# Customized Image Narrative Generation via Interactive Visual Question Generation and Answering

Andrew Shin[1]    Yoshitaka Ushiku[1]    Tatsuya Harada[1,2]
[1]The University of Tokyo, [2]RIKEN
{andrew,ushiku,harada}@mi.t.u-tokyo.ac.jp

## Abstract

*Image description task has been invariably examined in a static manner with qualitative presumptions held to be universally applicable, regardless of the scope or target of the description. In practice, however, different viewers may pay attention to different aspects of the image, and yield different descriptions or interpretations under various contexts. Such diversity in perspectives is difficult to derive with conventional image description techniques. In this paper, we propose a customized image narrative generation task, in which the users are interactively engaged in the generation process by providing answers to the questions. We further attempt to learn the user's interest via repeating such interactive stages, and to automatically reflect the interest in descriptions for new images. Experimental results demonstrate that our model can generate a variety of descriptions from single image that cover a wider range of topics than conventional models, while being customizable to the target user of interaction.*

## 1. Introduction

Recent advances in visual language field enabled by deep learning techniques have succeeded in bridging the gap between vision and language in a variety of tasks, ranging from describing the image [15, 7, 28, 29] to answering questions about the image [2, 5]. Such achievements were possible under the premise that there exists a set of ground truth references that are universally applicable regardless of the target, scope, or context. In real-world setting, however, image descriptions are prone to an infinitely wide range of variabilities, as different viewers may pay attention to different aspects of the image in different contexts, resulting in a variety of descriptions or interpretations. Due to its subjective nature, such diversity is difficult to obtain with conventional image description techniques.

In this paper, we propose a customized image narrative generation task, in which we attempt to actively engage the
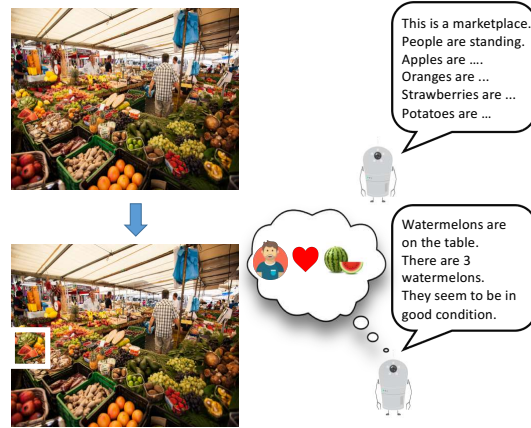


Figure 1: Example of conventional image description (top) and customized image narrative (bottom).

users in the description generation process by asking questions and directly obtaining their answers, thus learning and reflecting their interest in the description. We use the term *image narrative* to differentiate our image description from conventional one, in which the objective is fixed as depicting factual aspects of global elements. In contrast, *image narratives* in our model cover a much wider range of topics, including subjective, local, or inferential elements.

We first describe a model for automatic image narrative generation from single image without user interaction. We develop a self Q&A model to take advantage of wide array of contents available in visual question answering (VQA) task, and demonstrate that our model can generate image descriptions that are richer in contents than previous models. We then apply the model to interactive environment by directly obtaining the answers to the questions from the users. Through a wide range of experiments, we demonstrate that such interaction enables us not only to customize the image description by reflecting the user's choice in the current image of interest, but also to automatically apply the learned preference to new images (Figure 1 ).

1

## 2. Related Works

**Visual Language:** The workflow of extracting image features with convolutional neural network (CNN) and generating captions with long short-term memory (LSTM) [11] has been consolidated as a standard for image captioning task. [15] generated region-level descriptions by implementing alignment model of region-level CNN and bidirectional recurrent neural network (RNN). [13] proposed DenseCap that generates multiple captions from an image at region-level. [12] built SIND dataset whose image descriptions display a more casual and natural tone, involving aspects that are not factual and visually apparent. While this work resembles the motivation of our research, it requires a sequence of images to fully construct a narrative.

Visual question answering (VQA) has escalated the interaction of language and vision to a new stage, by enabling a machine to answer a variety of questions about the image, not just describe certain aspects of the image. A number of different approaches have been proposed to tackle VQA task, but classification approach has been shown to outperform generative approach [1, 14]. [8] proposed multimodal compact bilinear pooling to compactly combine the visual and textual features. [24] proposed an attention-based model to select a region from the image based on text query. [19] introduced co-attention model, which not only employs visual attention, but also question attention.

**User Interaction:** Incorporating interaction with users into the system has rapidly become a research interest. Visual Dialog [5] actively involves user interaction, which in turn affects the responses generated by the system. Its core mechanism, however, functions in an inverse direction from our model, as the users ask the questions about the image, and the system answers them. Thus, the focus is on extending the VQA system to a more context-dependent, and interactive direction. On the other hand, our model's focus is on generating customized image descriptions, and user interaction is employed to learn the user's interest, whereas Visual Dialog is not concerned about the users themselves.

[6] introduces an interactive game, in which the system attempts to localize the object that the user is paying attention to by asking relevant questions that narrow down the potential candidates, and obtaining answers from the users. This work is highly relevant to our work in that user's answers directly influence the performance of the task, but our focus is on contents generation instead of object localization or gaming. Also, our model not only utilizes user's answer for current image, but further attempts to apply it to new images. Recent works in reinforcement learning (RL) have also employed interactive environment by allowing the agents to be taught by non-expert humans [4]. However, its main purpose is to assist the training of RL agents, while our goal is to learn the user's interest specifically.
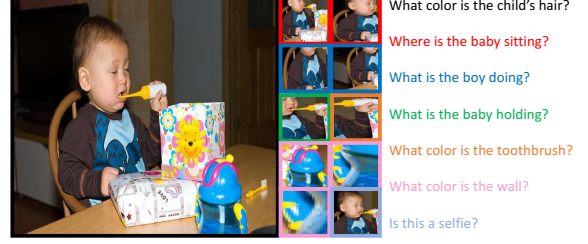


Figure 2: Example of regions extracted from the image, and the questions generated from each region.

## 3. Automatic Image Narrative Generation

We first describe a model to generate image narrative that covers a wide range of topics without user interaction. We propose a self Q&A model where questions are generated from multiple regions, and VQA is applied to answer the questions, thereby generating image-relevant contents.

**Region Extraction:** Following [9], we first extract region candidates from the feature map of an image, by applying linear SVM trained on annotated bounding boxes at multiple scales, and applying non-maximal suppression. The region candidates then go through inverse cascade from upper, fine layer to lower, coarser layers of CNN, in order to better-localize the detected objects. This results in region proposals that are more contents-oriented than selective search [26] or Edge Boxes [17]. We first extracted top 10 regions per image. Figure 2 shows an example of the regions extracted in this way. In the experiments to follow, we set the number of region proposals $K$ as 5, since the region proposals beyond top 5 tended to be less congruent, thus generating less relevant questions.

**Visual Question Generation:** In image captioning task, it is conventional to train an LSTM with human-written captions as ground truth annotations. On the other hand, in VQA task, questions are frequently inserted to LSTM in series with fixed image features, and the answers to the questions become the ground truth labels to be classified. Instead, we replace the human-written captions with human-written questions, so that LSTM is trained to predict the question, rather than caption.

