

Motivation

Task: Given an image as input, generate diverse questions **Applications:**

- AI assistants using visual cues
- VQA and Robotics
- Engaging machine-human interaction
- Active teaching of content

Contributions:

- Generative capability of VAE + Representative from LSTM nets
- Low dimensional latent space embedding of questions
- Sample from 'latent space' to generate new questions





Introduction & Related Work

Previous methods:

- Mostafazadeh et al. (ACL '16): Pretrained CNN + GRU Appealing model, one question per image
- Vijaykumar et al. (ICLR '17): Diverse beam search Samples from complex energy landscape, a hard task in general

Our Approach: (Inference)

- Sample from an embedding space *z*
- Decode sample to form a question *x*

How to construct the embedding space?

Variational Auto-encoder:

- Construct image embedding
- Encode image embedding and sentence
- Sample from constructed embedding
- Reconstruction sentence while taking into account image



Maximizing likelihood: parametric distribution $p_{\theta}(x)$ Introduce posterior $q_{\phi}(z|x)$ (manifold assumption)

$$\ln p_{\theta}(x) = \sum_{z} q_{\phi}(z|x) \ln p_{\theta}(x) = \mathcal{L}(q_{\phi}, p_{\theta}(x, z)) + \mathrm{KL}(z)$$

Lower bound: $\mathcal{L}(q_{\phi}, p_{\theta}(x, z)) = \sum_{z} q_{\phi}(z|x) \ln \frac{p_{\theta}(x|z)p(z)}{q_{\phi}(z|x)}$

Generating Diverse Questions using VAE Unnat Jain^{*1} Ziyu Zhang^{*2} Alexander Schwing¹

²Snapchat Research

Results (continued)

Accuracy metrics: Bleu, Meteor



Generated questions:



What is the num train?
Is this a mode station
Is this train ir setting
Is this train in t states?







Error Analysis



Recognition failures: Pre-learned visual features are incapable of capturing fine information

Co-occurrence based failures: Objects not in the image appear in questions.





