw) image categories. And since models performed best on "unanswerable" question types, which explained why flawed images have better overall performance.

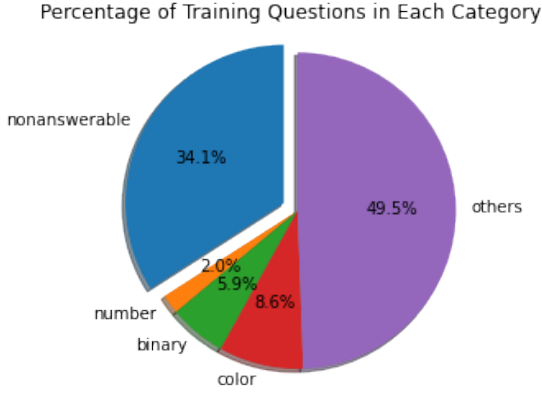


Figure 8: Question Type Distribution

This also suggests that allowing models first to figure out if an image is a flawed image will be beneficial as it changes the background probability for unanswerable output.

## 5 Proposed Approach

We categorize questions into five types as in Figure 8. We observe that models tend to have difficulties in distinguishing answer types as in Table 2. Based on our error analysis in previous section, we believe a pipeline approach would benefit by allowing models to first figure out which type of visual scene and question they are facing. Although pipelined model has been widely adopted in uni-modal tasks such as Question Answering (Kwiatkowski et al., 2019; Wang et al., 2020b), Handwritten Text Recognition (Chung and Delteil, 2019) and Data Collection and Processing (Di Mitri et al., 2019), it was rarely used in multi-modal settings. To the best of our knowledge, we are the first to experiment with the hierarchical pipeline approach for multi-modal Visual Question Answering task.

Inspired by (Wang et al., 2020b) that proposed a two-stage training procedure for Question Answering task, we propose to use a separate classifier that predicts the question types on top of the main VQA model. We then feed the output of the classifier to the Cross-Modality Encoder of the main VQA model. A more detailed discussion would be in Section 5.2. With this additional information on question type explicitly injected, we believe the model would be less likely to mix different kinds of questions.

### 5.1 Training Data Generation

As we can see from the Figure 9, the original vizwiz dataset has an uneven distribution over five ques-

tion types. So models trained with the original dataset will have a huge bias over "unanswerable" and "others" question types as these two types consist of the majority training dataset. This would lead the model to favor "unanswerable" and "others" as output. In order to eliminate such bias, two methods could be adopted: firstly, rebalance the dataset and let the five categories have an equal portion. Secondly, we can use Focal Loss (Lin et al., 2018), a weighted loss function, to ensure majority labels do not dominate the loss.

### 5.2 Model Architecture

As shown in Figure 9, our model now has two parts. The left is a classifier, and the right is the main VQA model.

**Classifier** In our work, we are motivated to utilize a prior question type distribution as an augmentation to the main VQA model. For the choice of our classifier architecture, we could in principle use any multi-modal models that can conduct classification tasks given a pair of image and question, such as LSTM+CNN, Pythia (Singh et al., 2019) and Uniter (Chen et al., 2020). We could even replicate the main VQA model with the training objective modified. In this project, we would use LXMERT as the classifier to generate this "prior" distribution. As a stretch goal, we would also like to experiment and compare other language-vision models if time and resource allowed.

Following (Tan and Bansal, 2019), we use the image features generated by Faster-RCNN and word embeddings generated by Bert as inputs and run three encoders: object relationship encoder, language encoder and a cross-modality encoder. The classifier has two intermediate outputs for language and vision respectively before the final MLP classification layer is applied. This final layer will predict the most probable question type.

**VQA Model** Besides the classifier, we propose to use LXMERT as our main VQA model as LXMERT outperformed LSTM+CNN in our baseline result analysis. Due to the fact that the classifier and the main VQA models are inherently multi-modal, both models will generate image, language and cross-modality representations as intermediate outputs, making further fusion desirable. This also makes our approach uniquely multimodal.

After the classifier outputs a prior probability distribution, the VQA model will adapt the image, language and cross-modality intermediate outputs