Given an image $I$ and a question $Q = (q_0,...q_N)$, the training proceeds as in [28]:

$$x_{-1} = CNN(I), x_t = W_e q_t, p_{t+1} = LSTM(x_t) \quad (1)$$

where $W_e$ is a word embedding, $x_t$ is the input features to LSTM at $t$, and $p_{t+1}$ is the resulting probability distribution for the entire dictionary at $t$. In the actual generation of questions, it will be performed over all region proposals $r_0,...,r_N \in I$:

$$x_{-1} = CNN(r_i), x_t = W_e q_{t-1}$$
$$q_t = \max_{q \in p} p_{t+1} = \mathrm{argmax} LSTM(x_t) \quad (2)$$

for $q_0,...q_N \in Q_{r_i}$. Figure 2 shows examples of questions generated from each region including the entire image. As

Table 1: Examples of questions generated using non-visual questions in VQG dataset.

| Image | Generated Questions |
|---|---|
|  | • What is the player's name?<br>• What is he speaking about?<br>• What is the score?<br>• Is this costume for a race?<br>• Has he worked there? |
|  | • Can the boy win the prize?<br>• Was this a charity event?<br>• What is she looking at?<br>• What are they waiting for?<br>• Who is that guy?<br>• What is he looking at? |

shown in the figure, by focusing on different regions and extracting different image features, we can generate multiple image-relevant questions from single image.

So far, we were concerned with generating "visual" questions. We also seek to generate "non-visual" questions. [21] generated questions that a human may naturally ask and require common-sense and inference. We examined whether we can train a network to ask multiple questions of such type by visual cues. We replicated the image captioning process described above, with 10,000 images of MS COCO and Flickr segments of VQG dataset, with 5 questions per image as the annotations. Examples of questions generated by training the network solely with non-visual questions are shown in Table 16.

**Visual Question Answering:** We now seek to answer the questions generated. We train the question answering system with VQA dataset [2]. Question words are sequentially encoded by LSTM as one-hot vector. Hyperbolic tangent non-linearity activation was employed, and element-wise multiplication was used to fuse the image and word features, from which softmax classifies the final label as the answer for visual question. We set the number of possible answers as 1,250.

As we augmented the training data with "non-visual" questions, we also need to train the network to "answer" those non-visual answers. Since [21] provides the questions only, we collected the answers to these questions on Amazon Mechanical Turk. Since many of these questions cannot be answered without specific knowledge beyond what is seen in the image (*e.g.* "*what is the name of the dog?*"), we encouraged the workers to use their imagination, but required them to come up with answers that an average person might also think of. For example, people frequently answered the question "*what is the name of the man?*" with "*John*" or "*Tom*." Such non-visual elements add vividness and story-like characteristics to the narrative as long as they are compatible with the image, even if not entirely verifiable.
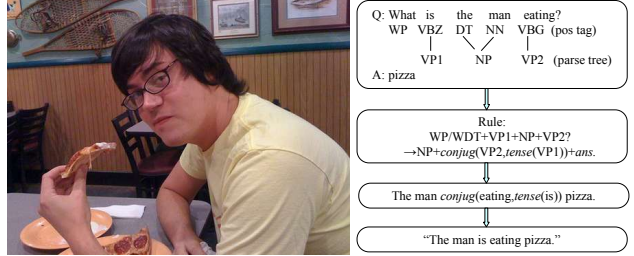


Figure 3: Example of question and answer converted to a declarative sentence by conversion rule.

**Natural Language Processing:** We are now given multiple pairs of questions and answers about the image. By design of the VQA dataset, which mostly comprises simple questions regarding only one aspect with the answers mostly being single words, the grammatical structure of most questions and answers can be reduced to a manageable pool of patterns. Exploiting these design characteristics, we combine the obtained pairs of questions and answers to a declarative sentence by application of rule-based transformations, as in [23, 25].

We first rephrase the question to a declarative sentence by switching word positions, and then insert the answers to its appropriate position, mostly replacing *wh*-words. For example, a question "*What is the man holding?*" is first converted to a declarative statement "*The man is holding what*" and the corresponding answer "*frisbee*" replaces "*what*" to make "*The man is holding frisbee.*" Part-of-speech tags with limited usage of parse tree were used to guide the process, particularly conjugation according to tense and plurality. Figure 3 illustrates the workflow of converting question and answer to a declarative sentence. See Supplemental Material for specific conversion rules. Part-of-speech tag notation is as used in PennTree I Tags [20].

# 4. Interactive Image Narrative Generation

We now extend the automatic image narrative generation model described in Section 3 to interactive environment, in which users participate in the process by answering questions about the image, so that generated narrative varies depending on the user input provided.

## 4.1. Applying Interaction within the Same Images

### 4.1.1 Question with Multiple Possible Answers

As discussed earlier, we attempt to reflect user's interest by asking questions that provide visual context. The foremost prerequisite for the interactive questions to perform that function is the possibility of various answers or interpretations. In other words, a question whose answer is so obvious that it can be answered in an identical way would not be valid as an interactive question. In order to make sure that each generated question allows for multiple possible
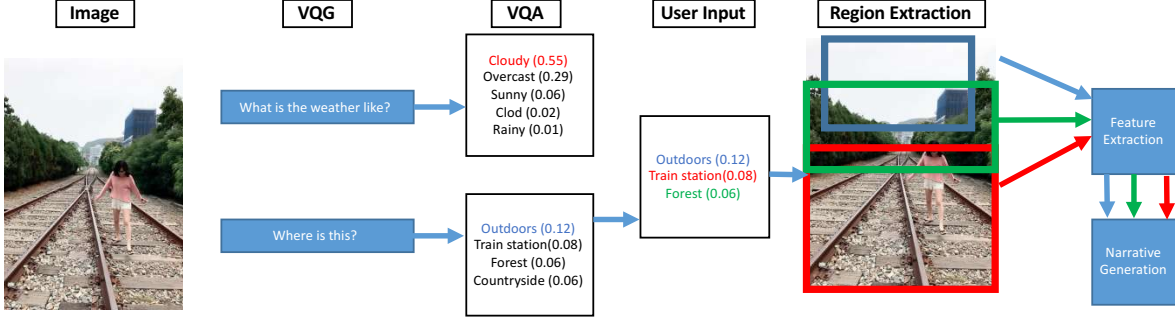
Figure 4: Questions that allow for multiple responses are generated to reflect user's interest and corresponding regions proceed to image narrative generation process.

answers, we internally utilize the VQA module. The question generated by the VQG module is passed on to VQA module, where the probability distribution $p_{ans}$ for all candidate answers $C$ is determined. If the most likely candidate $c_i = \max p_{ans}$, where $c_i \in C$, has a probability of being answer over a certain threshold $\alpha$, then the question is considered to have a single obvious answer, and is thus considered ineligible. The next question generated by VQG is passed on to VQA to repeat the same process until the the following requirement is met:

$$c_i < \alpha, c_i = \max p_{ans} \qquad (3)$$

In our experiments, we set $\alpha$ as 0.33. We also excluded the yes/no type of questions. Figure 4 illustrates an example of a question where the most likely answer had a probability distribution over the threshold (and is thus ineligible), and another question whose probability distribution over the candidate answers was more evenly distributed (and thus proceeds to narrative generation stage).

### 4.1.2 Region Extraction

Once the visual question that allows for multiple responses is generated, a user inputs his answer to the question, which is assumed to reflect his interest. We then need to extract a region within the image that corresponds to the user's response. We slightly modify the attention networks introduced in [30] in order to obtain the coordinates of the region that correspond to the user response. In [30], the question itself was fed into the network, so that the region necessary to answer that question is "*attended to.*" On the other hand, we are already given the answer to the question by the user. We take advantage of this by making simple yet efficient modification, in which we replace the *wh-* question terms with the response provided by the user. For example, a question "what is on the table?" with a user response "pizza" will be converted to a phrase "pizza is on the table," which is fed into attention network. This is similar to the rule-based NLP conversion in Section 3. We obtain the co-ordinates of the region from the second attention layer, by obtaining minimum and maximum values for *x*-axis and *y*-

axis in which the attention layer reacts to the input phrase. Since the regions are likely to contain the objects of interest at very tight scale, we extracted the regions at slightly larger sizes than coordinates. A region $r_i$ of size $(w_{r_i}, h_{r_i})$ with coordinates $x_{0_i}, y_{0_i}, x_{max_i}, y_{max_i}$ for image $I$ of size $(W, H)$ is extracted with a magnifying factor $\alpha$ (set as 0.25):

$$r_i' = (\max(0, x_{0_i} - w_{r_i}\alpha), \max(0, y_{0_i} - h_{r_i}\alpha),$$
$$\min(W, x_{max_i} + w_{r_i}\alpha), \min(H, y_{max_i} + h_{r_i}\alpha)) \qquad (4)$$

Given the region and its features, we can now apply the image narrative generation process described in Section 3 with minor modifications in setting. Regions are further extracted, visual questions are generated and answered, and rule-based natural language processing techniques are applied to organize them. Figure 4 shows an overall workflow of our model.

