

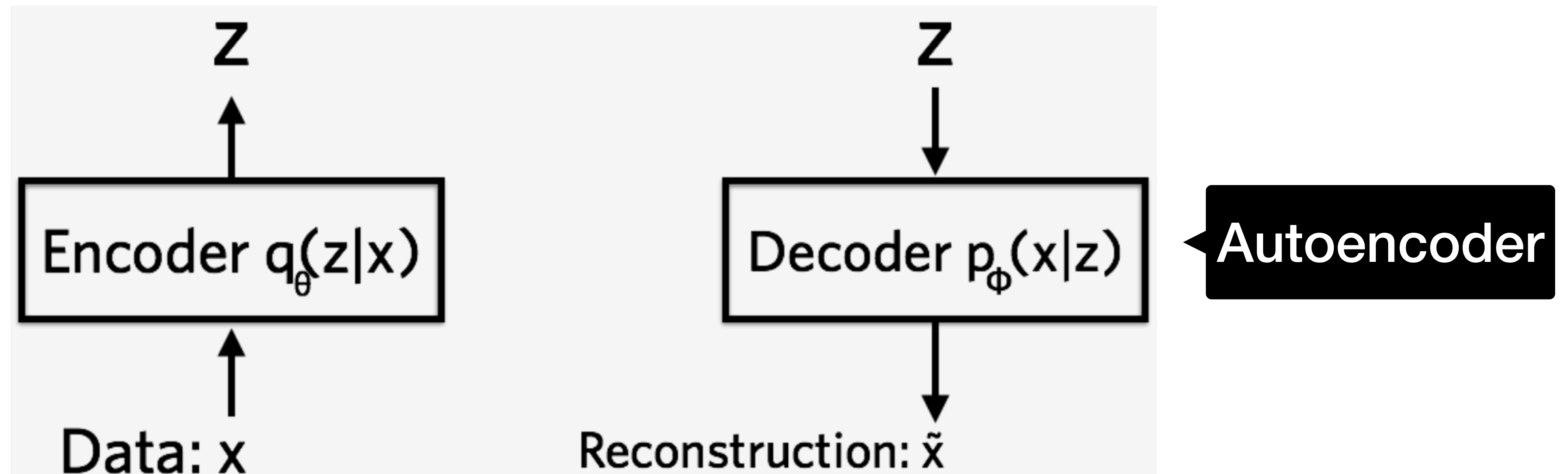
# **Variational Auto-encoding**

**Advanced NLP: Summer 2023**

**Anoop Sarkar**

# Encoder-Decoder neural nets

- An encoder  $q$  takes an input  $x$  and encodes it into a hidden representation  $z$  using some parameters  $\theta$ . Encoder:  $q_{\theta}(z | x)$
- An decoder  $p$  takes a hidden layer  $z$  and decodes it into an output  $\tilde{x}$  using some parameters  $\phi$ . Decoder:  $\tilde{x} \sim p_{\phi}(x | z)$
- The output  $\tilde{x}$  should be similar to but not necessarily identical to the true  $x$



# Autoencoder loss

- How much information is lost by going from  $x$  to  $z$  and then back to  $\tilde{x}$ ?
- We measure the information loss by representing using  $z$  using reconstruction log-likelihood
- $\log p_{\phi}(x | z)$  measured in nats (bits are base 2, nats are base  $e$ )
- The loss function for an *variational* autoencoder is the negative log likelihood with a regularizer
- For single data point  $x_i$  we compute the above loss  $l_i$ .
- Total loss for the dataset:  $\sum_i l_i$

# Variational autoencoder loss

- Loss function  $l_i$  for datapoint  $x_i$  is
- $l_i(\theta, \phi) = -E_{z \sim q_\theta(z|x_i)}[\log p_\phi(x_i | z)] + KL(q_\theta(z | x_i) || p(z))$
- First term is the expected negative log-likelihood of the data point  $x_i$
- We want to place the most probability mass on the true output  $x_i$
- Second term is the regularizer: the Kullback-Leibler divergence between the encoder distribution  $q_\theta(z | x)$  and  $p(z)$
- $p(z)$  is used to reward "good" values of the hidden representation that are efficient, can be sampled from easily and do not memorize the dataset.
- $z = \mu + \sigma \circ \epsilon$  where  $\epsilon \sim \text{Normal}(0,1)$  and  $\circ$  is element wise multiplication