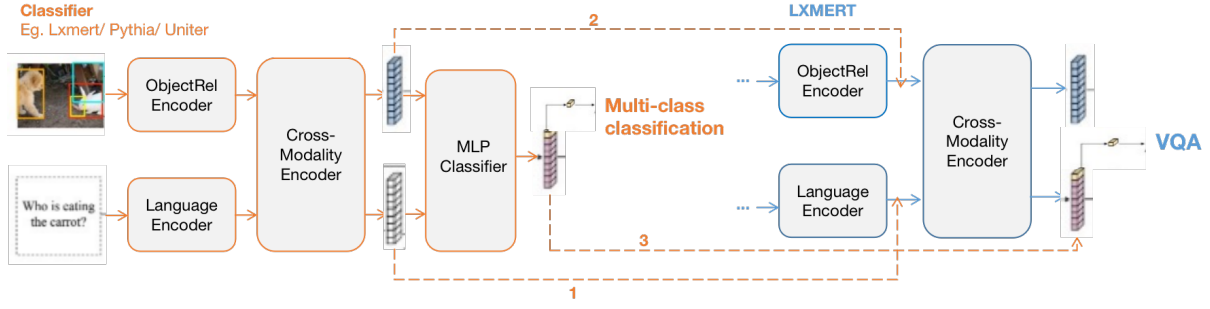


Figure 9: Proposed model architecture. Our model consists of a vision-language multi-class classifier and a main VQA model.

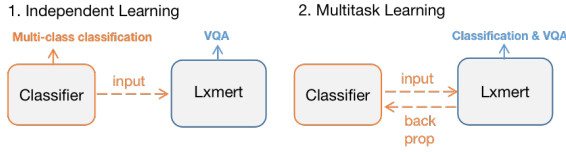


Figure 10: Finetune Strategy. The left figure represents two models learning independently and the right represents multi-task learning.

from the classifier. As shown by the dashed lines 1 and 2 in Figure 9, we augment the language and image representations generated by the uni-modal encoders in the main LXMERT model with the two representations generated by the cross-modality encoders in the classifier model. Additionally, as indicated by dashed line 3, we fuse the hidden states of the final layer in the classifier into the output representation of the main LXMERT’s cross-modality encoder.

### 5.3 Fusion Methods

We would experiment on different fusion methods to integrate the image, language and the final representations described above. Specifically, we could concatenate them directly, take their element-wise product, or use Multi-Layer Perceptrons to integrate their dimensions.

### 5.4 Finetune Strategy

In terms of finetune strategy, we would like to experiment training the classifier model and the main VQA model jointly and independently as shown in Figure 10.

**Multitask Learning** We could train the classifier and the main VQA model simultaneously with loss combined for the two tasks. We expect it to generate more robust intermediate representations but result in a more complicated model structure and prone to overfitting. A more detailed discus-

sion of Multi-Task Learning would be in section 5.5.

**Independent Learning** Alternatively, in order to prevent over-fitting, our classifier could be completely separate and independent of the main model, thus the classifier’s parameters would not be updated during the main model’s finetuning process.

### 5.5 Multi-Task Learning

As shown in Figure 10, in order to train the classifier and main VQA model jointly, we would define a new loss function to achieve this multi-task learning. First, following (Tan and Bansal, 2019), we would minimize the binary cross entropy for each task:

$$L = -y^* \log(prob) - (1 - y^*) \log(1 - prob) \quad (1)$$

Second, we will then train the entire model with the weighted loss:

$$L_{final} = \alpha L_{classifier} + \beta L_{VQA} \quad (2)$$

Here  $\alpha$  and  $\beta$  are two hyperparameters that we will also finetune. Furthermore, we could also replace  $L_{classifier}$  with Focal Loss (Lin et al., 2018) we mentioned earlier to eliminate dataset bias. Finally, this loss will be backpropagated through the entire model and the model parameters will be updated.

### 5.6 Training Issues

Despite the fact that multi-modal networks receive more information, they are more prone to over-fitting due to increased complexity (Wang et al., 2020a). Our proposed model adopts a pipeline approach to combine two complex multi-modal networks together, and therefore it is more likely to be overfitted.

We adopt two approaches to help us better understand and potentially mitigate this issue. First, we

will experiment with two training/finetune strategies: Independent Learning and Multitask Learning as mentioned in section 5.4.

Second, we could use a simpler network as our classifier if overfitting occurs. For example, we could replace LXMERT with a simple LSTM+CNN network or even use a shared parameter model for both question-type classification and VQA tasks through multi-task learning with a combined loss function.

