

# Category Consistent Cyclic Visual Question Generation

Shagun Uppal<sup>\* 1</sup>, Anish Madan<sup>\* 1</sup>, Sarthak Bhagat<sup>\* 1</sup>, Yi  
Yu<sup>2</sup>, Rajiv Ratn Shah<sup>1</sup>

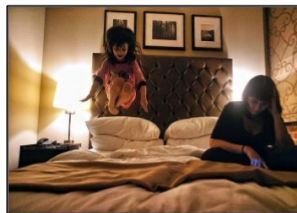
<sup>1</sup> IIT-Delhi, <sup>2</sup> NII Japan, <sup>\*</sup> Equal Contribution

May 24, 2020

**Visual Question Generation (VQG)** is the task of generating natural questions given an image.

Challenges in constructing a VQG system:

- ▶ Capturing various concepts in images.
- ▶ Relevance of generated questions to the image.
- ▶ Many-to-one mapping between the image and generated questions since multiple questions are possible for an image.
- ▶ Avoid questions which invoke generic answers like "yes" / "I don't know".



**Possible Category-Question pairs:**

**SPATIAL:** Where are the pictures hanging?

**ACTIVITY:** What is the little girl doing?

**BINARY:** Is the lamp on?

**COUNT:** How many pillows are there on the bed?

**COLOR:** What is the color of the girl's dress?

- ▶ Weaken supervision by removing the need for answers.
- ▶ Variational training using a single combined latent space for image and category by maximizing mutual information.
- ▶ Category consistency using cyclic training in two disjoint steps.
- ▶ Center loss for category-wise clustering.
- ▶ Hyper-prior on latent space for encapsulation of independent features.

# Training Framework

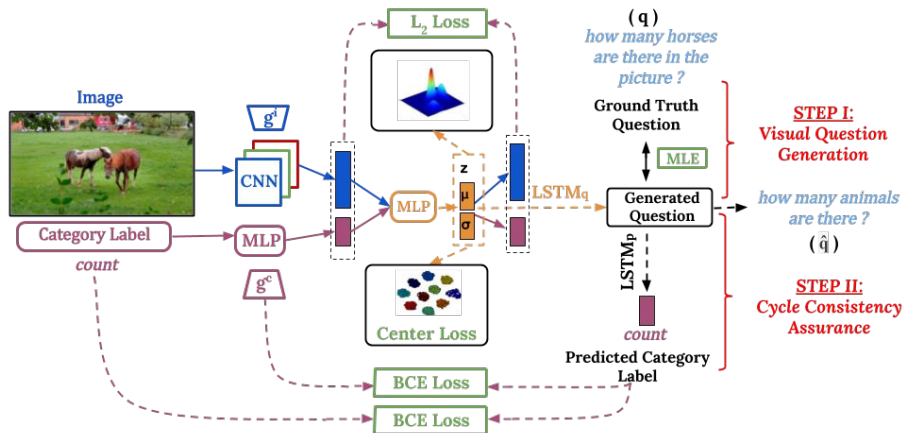


Figure 1: C3VQG Training

# Inference Framework

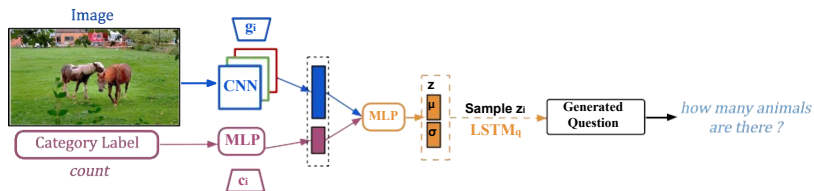


Figure 2: C3VQG Inference

$$\mathcal{L}_{center} = \|z_k - c_k\|_2^2$$

- ▶ Helps distinguish inter-category latent features by enforcing clustering.
- ▶ Centers are obtained by averaging the features of the corresponding classes updated based on mini-batches instead of the entire training data due to computational time constraints
- ▶ Update of these centers are scaled by a constant ( $< 1$ ) to avoid sudden fluctuations.

$$\mathcal{L}_{bayes} = \sum_{j=1}^d \mathbb{E}_{pd(x_k^{cc})} [KL(f(z_{k,j}|x_{k,j}^{cc})||\mathcal{N}(z_{k,j}; 0, \alpha^{-1}))]$$
$$+\lambda_{reg} \sum_{j=1}^d (\alpha_j^{-1} - 1)^2$$

- ▶ A hyper-prior on learning the inverse variance of the variational latent prior
- ▶ Helps to capture intrinsically independent visual features within the combined latent space.
- ▶ This helps us in generating more diverse questions.