# Reparametrize $z$

- We want to use gradient descent to learn  $q_{\theta}(z | x)$
- Need to take derivative of  $p(z)$  wrt  $\theta$
- We reparametrize  $z$
- $z = \mu + \sigma \circ \epsilon$  where  $\epsilon \sim \text{Normal}(0,1)$  and  $\circ$  is element wise multiplication
- Now we can take derivatives of  $p(z)$  wrt  $\mu$  and  $\sigma$
- Output of  $q_{\theta}(z | x)$  is a vector of  $\mu$ 's and  $\sigma$ 's

# Variational autoencoder loss

- The regularizer term keeps the representation of  $z$  sufficiently diverse
- Without the regularizer, given large enough  $z$  the encoder-decoder would simply memorize the entire dataset
- Two different  $x_i$  and  $x_j$  that are actually very close to each other would end up learning very different  $z_i$  and  $z_j$  which defeats the purpose of modeling similarity between inputs.
- The regularizer would make sure  $z_i$  and  $z_j$  cannot get too far from each other unless  $x_i$  is very different from  $x_j$
- The variational autoencoder (vae) is trained using gradient descent

# Variational autoencoder loss

- Unfortunately, gradient descent requires computing distribution  $q_{\theta}(z | x)$
- This is exponential because it is over all configurations of latent variable  $z$
- Variational inference approximates this using a distribution  $q_{\lambda}(z | x)$
- $\lambda$  is the variational parameter which indexes a family of distributions
- If  $q$  is a normal distribution then  $\lambda_{x_i}$  would be the mean  $\mu$  and variance  $\sigma^2$  for each data point  $x_i$
- $\lambda_{x_i} = (\mu_{x_i}, \sigma_{x_i}^2)$

# Tractable variational inference

- We want to measure how well does the variational distribution  $q_\lambda(z | x)$  approximate the true distribution  $q(z | x)$
- We use the KL divergence again:  $KL(q_\lambda(z | x) || q(z | x))$
- The optimal approximate distribution involves finding the optimal variational parameters  $\lambda$
- $q_\lambda^*(z | x) = \arg \min_{\lambda} KL(q_\lambda(z | x) || q(z | x))$
- Unfortunately, this is still intractable



# Tractable variational inference

- Define ELBO( $\lambda$ ) the Evidence Lower BOund of  $\lambda$
- $$\text{ELBO}(\lambda) = E_{z \sim q_\lambda} [\log p(x | z)] - E_{z \sim q_\lambda} [\log q_\lambda(z | x)]$$
- Minimizing  $\text{KL}(q_\lambda \| p)$  wrt  $\lambda$  is equivalent to maximizing ELBO( $\lambda$ )
- For each data point  $x_i$
- $$\text{ELBO}_i(\lambda) = E_{z \sim q_\lambda(z|x_i)} [\log p_\phi(x_i | z)] - \text{KL}(q_\lambda(z | x_i) \| p(z))$$
- Maximizing  $\text{ELBO}_i(\lambda)$  is equivalent to minimizing 
$$l_i(\theta, \phi) = - E_{z \sim q_\theta(z|x_i)} [\log p_\phi(x_i | z)] + \text{KL}(q_\theta(z | x_i) \| p(z))$$