### 4.2. Applying Interaction to New Images

We represent each instance of image, question, and user choice as a triplet consisting of image feature, question feature, and the label vector for the user's answer. In addition, collecting multiple choices from identical users enables us to represent any two instances by the same user as a pair of triplets, assuming source-target relation. With these pairs of triplets, we can train the system to predict a user's choice on a new image and a new question, given the same user's choice on the previous image and its associated question. User's choice $x_{ans_i}$ is represented as one-hot vector where the size of the vector is equal to the number of possible choices. We refer to the fused feature representation of this triplet consisting of image, question, and the user's choice as **choice vector**.

We now project the image feature $x_{img_j}$ and question feature $x_{q_j}$ for the second triplet onto the same embedding space as the choice vector. We can now train a softmax classification task in which the feature from the common embedding space predicts the user's choice $x_{ans_j}$ on new question. In short, we postulate that the answer with index $u$, which maximizes the probability calculated by LSTM, is to be chosen as $x_{ans_l}$ by the user who chose $x_{ans_k}$, upon
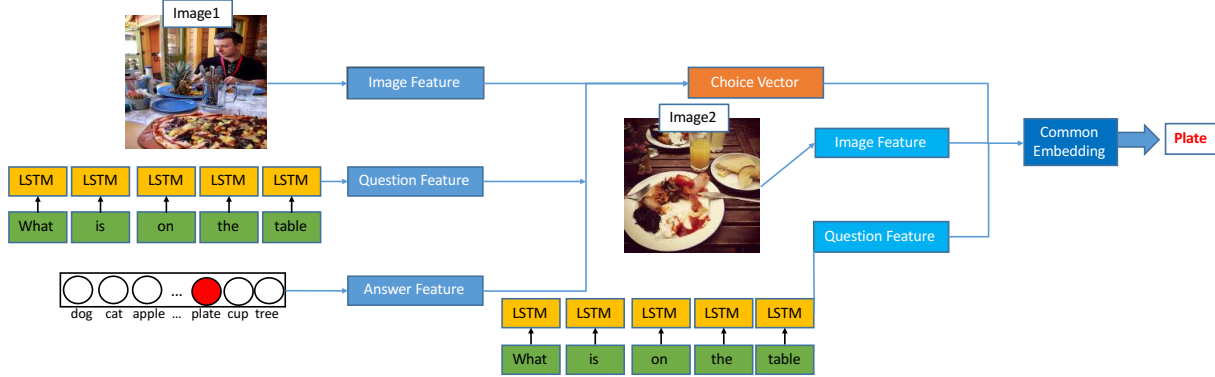
Figure 5: Training with pair of choices made by the same user. Given the choice vector for image 1 and new image feature and question feature for image 2, it is trained to predict the answer for the question on image 2.

seeing a tuple $(x_{img_l}, x_{q_l})$ of new image and new question:

$$u = \arg\max_v P(v; c_k, x_{img_l}, x_{q_l}) \qquad (5)$$

where $P$ is a probability distribution determined by softmax over the space of possible choices, and $c_k$ is the choice vector corresponding to $(x_{img_k}, x_{q_k}, x_{ans_k})$. This overall procedure and structure are essentially identical as in VQA task, except we augment the feature space to include choice vector. Figure 5 shows the overall workflow for training.

## 5. Experiments

### 5.1. Automatic Image Narrative Generation

#### 5.1.1 Setting

We applied the model described in Section 3 to 40,775 images in test 2014 split of MS COCO [18]. We compare our proposed model to three baselines as following:

**Baseline 1 (COCO):** general captioning trained on MS COCO applied to both images in their entireties and the region proposals

**Baseline 2 (SIND):** captions with model trained on MS SIND dataset [12], applied to both images in their entireties and the region proposals

**Baseline 3 (DenseCap):** captions generated by Dense-Cap [13] at both the whole images and regions with top 5 scores using their own region extraction implementation.

#### 5.1.2 Evaluation

**Automatic Evaluation**: It is naturally of our interest how humans would actually write image narratives. Not only can we perform automatic evaluation for reference, but we can also have a comprehension of what characteristics would be shown in actual human-written image narratives. We collected image narratives for a subset of MS COCO dataset [1]. We asked the workers to write a 5-sentence narrative about

---
[1] http://www.mi.t.u-tokyo.ac.jp/projects/narrative

Table 2: Examples of human-written image narratives collected on Amazon Mechanical Turk.

| Image | Human-written Narrative |
|---|---|
|  | This cat is having fun. She is very confused about the change in carpet. It is funny that this has interested her so much. Cats are very picky and they do not like changes. She is probably mad about this. |
|  | The pizza cook makes the pizza. The couple looks forward to pizza. The oven is very hot. He is a master at making pizza. He was born in italy. |

Table 3: Performances of the generated image narratives with human-written image narratives as ground truth.

| Model | BLEU1 | BLEU2 | BLEU3 | BLEU4 |
|---|---|---|---|---|
| COCO | 13.97 | 6.13 | 2.85 | 1.39 |
| SIND | 13.39 | 2.99 | 0.82 | 0.18 |
| DenseCap | 20.77 | 9.26 | 4.15 | 1.90 |
| Ours | **20.87** | 8.71 | 3.58 | 1.41 |

the image in a story-like way. We made it clear that the description can involve not only factual description of the main event, but also local elements, sentiments, inference, imagination, etc., provided that it can relate to the visual elements shown in the image. Table 2 shows examples of actual human-written image narratives collected and they display a number of intriguing remarks. On top of the elements and styles we asked for, the participants actively employed many other elements encompassing humor, question, suggestion, etc. in a highly creative way. It is also clear that conventional captioning alone will not be able to capture or mimic the semantic diversity present in them.

We performed automatic evaluation with BLEU [22] with collected image narratives as ground truth annotations.

Table 4: Each model's performance on DIANE.

| Metric | COCO | SIND | DenseCap | Ours |
|---|---|---|---|---|
| Diversity | 2.972 | 2.060 | 3.102 | **3.580** |
| Interesting | 2.875 | 2.100 | 3.336 | **3.489** |
| Accuracy | 2.812 | 2.105 | **3.188** | 3.132 |
| Naturalness | 2.754 | 2.059 | 3.146 | **3.374** |
| Expressivity | 2.819 | 2.141 | 3.257 | **3.381** |
| Overall | 2.846 | 2.093 | 3.201 | **3.391** |
| % of Win. | .300 | .195 | .357 | **.400** |

Table 5: Against each model on $\chi^2$ with 2 degrees of freedom, and one-sided $p$-value from binomial probability.

| vs. Model | > | = | < | $\chi^2$ | $p$-value |
|---|---|---|---|---|---|
| COCO | 2,208 | 1,222 | 1,570 | 133.37 | 1.4e-25 |
| SIND | 2,970 | 538 | 1,492 | 812.93 | 1.1e-11 |
| DenseCap | 1,890 | 1,454 | 1,656 | 271.33 | 4.5e-05 |

Table 3 shows the results. While resemblance to human-written image narratives may not necessarily guarantee better qualities, our model, along with DenseCap, showed highest resemblance to human-written image narratives. As we will see in human evaluation, such tendency turns out to be consistent, suggesting that resemblance to human-written image narratives may indeed provide a meaningful reference.