# Experimental Results - Qualitative Generations

 <p><b>BINARY</b> is the image colored?</p> <p><b>COLOR</b> what is the color of the animal?</p> <p><b>OBJECT</b> what is the animal eating?</p> <p><b>MATERIAL</b> what is the fence made of?</p>	 <p><b>BINARY</b> is the man wearing a hat?</p> <p><b>COUNT</b> how many people are in the photo?</p> <p><b>ACTIVITY</b> what is the man doing?</p> <p><b>OBJECT</b> what is the man holding?</p>	 <p><b>ATTRIBUTE</b> is the building on the right tall?</p> <p><b>COLOR</b> what is the color of the building?</p> <p><b>SPATIAL</b> is the tree on left of the building?</p> <p><b>MATERIAL</b> what is the building made of?</p>	 <p><b>BINARY</b> is this a kitchen?</p> <p><b>COUNT</b> how many layers is the cake?</p> <p><b>ATTRIBUTE</b> how does the cake taste?</p> <p><b>OBJECT</b> what are the people eating?</p>
 <p><b>COUNT</b> how many people are there?</p> <p><b>COLOR</b> what color is the ground?</p> <p><b>ACTIVITY</b> what is the woman doing?</p> <p><b>OBJECT</b> what is the woman holding?</p>	 <p><b>BINARY</b> is the car parked in a garage?</p> <p><b>COLOR</b> what color is the car?</p> <p><b>SPATIAL</b> where is the car parked?</p> <p><b>MATERIAL</b> what is the building made of?</p>	 <p><b>COUNT</b> how many cars are there?</p> <p><b>COLOR</b> what color is the truck?</p> <p><b>ACTIVITY</b> what is the man in middle doing?</p> <p><b>SPATIAL</b> where is the car parked?</p>	 <p><b>BINARY</b> is the train in the station?</p> <p><b>COLOR</b> what is the color of the train?</p> <p><b>ATTRIBUTE</b> is this a modern train?</p> <p><b>MATERIAL</b> what are the tracks made of?</p>
 <p><b>BINARY</b> is the television on or off?</p> <p><b>ANIMAL</b> what kind of animal is this?</p> <p><b>OTHERS</b> what brand is the computer?</p> <p><b>OBJECT</b> what is the cat sitting on?</p>	 <p><b>BINARY</b> is the man happy?</p> <p><b>COLOR</b> what color is the man's shirt?</p> <p><b>ACTIVITY</b> what is the man in the middle doing?</p> <p><b>OBJECT</b> what is the man carrying?</p>	 <p><b>COUNT</b> how many people are there?</p> <p><b>COLOR</b> what color is the man's shirt?</p> <p><b>ACTIVITY</b> what is the man doing?</p> <p><b>OBJECT</b> what is the man riding?</p>	 <p><b>COUNT</b> how many cars are there?</p> <p><b>COLOR</b> what color is the car?</p> <p><b>ATTRIBUTE</b> is the building tall?</p> <p><b>SPATIAL</b> where is the car on the left parked?</p>

Figure 3: Question generated for each image from multiple answer categories using our approach.



## Experimental Results - Qualitative Generations

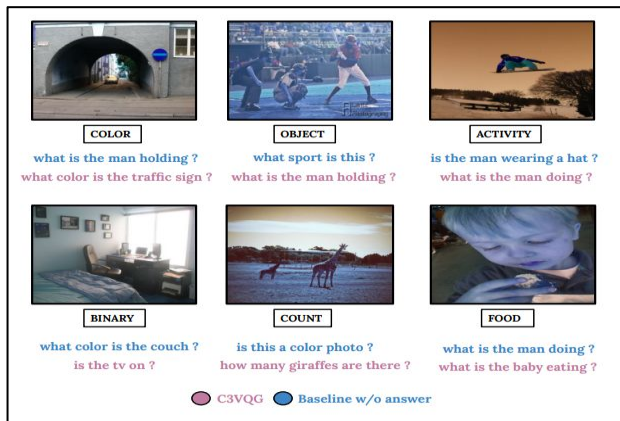


Figure 4: Qualitative results for C3VQG and Krishna et. al<sup>1</sup> without answers.