# Applications: Image generation

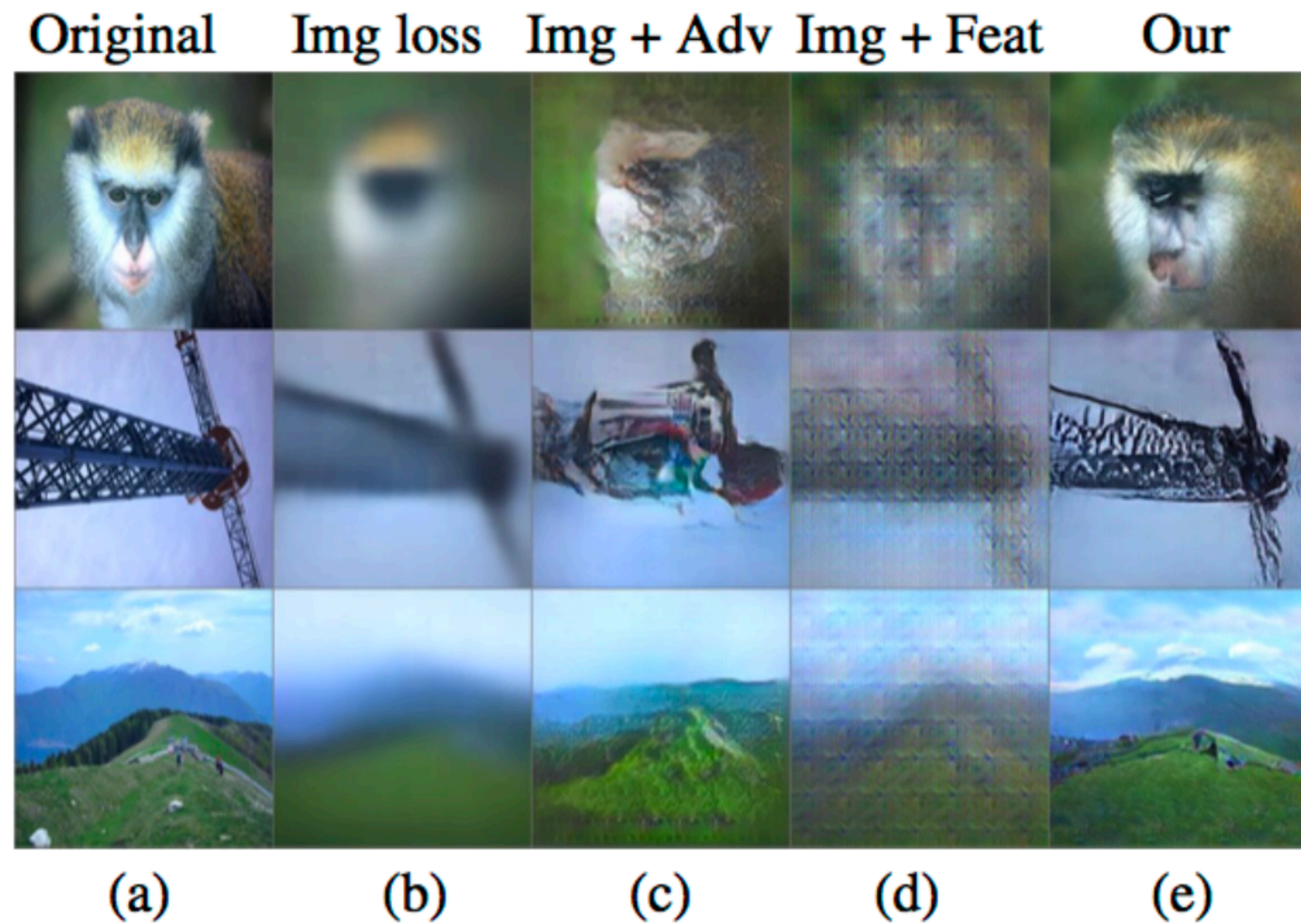


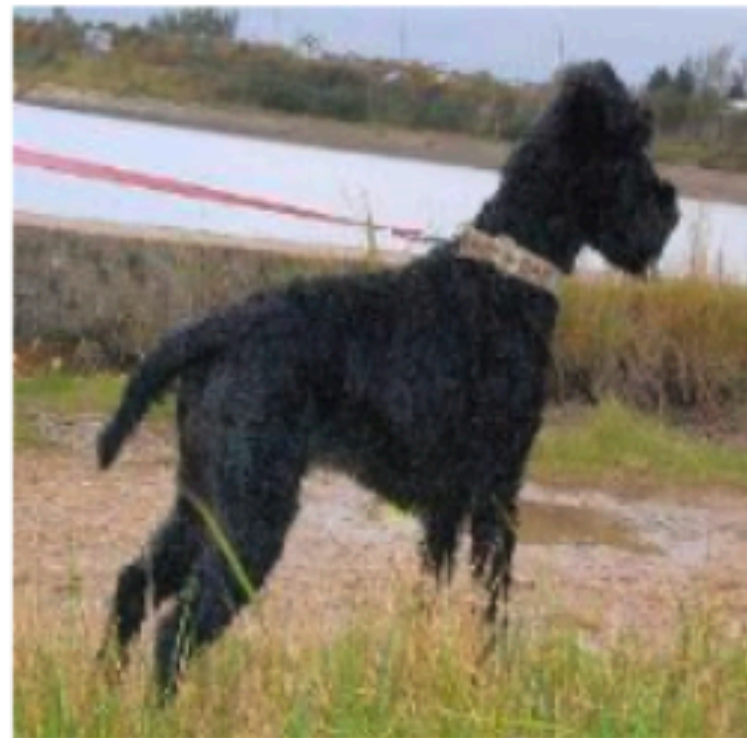
Figure 1: Reconstructions from AlexNet FC6 with different components of the loss.