**Human Evaluation:** We asked the workers to rate each model's narrative with 5 metrics that we find essential in evaluating narratives; *Diversity*, *Interestingness*, *Accuracy*, *Naturalness*, and *Expressivity* (DIANE). Evaluation was performed for 5,000 images with 2 workers per image, and all metrics were rated in the scale of 1 to 5 with 5 being the best performance in each metric. We asked each worker to rate all 4 models for the image on all metrics.

Table 6 shows example narratives from each model. Table 4 shows the performance of each model on the evaluation metrics, along with the percentage of each model receiving the highest score for a given image, including par with other models. Our model obtained the highest score on *Diversity*, *Interestingness* and *Expressivity*, along with the highest overall score and the highest percentage of receiving best scores. In all other metrics, our model was the second highest, closely trailing the models with highest scores. Table 5 shows our model's performance against each baseline model, in terms of the counts of wins, losses, and pars. $\chi^2$ values on 2 degrees of freedom are evaluated against the null hypothesis that all models are equally preferred. The rightmost column in Table 5 corresponds to the one-sided $p$-values obtained from binomial probability against the same null hypothesis. Both significance tests provide an evidence that our model is clearly preferred over others.

**Discussion:** General image captioning trained on MS COCO shows weaknesses in accuracy and expressivity. Lower score in accuracy is presumably due to quick di-

version from the image contents as it generates captions directly from regions. Since it is restricted by an objective of describing the entire image, it frequently generates irrelevant description on images whose characteristics differ from typical COCO images, such as regions within an image as in our case. Story-like captioning trained on MS SIND obtained the lowest scores in all metrics. In fact, examples in Table 6 display that the narratives from this model are almost completely irrelevant to the corresponding images, since the correlation between single particular image and assigned caption is very low. DenseCap turns out to be the most competitive among the baseline models. It demonstrates the highest accuracy among all models, but shows weaknesses in interestingness and expressivity, due to their invariant tone and design objective of factual description. Our model, highly ranked in all metrics, demonstrates superiority in many indispensable aspects of narrative, while not sacrificing the descriptive accuracy.

## 5.2. Interactive Image Narrative Generation

### 5.2.1 Setting

We first need to obtain data that reflect personal tendencies of different users. Thus, we not only need to collect data from multiple users so that individual differences exist, but also to collect multiple responses from each user so that individual tendency of each user can be learned.

We generated 10,000 questions that allow for multiple responses following the procedure described in Section 4. We grouped every 10 questions into one task, and allowed 3 workers per task so that up to 3,000 workers can participate. Since multiple people are participating for the same group of images, we end up obtaining different sets of responses that reflect each individual's tendency.

We have permutation of 10 choose 2, $P(10, 2) = 90$ pairs of triplets for each user, adding up to 270,000 pairs of training data. Note that we are assuming a source-to-target relation within the pair, so the order within the pair does matter. We randomly split these data into 250,000 and 20,000 for training and validation splits, and performed 5-fold validation with training procedure described in Section 4. With 705 labels as possible choices, we had an average of 68.72 accuracy in predicting the choice on new image, given the previous choice by the same user. Randomly matching the pairs with choices from different users seemingly drops the average score down to 45.17, confirming that the consistency in user choices is a key point in learning preference.

### 5.2.2 Evaluation

**Question Generation:** For question generation, our interest is whether our model can generate questions that allow for various responses, rather than single fixed response. We

Table 6: Examples of image narratives. See Supplemental Material for many more examples.

| Image | COCO | SIND | DenseCap | Ours |
|-------|------|------|----------|------|
|  | An elephant standing in a field of grass. A large elephant standing in a field of grass. A bathroom with a sink and a mirror. A large building with a clock on it. | The dog was very happy to see the animals. We had a great time. I went to the museum today. We went to the city to see the sights. We saw a lot of old buildings. | An elephant standing in a field of grass. A gray elephant. Elephant trunk is curled. Elephant in the photo. Trunk of an elephant. Elephants walking on the road. | This is a baby elephant. The elephants are standing on grass. They are bored. The elephant is sitting. This is not a zoo. |
|  | A teddy bear sitting on top of a wooden table. A teddy bear sitting on top of a bed. A stuffed bear is sitting on a bed. | The cake was delicious. I had a great time. The food was delicious. | A teddy bear sitting on top of a wooden table. Teddy bear on a table. A brown teddy bear. A teddy bear. A teddy bear on a table. The head of a teddy bear. | Bear is on the cake. That stuffed animal has a funny face. It is mine. The bear is wearing hat. The cake is white. |

Table 7: Evaluation results on whether the generated questions allow for multiple responses.

| Model | # Overall | # Yes | # No |
|-------|-----------|-------|------|
| Ours | 1,000 | **664** | 336 |
| VQG | 1,000 | 217 | 783 |
| Overall | 2,000 | 881 | 1119 |

Table 8: Examples of generated questions using our proposed model and VQG respectively.

| Image | Generated Questions |
|-------|---------------------|
|  | **Ours** |
| | What is the color of the shirt? |
| | **VQG** |
| | How many children are there? |
|  | **Ours** |
| | What is on the table? |
| | **VQG** |
| | What is the table made of? |

asked the workers on Amazon Mechanical Turk to decide whether the question can be answered in various ways or has multiple answers, given an image. 1,000 questions were generated with our proposed model using both VQG and VQA, and another 1,000 questions were generated using VQG only.

Table 7 shows the number of votes for each model. It is very clear that the questions generated from our proposed model of parallel VQG and VQA outperformed by far the questions generated from VQG only. This is inevitable in a sense that VQG module was trained with human-written questions that were intended to train the VQA module, i.e. with questions that mostly have clear answers. On the other hand, our model deliberately chose the questions from VQG that have evenly distributed probabilities for answer labels, thus permitting multiple possible responses. Table 20 shows examples of visual questions generated from our

model and VQG only respectively. In questions generated from our model, different responses are possible, whereas the questions generated from VQG only are restricted to single obvious answer.

**Reflection of User's Choice on the Same Image:** Our next experiment is on the user-dependent image narrative generation. We presented the workers with 3,000 images and associated questions, with 3 possible choices as a response to each question. Each worker freely chooses one of the choices, and is asked to rate the image narrative that corresponds to the answer they chose, considering how well it reflects their answer choices. As a baseline model, we examined a model where the question is absent in the learning and representation, so that only the image and the user input are provided. Rating was performed over scale of 1 to 5, with 5 indicating highly reflective of their choice. Table 11 shows the result. Agreement score among the workers was calculated based on [3]. Agreement score for our model falls into the range of 'moderate' agreement, whereas, for baseline model, it is at the lower range of 'fair' agreement, as defined by [16], demonstrating that the users more frequently agreed upon the reliability of the image narratives for our model. Our model clearly has an advantage over using image features only with a margin considerably over standard deviation. Table 26 shows examples of images, generated question, and image narratives generated depending on the choice made for the question respectively.