<sup>1</sup>Krishna, Bernstein, and Fei-Fei, "[Information Maximizing Visual Question Generation](#)".

We evaluate the efficacy of our approach using a set of evaluation metrics.

- ▶ Language Modelling Metrics: BLEU, METEOR, CIDEr, ROUGE-L
- ▶ Diversity Based Metrics
- ▶ Relevance Based Metrics (Crowd Sourced Metrics)

## Experimental Results - Quantitative metrics

Supervision	Models	Bleu-1	Bleu-2	Bleu-3	Bleu-4	METEOR	CIDEr	ROUGE-L
Supervised (w A)	IA2Q [24]	32.43	15.49	9.24	6.23	11.21	36.22	-
	V-IA2Q [9]	36.91	17.79	10.21	6.25	12.39	36.39	-
	Krishna <i>et al.</i> [14]	47.40	28.95	19.93	14.49	18.35	85.99	49.10
Weakly Supervised (w/o A)	IC2Q [24]	30.42	13.55	6.23	4.44	9.42	27.42	-
	V-IC2Q [9]	35.40	<b>25.55</b>	14.94	<b>10.78</b>	13.35	42.54	-
	Krishna <i>et al.</i> [14] w/o A	31.20	16.20	11.18	6.24	12.11	35.89	40.27
	<i>I</i>	38.44	19.83	12.02	7.69	13.27	45.19	40.90
	<i>I + II</i>	38.80	20.12	12.32	7.96	13.40	46.42	41.27
	<i>I + CL</i>	38.81	20.14	12.30	7.91	13.41	46.96	41.21
	<i>I + II + CL</i>	38.94	20.30	12.47	8.10	13.47	<b>47.32</b>	41.27
	<i>I + II + Bayes</i>	38.71	19.89	12.14	7.87	13.23	42.47	41.32
	<i>I + CL + Bayes</i>	38.64	20.06	12.28	7.95	13.32	45.83	41.16
<i>I + II + CL + Bayes</i>	<b>41.87</b>	22.11	<b>14.96</b>	10.04	<b>13.60</b>	46.87	<b>42.34</b>	

Table 1: Ablation study for different components of C3VQG using different language modeling quantitative metrics against other baselines in VQG. We compare our approach against previous works using answers as well as without answers.

## Experimental Results - Quantitative metrics

Categories	V-IC2Q [9]		Krishna <i>et al.</i> [14]		C3VQG w/o Bayes		C3VQG	
	Strength	Inventiveness	Strength	Inventiveness	Strength	Inventiveness	Strength	Inventiveness
count	15.77	30.91	26.06	41.30	58.33	55.20	65.21	61.84
binary	18.15	41.95	28.85	54.50	58.39	36.32	65.12	38.55
object	11.27	34.84	24.19	43.20	57.77	51.51	65.58	58.85
color	4.03	13.03	17.12	23.65	58.38	48.97	65.21	54.34
attribute	37.76	41.09	46.10	52.03	60.05	58.38	64.59	63.02
materials	36.13	31.13	45.75	40.72	57.93	56.79	64.87	63.48
spatial	61.12	62.54	70.17	68.18	57.90	57.80	65.18	64.96
food	21.81	20.38	33.37	31.19	58.49	55.42	65.20	62.21
shape	35.51	44.03	45.81	55.65	58.85	58.75	66.01	65.98
location	34.68	18.11	45.25	27.22	58.39	58.10	65.09	64.72
predicate	22.58	17.38	36.20	31.29	57.05	57.05	65.67	65.67
time	25.58	15.51	34.43	25.30	58.13	58.10	65.00	64.96
activity	7.45	13.23	21.32	26.53	58.00	56.78	64.98	63.67
Overall	12.97	38.32	26.06	52.11	58.23	54.99	<b>65.24</b>	<b>61.55</b>

Table 2: Quantitative evaluation of C3VQG against other baselines using diversity-based metrics.

Model	Relevance	
	Image	Category
V-IC2Q [9]	90.10	39.00
Krishna <i>et al.</i> [14] w/o A	<b>98.10</b>	42.70
C3VQG w/o Bayes, CL	98.00	58.40
C3VQG	97.80	<b>60.50</b>

**Table 3: Quantitative evaluation of C3VQG against other weakly supervised baselines using crowd-sourced metrics.**

Thank You !

For more details, please check our paper:

[C3VQG: Category Consistent Cyclic Visual Question Generation](#)