## References

- Aishwarya Agrawal, Jiasen Lu, Stanislaw Antol, Margaret Mitchell, C. Lawrence Zitnick, Dhruv Batra, and Devi Parikh. 2016. [Vqa: Visual question answering](#).
- Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. [Bottom-up and top-down attention for image captioning and visual question answering](#).
- Hedi Ben-younes, Rémi Cadene, Matthieu Cord, and Nicolas Thome. 2017. [Mutan: Multimodal tucker fusion for visual question answering](#).
- Soravit Changpinyo, Bo Pang, Piyush Sharma, and Radu Soricut. 2019. [Decoupled box proposal and featurization with ultrafine-grained semantic labels improve image captioning and visual question answering](#).
- Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. [Uniter: Universal image-text representation learning](#).
- Tai-Yin Chiu, Yinan Zhao, and Danna Gurari. 2020. [Assessing image quality issues for real-world problems](#).
- Jonathan Chung and Thomas Delteil. 2019. [A computationally efficient pipeline approach to full page of-line handwritten text recognition](#). In *2019 International Conference on Document Analysis and Recognition Workshops (ICDARW)*, volume 5, pages 35–40.
- Jean-Benoit Delbrouck, Noé Tits, Mathilde Brous-miche, and Stéphane Dupont. 2020. [A transformer-based joint-encoding for emotion recognition and sentiment analysis](#). *Second Grand-Challenge and Workshop on Multimodal Language (Challenge-HML)*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [Bert: Pre-training of deep bidirectional transformers for language understanding](#).
- Daniele Di Mitri, Jan Schneider, MM Specht, and HJ Drachsler. 2019. [Multimodal pipeline: a generic approach for handling multimodal data for supporting learning](#). In *First workshop on AI-based Multimodal Analytics for Understanding Human Learning in Real-world Educational Contexts*.
- Jeff Donahue, Lisa Anne Hendricks, Marcus Rohrbach, Subhashini Venugopalan, Sergio Guadarrama, Kate Saenko, and Trevor Darrell. 2016. [Long-term recurrent convolutional networks for visual recognition and description](#).
- Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus Rohrbach. 2016. [Multimodal compact bilinear pooling for visual question answering and visual grounding](#).
- Peng Gao, Haoxuan You, Zhanpeng Zhang, Xiaogang Wang, and Hongsheng Li. 2019. [Multi-modality latent interaction network for visual question answering](#).
- Danna Gurari, Qing Li, Abigale J. Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P. Bigham. 2018. [Vizwiz grand challenge: Answering visual questions from blind people](#).
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. [Deep residual learning for image recognition](#).
- Vahid Kazemi and Ali Elqursh. 2017. [Show, ask, attend, and answer: A strong baseline for visual question answering](#).
- Jin-Hwa Kim, Jaehyun Jun, and Byoung-Tak Zhang. 2018. [Bilinear attention networks](#).
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. [Natural questions: a benchmark for question answering research](#). *Transactions of the Association for Computational Linguistics*, 7:453–466.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. [Visualbert: A simple and performant baseline for vision and language](#).
- Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. 2018. [Focal loss for dense object detection](#).
- Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. 2015. [Microsoft coco: Common objects in context](#).
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. [Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks](#).

- Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. 2017. [Hierarchical question-image co-attention for visual question answering](#).
- Hyeonseob Nam, Jung-Woo Ha, and Jeonghee Kim. 2017. [Dual attention networks for multimodal reasoning and matching](#).
- Duy-Kien Nguyen and Takayuki Okatani. 2018. [Improved fusion of visual and language representations by dense symmetric co-attention for visual question answering](#).
- Gao Peng, Zhengkai Jiang, Haoxuan You, Pan Lu, Steven Hoi, Xiaogang Wang, and Hongsheng Li. 2019. [Dynamic fusion with intra- and inter- modal- ity attention flow for visual question answering](#).
- Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. 2016. [You only look once: Unified, real-time object detection](#).
- Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2016. [Faster r-cnn: Towards real-time object detection with region proposal networks](#).
- Karen Simonyan and Andrew Zisserman. 2015. [Very deep convolutional networks for large-scale image recognition](#).
- Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. 2019. [Towards vqa models that can read](#).
- Shuyang Sun, Jiangmiao Pang, Jianping Shi, Shuai Yi, and Wanli Ouyang. 2019. [Fishnet: A versatile backbone for image, region, and pixel level prediction](#).
- Hao Tan and Mohit Bansal. 2019. [LXMERT: Learning cross-modality encoder representations from transformers](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5100–5111, Hong Kong, China. Association for Computational Linguistics.
- Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. 2015. [Learning spatiotemporal features with 3d convolutional networks](#).
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#).
- Weiyao Wang, Du Tran, and Matt Feiszli. 2020a. [What makes training multi-modal classification networks hard?](#)
- Xuguang Wang, Linjun Shou, Ming Gong, Nan Duan, and Daxin Jiang. 2020b. [No answer is better than wrong answer: A reflection model for document level machine reading comprehension](#).
- Huijuan Xu and Kate Saenko. 2016. [Ask, attend and answer: Exploring question-guided spatial attention for visual question answering](#).
- Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, and Yoshua Bengio. 2016. [Show, attend and tell: Neural image caption generation with visual attention](#).
- Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. 2019. [Deep modular co-attention networks for visual question answering](#).
- Zhou Yu, Jun Yu, Chenchao Xiang, Zhou Zhao, Qi Tian, and Dacheng Tao. 2018. [Rethinking diversified and discriminative proposal generation for visual grounding](#).
- Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. 2021. [Multi-stage progressive image restoration](#).
- Yukun Zhu, Ryan Kiros, Richard Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. [Aligning books and movies: Towards story-like visual explanations by watching movies and reading books](#).