**Reflection of User's Choice on New Images:** Finally, we experiment with applying user's interest to new images. As in the previous experiment, each worker is presented with an image and a question, with 3 possible choices as an answer to the question. After they choose an answer, they are presented with a new image and a new image narrative. Their task is to determine whether the newly presented image narrative reflects their choice and interest. As a base-

Table 9: Examples of image narratives generated depending on the user choices.

| Image | Answers, Regions and Narratives | | |
|---|---|---|---|
| | Skateboard | Motorcycle | Car |
| |  |  |  |
| **Generated Question** | The man is riding skateboard. | The man is riding motorcycle. | The man is riding car. |
| What is the man riding? | The man is skateboarding. | It is white. | This is a modern car. |
| | The color of the jacket is red. | The motorcycle is honda. | It is a black and white photo. |

Table 10: Examples of image narratives generated on new images, depending on the choices made.

| Image & Question | Choice | New Image | Image Narrative |
|---|---|---|---|
|  What animal is this? | giraffe |  | The giraffe is standing. The weather is sunny. |
| | zebra | | Zebra is thinking. It is not in a zoo. |
| | rhino | | 2 animals are in the picture. The sky is blue. |
|  What kind of animal is that? | dog |  | The horse is running. The car is white. |
| | sheep | | The boy is wearing red shirt. Tree is in the background. |
| | person | | The man is riding horse. The man is wearing hat. |

Table 11: Evaluation results on how well the generated image narrative reflects the choices they made.

| Model | Avg. Score | Agreement |
|---|---|---|
| Ours | **3.851±1.12** | **.601** |
| image only | 2.636±1.01 | .432 |

Table 12: Evaluation results on how well the generated image narrative for new images reflects their interest.

| Model | Avg. Score | Agreement |
|---|---|---|
| Ours | **3.455±0.93** | **.527** |
| random match | 2.772±0.79 | .489 |
| image only | 2.238±1.24 | .428 |

line, we again examined a model where the question is absent in the learning and representation stages. In addition, we performed an experiment in which we trained preference learning module with randomly matched choices. This allows us to examine whether there exists a consistency in user choices that enables us to apply the learned preferences to new image narratives.

Table 12 shows the result. As in previous experiment, our model clearly has an advantage over using image features only. Inter-rater agreement score is also more stable for our model. Training preference learning module with randomly matched pairs of choices resulted in a score below our proposed model, but above using the image features only. This may imply that, even with randomly matched

pairs, it is better to train with actual choices made by the users with regards to specific questions, rather than with conspicuous objects only. Overall, the result confirms that it is highly important to provide a context, in our case by generating visual questions, for the system to learn and reflect the user's specific preferences. It also shows that it is important to train with consistent choices made by identical users. Table 27 shows examples of image narratives generated for new images, depending on the choice the users made for the original image, given the respective questions.

## 6. Conclusion

We proposed a customized image narrative generation task, where we proposed a model to engage the users in image description generation task, by directly asking questions to the users, and collecting answers. Experimental results demonstrate that our model can successfully diversify the image description by reflecting the user's choice, and that user's interest learned can be further applied to new images.

## Acknowledgments

# References

[1] A. Agrawal, D. Batra, and D. Parikh. Analyzing the behavior of visual question answering models. In *EMNLP*, 2016.

[2] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. Lawrence Zitnick, and D. Parikh. VQA: Visual Question Answering. In *ICCV*, 2015.

[3] E. Bennett, R. Alpert, and A. Goldstien. Communications through limited-response questioning. *Public Opinion Quarterly*, 18:303–308, 1954.

[4] P. Christiano, J. Leike, T. B. Brown, M. Martic, S. Legg, and D. Amodei. Deep reinforcement learning from human preferences. *https://arxiv.org/abs/1706.03741*, 2017.

[5] A. Das, S. Kottur, A. Singh, D. Yadav, J. M. Moura, D. Parikh, and D. Batra. Visual dialog. In *CVPR*, 2017.

[6] H. de Vries, F. Strub, S. Chandar, O. Pietquin, H. Larochelle, and A. C. Courville. Guesswhat?! visual object discocery through multi-modal dialogue. In *CVPR*, 2017.

[7] H. Fang, S. Gupta, F. N. Iandola, R. K. Srivastava, L. Deng, P. Dollár, J. Gao, X. He, M. Mitchell, J. C. Platt, C. L. Zitnick, and G. Zweig. From captions to visual concepts and back. In *CVPR*, 2015.

[8] A. Fukui, D. H. Park, D. Yang, A. Rohrbach, T. Darrell, and M. Rohrbach. Emnlp. 2016.

[9] A. Ghodrati, A. Diba, M. Pedersoli, T. Tuytelaars, and L. J. V. Gool. Deepproposal: Hunting objects by cascading deep convolutional layers. In *ICCV*, 2015.

[10] R. Girshick. Fast r-cnn. In *ICCV*, 2015.

[11] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.

[12] T.-H. Huang, F. Ferraro, N. Mostafazadeh, I. Misra, A. Agrawal, J. Devlin, R. Girshick, X. He, P. Kohli, D. Batra, C. L. Zitnick, D. Parikh, L. Vanderwende, M. Galley, and M. Mitchell. Visual storytelling. In *NAACL*, 2016.

[13] J. Johnson, A. Karpathy, and L. Fei-Fei. DenseCap: Fully Convolutional Localization Networks for Dense Captioning. In *CVPR*, 2016.

[14] K. Kafle and C. Kanan. Visual question answering: Datasets, algorithms, and future challenges. *https://arxiv.org/abs/1610.1465*, 2016.

[15] A. Karpathy and F.-F. Li. Deep Visual-Semantic Alignments for Generating Image Descriptions. In *CVPR*, 2015.

[16] J. R. Landis and G. G. Koch. The measurement of observer agreement for categorical data. *Biometrics*, 33(1), 1977.

[17] P. D. Larry Zitnick. Edge boxes: Locating object proposals from edges. In *ECCV*, 2014.

[18] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollr, and C. L. Zitnick. Microsoft COCO: Common Objects in Context. In *ECCV*, 2014.

[19] J. Lu, J. Yang, D. Batra, and D. Parikh. Hierarchical question-image co-attention for visual question answering. In *NIPS*, 2016.

[20] M. Marcus, G. Kim, M. A. Marcinkiewicz, R. MacIntyre, A. Bies, M. Ferguson, K. Katz, and B. Schasberger. The penn treebank: annotating predicate rgument structure. In *HLT*, 1994.

[21] N. Mostafazadeh, I. Misra, J. Devlin, M. Mitchell, X. He, and L. Vanderwende. Generating natural questions about an image. In *ACL*, 2016.

[22] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. Bleu: A method for automatic evaluation of machine translation. In *ACL*, 2002.

[23] M. Ren, R. Kiros, and R. Zemel. Exploring Models and Data for Image Question Answering. In *NIPS*, 2015.

[24] K. Shih, S. Singh, and D. Hoiem. Where to look: Focus regions for visual question answering. In *CVPR*, 2016.

[25] A. Shin, Y. Ushiku, and T. Harada. The color of the cat is gray: 1 million full-sentences visual question answering (fsvqa). *arXiv:1609.6657*, 2016.

[26] J. R. Uijlings, K. E. van de Sande, T. Gevers, and A. W. Smeulders. Selective search for object recognition. *IJCV*, 104:154–171, 2013.

[27] R. Vedantam, C. L. Zitnick, and D. Parikh. Cider: Consensus-based image description evaluation. In *CVPR*, 2015.

[28] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and Tell: A Neural Image Caption Generator. In *CVPR*, 2015.

[29] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio. Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. In *ICML*, 2015.

[30] Z. Yang, X. He, J. Gao, L. Deng, and A. Smola. Stacked attention networks for image question answering. In *CVPR*, 2016.

Table 13: Examples of captions and questions for the same image. While captions essentially describe the same contents, questions widely vary in terms of the topics.

| Image |
|:---:|
|  |
| **COCO Captions** |
| • A group of people sitting on the back |
| • Several people are taking a ride on elephants |
| • Some people are riding elephants in the jungle |
| • The people are riding on the two elephants |
| • People riding on elephants in the jungle |
| **VQA Questions** |
| • Is this standard transportation in the United States? |
| • Are they on a paved roadway? |
| • How many people are riding elephants? |



Figure 6: Viewer's attention varies depending on the context provided.

# Appendices

## A. Why generate quesetions?

A question may arise as to why not to simply ask the users to select the region or part of the image that stands out the most to them. In such case, there would be no need to *generate* the questions for each image, as the question '*what stands out the most?*' would suffice for all images. This, however, would be equivalent to a simple saliency annotation task, and would not allow for any meaningful customization or optimization per user. Thus, as discussed above, generating a question for each image is intended to

provide a context in which each user can apply their own specific interest. Figure 6 shows how providing context via questions can diversify people's attention. Apart from simply generating diverse image narratives based on the user input, many potential applications can be conceived of. For example, in cases where thorough description of an entire scene results in a redundant amount of information both quality and quantity-wise, application of our model can be applied to describe just the aspect that meets the user's interest that was learned.

Table 14: Statistics from the crowd-sourcing task on collecting answers to non-visual questions.

| # of answers collected | 48,090 |
|:---|:---:|
| # of unique answers | 15,469 |
| # of workers participated | 187 |
| max. # assignments by worker | 1609 |
| avg. # assignments per worker | 51.43 |
| rewards per assignment | $.10 |
| 10 most common answers | '*yes*','*no*','*tom*', '*london*','*mine*', '*downtown*','*john*','*me*' '*halloween*','*new york*' |

Table 15: Examples of answers collected on VQG.

| Question | Answer |
|:---:|:---:|
| '*What is the name of the man?*' | '*Tom*' |
| '*What is the score in the game?*' | '*0-0*' |
| '*What kind of record is being played?*' | '*rap records*' |
| '*How long until the bathroom is fixed?*' | '*1 week*' |
| '*Why is he making weird face?*' | '*he's drunk*' |
| '*What's the cat's name?*' | '*Moni*' |
| '*How much did that cost?*' | '*10 dollars*' |
| '*What destroyed this town?*' | '*bomb*' |
| '*Why are the trees lit up?*' | '*It's Christmas time*' |
| '*What are the ingredients?*' | '*fish,bread,broccoli*' |

## B. Clarification of DIANE

Few works tackled the task of narrative evaluation, hardly taking visual information into consideration. Although we could not find an authoritative work on the topic of narrative evaluation, this was our best attempt at not only reflecting precision/recall, but various aspects contributing to the integrity of the image narrative. *Diversity* deals with the coverage of diction and contents in the narrative, roughly corresponding to *recall*. *Interestingness* measures the extent to which the contents of the narrative grasp the user's attention. *Accuracy* measures the degree to which the description is relevant to the image, corresponding to *pre-*
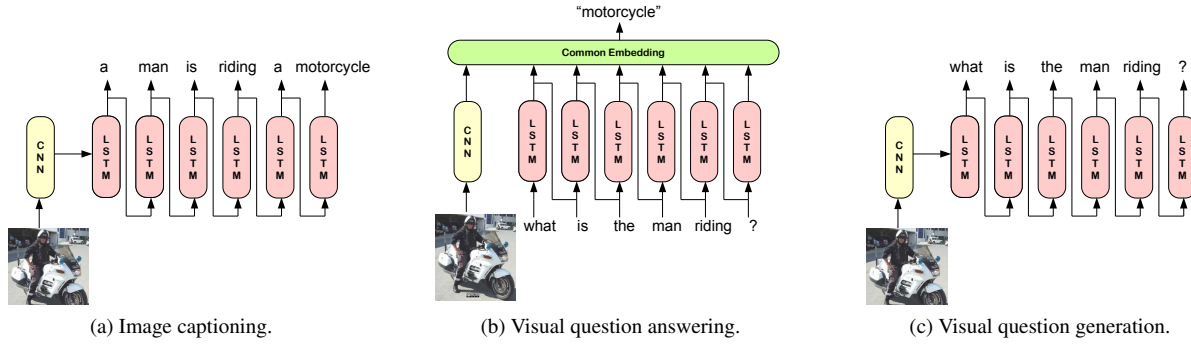
(a) Image captioning.

(b) Visual question answering.

(c) Visual question generation.

Figure 7: Illustration of the overall workflow for each task.

Table 16: Examples of questions generated using non-visual questions in VQG dataset.

| Image | Generated Questions |
|---|---|
|  | • Is this a hotel room?<br>• Is that a picture of your house?<br>• Where did you get the pillows?<br>• Is this new tile?<br>• Was that clean there?<br>• How big is that room? |
|  | • Is the woman drunk?<br>• Is this a church?<br>• Is this structure in a museum?<br>• What city was this in?<br>• Are they protesting? |
|  | • What kind of pizza is that?<br>• Is it for dinner?<br>• What kind of topping is this on the pizza?<br>• What does the plate say? |
|  | • What is this bird staring at?<br>• How long will it be there?<br>• Is that a real bird?<br>• What sort of bird is that?<br>• What kind of flower is that? |

Table 17: Examples of human-written image narratives collected on Amazon Mechanical Turk.

| Image | Human-written Narrative |
|---|---|
|  | The food truck looks good.<br>I bet they have good food.<br>Does everyone in a food truck have a beard?<br>I am so done with the beard thing.<br>Hope his beard does not get into the food. |
|  | Car is in very good shape for the age.<br>This is a prefect car for California.<br>I think i see that this is in Huntington beach.<br>This would attract a lot of attention.<br>Great way to pick up girls or guys. |
|  | A dad and his daughter are sitting on the couch.<br>They have just woken up.<br>They each have a cup of juice.<br>They use cups with lids so they don't spill on the couch. |
|  | Tom is playing with a frisbee.<br>He is practicing new moves.<br>He jumped up in the air.<br>He is trying to catch it between his legs.<br>He was successful in his attempt. |

*cision*. Contents that are not visually verifiable are considered accurate only if they are compatible with salient parts of the image. *Naturalness* refers to the narrative's overall resemblance to human-written text or human-spoken dialogue. *Expressivity* deals with the range of syntax and tones in the narrative.

# C. Additional Experiments

We also performed an experiment in which we generate image narratives by following conventional image captioning procedure with human-written image narratives collected on Amazon Mechanical Turk. In other words, we trained LSTM with CNN features of images and human-written image narratives as ground truth captions. If such setting turns out to be successful, our model would not have

Table 18: Statistics for human-written image narratives collected on Amazon Mechanical Turk.

| | |
|---|---|
| # of answers collected | 13,221 |
| rewards per assignment | $.20 |
| minimum length of image narrative | 10 |
| maximum length of image narrative | 83 |
| average length of image narrative | 31.629 |

Table 19: Examples of image narratives generated by training with human-written image narratives.

| Image | Generated Narrative |
|---|---|
|  | a man is sitting on a chair he is wearing a white shirt he seems to |
|  | a man is holding a hot dog he is wearing a white shirt he seems to |

Table 20: Examples of generated questions for user interaction using our proposed model and VQG only respectively.

| Image | Generated Questions |
|---|---|
|  | **Ours** |
| | What is the dog doing? |
| | **VQG** |
| | What is the color of the couch? |
|  | **Ours** |
| | What is the color of the car? |
| | **VQG** |
| | What is the weather like? |

much comparative merit.

We trained an LSTM with collected image-narratives for training split of MS COCO. We retained the experimental conditions identically as previous experiments, and trained for 50 epochs. Table 19 shows example narratives generated. Not only does it utterly fail to learn the structure of image narratives, but it hardly generates text over one sentence, and even so, its descriptive accuracy is very poor. Since LSTM now has to adjust its memory cells' dependency on much longer text, it struggles to even form a complete sentence, not to mention inaccurate description. This tells us that simply training with human-written image narratives does not result in reliable outcomes.

With reference human-written image narratives, we further performed CIDEr [27] evaluation as shown in Table 25.

# D. Discussion

It was shown via the experiments above that there exists a certain consistency over the choices made by the same user, and that it is thus beneficial to train with the choices made by the same users. Yet, we also need to investigate whether such consistency exists across different categories of images. We ran Fast-RCNN [10] on the images used in our experiment, and assigned the classes with probability over 0.7 as the labels for each image. We then define any two images to be in the same category if any of the assigned labels overlaps. Of 3,000 pairs of images used in the experiment, 952 pairs had images with at least one label overlapping. Our proposed model had average human evaluation score of 4.35 for pairs with overlapping labels and 2.98 for pairs without overlapping labels. Baseline model with image features only had 2.57 for pairs with overlapping labels and 2.10 for pairs without overlapping labels. Thus, it is shown that a large portion of the superior performance of our model comes from the user's consistency for the images of the same category, which is an intuitively correct conclusion.

However, our model also has superiority over baseline model for pairs without overlapping labels. This may seem more difficult to explain intuitively, as it is hard to see any explicit correlation between, for example, a car and an apple, other than saying that it is somebody's preference. We manually examined a set of such examples, and frequently found a pattern in which the color of the objects of choices was identical; for example, a red car and an apple. It is difficult to attribute it to a specific cause, but it is likely that there exists some degree of consistency in user choices over different categories, although to a lesser extent than for images in the same category. Also, it is once again confirmed that it is better to train with actual user choices made on specific questions, rather than simply with most conspicuous objects.

# E. Additional Figures & Tables

Table 13 shows the contrast between semantic diversity of captions and questions. Figure 7 shows overall architecture each of image captioning, visual question answering, and visual question generation task. Table 14 shows statistics for crowd-sourcing task on collecting answers to non-visual questions in VQG dataset. Table 15 shows examples of answers to VQG questions collected on crowd-sourcing. Table 16 shows examples of generated questions using VQG dataset. Table 17 shows examples of human-written image narratives. Table 18 shows statistics for human-written image narratives collection. Table 21 shows conversion rules for natural language processing stage for narrative generation process as used in Section 3. Table 22 to Table 24 show more examples of image narratives. Ta-

Table 21: Conversion rules for transforming question and answer pairs to declarative sentences.

| Type | Rule (Q→A) | Question | Ans. | Converted Ans. |
|---|---|---|---|---|
| yes/no | VB1+NP+VB2/JJ? <br> →NP+*conjug* (VB2/JJ,*tense*(VB1)) | – <br> - <br> *Did he get hurt?* | – <br> - <br> *yes* | – <br> <br> *He got hurt.* |
| | **or**, NP +*negate*(*conjug* (VB2/JJ,*tense*(VB1))) | - <br> - <br> *Is she happy?* | <br> - <br> *no* | <br> <br> *She is not happy.* |
| | MD+ NP+VB? <br> →NP+MD+VB **or**, <br> NP+*negate*(MD)+VB | – <br> *Will the boy fall asleep?* <br> *May he cross the road?* | – <br> *yes* <br> *no* | – <br> *The boy will fall asleep.* <br> *He may not cross the road.* |
| number | "*How many*"+NP+ /*is/are*+EX? <br> →EX+*is/are*+*ans*+NP | - <br> - <br> *How many pens are there?* | - <br> - <br> *2* | <br> <br> *There are 2 pens.* |
| | "*How many*"+NP1(+MD) +VB(+NP2)? <br> →*ans*(+MD)+VB(+NP2) | - <br> – <br> *How many people are walking?* | - <br> – <br> *3* | <br> – <br> *3 people are walking.* |
| | "*How many*"+NP1+ VB1/MD+NP2+VB2? <br> →NP2 +(MD+VB2)/*conjug* (VB2,*tense*(VB1)) +*ans*+NP1 | - <br> – <br> - <br> - <br> - <br> *How many pens does he have?* | - <br> – <br> - <br> - <br> - <br> *4* | <br> – <br> <br> <br> <br> *He has 4 pens.* |
| others | WP/WRB/WDT+ "*is/are*"+NP? <br> → NP+"*is/are*"+*ans*. | - <br> - <br> *Who are they?* | - <br> - <br> *students* | <br> <br> *They are students.* |
| | WP+NP+VP? → *ans*.+VP | *What food is on the table?* | *apple* | *Apple is on the table.* |
| | WDT+NP+VP(+NP2)? <br> →*ans*.(+NP)+VP(+NP2) | - <br> *Which hand is holding it?* | - <br> *left* | <br> *Left hand is holding it.* |
| | WP/WDT+MD+VB? <br> →*ans*.+MD+VB | - <br> *Who would like this?* | - <br> *dog* | <br> *Dog would like this.* |
| | WP/WDT+MD+NP+VB? <br> →NP+MD+VB+*ans*. | - <br> *What would the man eat?* | - <br> *apple* | <br> *The man would eat apple.* |
| | WP/WDT+VP(+NP)? <br> →*ans*.+VP(+NP) | - <br> *Who threw the ball?* | - <br> *pitcher* | <br> *Pitcher threw the ball.* |
| | WP/WDT+VB1+NP+VB2? <br> →NP+*conjug* (VB2,*tense*(VB1))+*ans*. | – <br> - <br> *What is the man eating?* | – <br> - <br> *apple* | – <br> <br> *The man is eating apple.* |

ble 20 shows examples of questions for user interaction that were generated using our proposed model of combining VQG and VQA, and the baseline of using VQG only. Table 26 shows another example of customized image narratives generated depending on the choices made by user upon the question. Table 27 shows examples of how the choices made by user upon the question were reflected in new images.

## F. Additional Clarifications

**Why were yes/no questions excluded?** Yes/no questions are less likely to induce multiple answers. The number of possible choices is limited to 2 in most cases, and rarely correspond well to particular regions.

**Failure cases for rule-based conversion:** Since both questions and answers are human-written, our conversion rule frequently fails with typos, abridgments, words with multiple POS tags, and grammatically incorrect questions. We either manually modified them or left them as they are.

| Image | COCO | SIND | DenseCap | Ours |
|---|---|---|---|---|
| | A teddy bear sitting on a wooden bench. A teddy bear sitting on top of a tree. A teddy bear is sitting on the ground. A train traveling down tracks next to a forest. A teddy bear is sitting on a tree branch. | The kids had a great time. The dog was very happy to see me. I had a great time. The view was amazing. We had a great time. | A teddy bear sitting on a wooden bench. A teddy bear in a red hat. Red teddy bear. A white teddy bear. The nose of a sheep. The teddy bear is s itting on the ground. | These are stuffed animals. That teddy bear can be scary when you see it at night. The animals are there for fun. The bear is sleeping. This is not a real bear. |
| | A man is eating a hot dog in a restaurant. A man holding a hot dog in his hand. A man holding a hot dog in a bun. A man in a suit and tie standing in front of a building. A man in a hat is holding a hot dog. | We had a great time. | A man is eating a hot dog in a restaurant. Woman holding a sandwich. Woman has brown hair. Woman in black jacket. A sandwich on a white plate. A brown wooden wall. | The girl is eating sandwich. Her name is mary. She is hungry. She eats a lot. She is smiling. |
| | A cat laying on top of a bed next to a remote control. A cat laying on top of a bed next to a laptop. A cat laying on top of a bed next to a window. A person is holding a piece of broccoli. A large building with a clock on it. | The dog was very happy to see me. I had a great time. We saw a lot of old buildings. The view from the top was amazing. | A cat laying on top of a bed next to a remote control. A cat laying on a bed. The head of a cat. Ear of a cat. The cat is brown. The ear of a cat. | There is 1 cat. This cat looks lonely. The cat is not sleeping. The weather is sunny. |
| | A man riding skis down a snow covered slope. A man in a red jacket is snowboarding. A man wearing a hat and a tie. A street light with a building in the background. A group of people standing on a beach with a kite. | We had a great time. I was so happy to see me. We went to the city to see the sights. I was so excited to see my friends. I went to the city last weekend. | A man riding skis down a snow covered slope. Man in red jacket. Snow covered mountain. The woman is wearing a helmet. The man is wearing black pants. Black and white jacket. | The person is skiing. This is alps. This person is having fun. He won the competition. The person is holding ski poles. |
| | A plate with a piece of cake on it. A close up of a pair of scissors on a table. A plate with a sandwich and a salad on it. A piece of cake on a plate with a fork. A close up of a plate of food on a table. | The food was delicious. The flowers were so beautiful. I had a great time. | A plate with a piece of cake on it. A yellow piece of donut. A red basket on the table. A box of donuts. A table with a wooden table. A donut with sprinkles. | There are 3 different types of food. This is not a healthy breakfast. This cake looks so fun. The orange color is bread. |

Table 22: More examples of image narratives.

**Experiments with different VQA models.** Most of well-known VQA models' performances are currently in a relatively tight range. In fact, we tried [8], SOTA at the time of experiment, but did not see any noticeable improvement.

**Is attention network retrained to handle sentences?** No, but we found that attention network trained for questions works surprisingly well for sentences, which makes sense since key words that provide attention-wise clue are
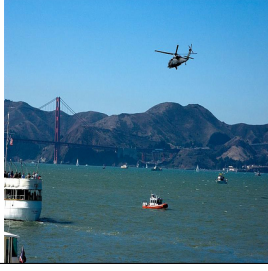
| Image | COCO | SIND | DenseCap | Ours |
|---|---|---|---|---|
|  | A baseball player swinging a bat at a ball. A man is standing in the grass with a kite. A bench sitting on top of a lush green hillside. A small boat in a body of water. A park bench with a tree in the background. | The man is giving a speech. We went to the park to see the sights. The view from the top was amazing. The building was very tall. We saw a lot of old buildings | A baseball player swinging a bat. Dirt on the ground. The jersey is blue. Player holding a bat. The helmet is black. Green grass on the field. | The boy is playing baseball. The score of the match is tied at zero. The name of that who is playing is john. The weather is like sunny. The bat is black. Trees are in the background. |
|  | A city street filled with lots of traffic. A red double decker bus driving down a street. A bus is parked on the side of the road. A bus that is driving down the street. A bathroom with a toilet and a sink. | We had a great time. The car was covered in snow. The inside of the building was very tall. The view from the top was amazing. | A traffic light. Trees with no leaves. A tree with no leaves. The front wheel of a bus. A street light. A traffic light on a pole. | That bus looks loud together. The traffic light is green. The man is walking. The color of the car is white. The weather is like cold. This is not a busy road. |
|  | A black bear walking across a lush green field. A polar bear walking in the snow near a tree. A black bear is walking in the water. A black and white photo of a bird on a tree branch. | The dog was so excited to see the water. The dog was very tired. The dog was very happy to see me. We had a great time. We saw a lot of old buildings. We went to the lake to see the sights. | Bear walking on the ground. A bear in the water. A large tree. A bear in the snow. Trees covered with snow. The bear is brown. | The bear is walking. This bear does not look silly. The bear is not eating. The bear is in the wild. It is cold outside. The weather is like cold. |
|  | A group of people flying kites on a beach. A large boat floating on top of a body of water. A bathroom with a sink and a mirror. A person on a snowboard in the snow. A plane flying in the sky over mountains. | The view from the top was amazing. We had a great time. The view was amazing. I went to the museum today. The beach was beautiful. | Plane in the sky. A blue sky with no clouds. Photo taken during the day. Red and white boat in water. The water is blue. The sky is blue. | The water is calm. This is daytona beach. The weather is like cloudy. The plane is flying. The water is blue. It is day. |
|  | A bathroom with a toilet sink and mirror. A bathroom with a toilet and a sink. A pair of scissors sitting on top of a wooden table. A street sign on a pole in front of a building. A pair of scissors sitting on top of a table. | I had a great time. I went to the park to see the sights. I went to the museum today. | A white bathroom sink. A silver metal towel rack. Silver metal faucet. A silver faucet. Toilet paper holder on wall. A white toilet paper. | This bathroom is in a hotel. The bathroom is clean. This room is very good. The color of the wall is white. There is a reflection in the mirror. The light is on. |

Table 23: More examples of image narratives.

likely limited, and hardly inquisitive words.

**Why not train with "I dont know?"** We were concerned that answers like "I don't know" would likely overfit.

It would also undermine creative aspect of image narrative, without adding much to functional aspect.

| Image | COCO | SIND | DenseCap | Ours |
|---|---|---|---|---|
| | A herd of zebras grazing in a field. A herd of zebra standing on top of a lush green field. A bird flying over a building with a clock. A man standing on a sidewalk next to a street sign. A group of zebras are standing in a field. | We had a great time. I went to the museum today. We saw a lot of interesting things. We went to the city to see the sights. We saw many different types of animals. We went to the museum. | A herd of zebras grazing in a field. A field of grass. Two zebras in a field. The photo was taken in the daytime. White clouds in blue sky. The grass is tall. | The zebras like each other. These animals are related. The zebras are not in a zoo. The animal is grazing. |
| | A close up of a pizza on a plate. A close up of a sandwich on a plate. A cat sitting on top of a window sill. A bathroom with a toilet and a sink. A plate of food with a sandwich and french fries. A person holding a hot dog in a bun. | The food was delicious. We had a great time. I went to the museum today. | A close up of a pizza on a plate. Pizza on a plate. Pizza on a table. The hand of a person. A cup of coffee. The pizza has red sauce. | 500 calories are in the meal. This is a pizza. This is not a healthy meal. This is not for vegetarian. |
| | A street with cars parked on the side of it. A car parked in front of a parking meter. A street sign on a pole on a street. A car parked on the side of a road. A street sign that is on a pole. | We went to the city to see the sights. The car was covered in snow. I went to the museum today. We went to the museum. We had a great time. We went to the location. | A street with cars parked on the side of it. A silver car parked on the street. A black car parked on the street. A white truck. Blue sky with no clouds. A black truck. | The car is gray. The car is parked illegally. Where the car is is inappropriate. That is pine tree behind. |
| | A plate of food with a fork and knife. A pizza with a lot of toppings on it. A plate with a sandwich and a salad. A close up of a plate of food with broccoli. | The food was delicious. | A plate of food with a fork and knife. Pizza on a table. A pizza on a plate. A slice of pizza. The pizza has red sauce. A slice of tomato. | This is a vegetarian pizza. This is not a cheese pizza. The green vegetable is spinach. This is a healthy meal. |
| | A dog that is sitting on a bench. A street sign on a pole in front of a building. A man riding a skateboard down a street. A dog is running with a frisbee in its mouth. A large building with a clock on it. | The dog was very happy to see me. We had a great time. The house was very nice. I went to the museum today. | A dog that is sitting on a bench. A brown dog. A brick sidewalk. Man walking on sidewalk. Dog walking on sidewalk. A white line on the ground. | There is a dog. The dog is sad. The color of the wall is white. The color of the fire hydrant is gray. |

Table 24: More examples of image narratives.

| Model | COCO | SIND | DenseCap | Ours |
|---|---|---|---|---|
| **CIDEr** | 18.0 | 9.9 | **28.0** | **27.7** |

Table 25: Each model's performance on CIDEr with human-written image narratives as ground truths.

Table 26: Examples of image narratives generated depending on the user choices.

| Image | Answers, Regions and Narratives | | |
|---|---|---|---|
| | Pizza | Pine Apple | Plate |
|  |  |  |  |
| **Generated Question** | Pizza is on the table. | Pine apple is on the table. | Plate is on the table. |
| What is on the table? | The man is eating pizza. | The man is vegetarian. | The man is eating. |
| | The pizza is thin crust. | | The man is eating more than one person would. |

Table 27: Examples of image narratives generated on new images, depending on the choices made.

| Image & Question | Choice | New Image | Image Narrative |
|---|---|---|---|
|  What is the man riding? | skateboard |  | No one is riding bicycle. The man is standing. |
| | motorcycle | | The motorcycle is red. No one is riding motorcycle. |
| | car | | This is not a modern building. The image is not in black and white. |
|  What color is the car? | white |  | The white object is bus. The car is white. |
| | green | | The bus is green. The train is headed to Washington. |
| | yellow | | The train is yellow. The image is not in black and white. |