A. Dosovitskiy and T. Brox. Generating images with perceptual similarity metrics based on deep networks. *arXiv preprint arXiv:1602.02644*, 2016.

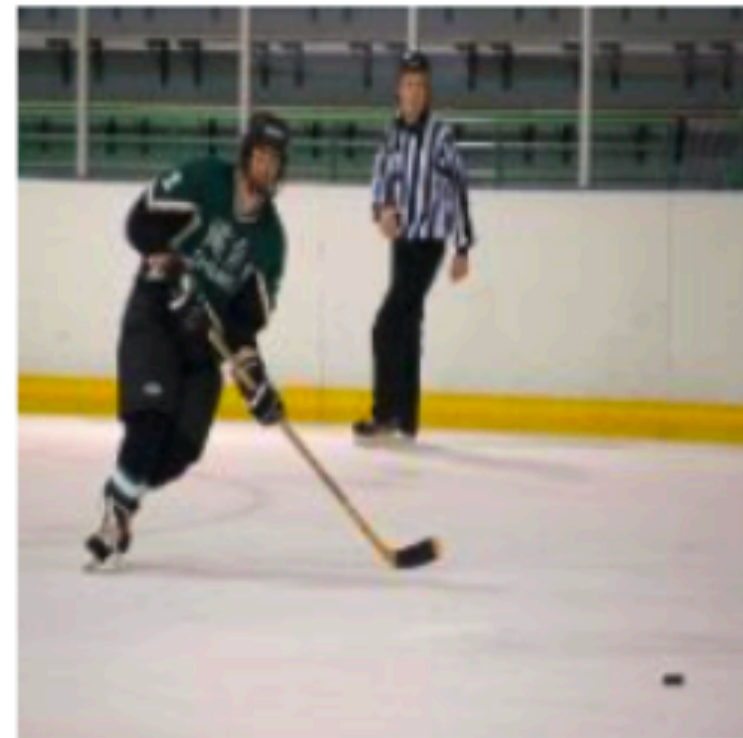
# Applications: caption generation



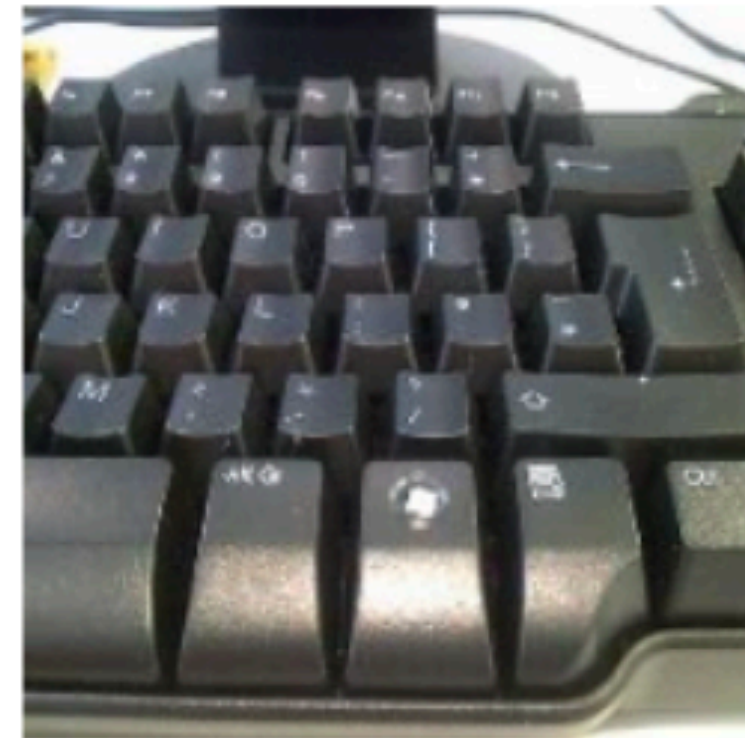
a man with a snowboard  
next to a man with glasses



a big black dog standing on  
the grass



a player is holding a  
hockey stick



a desk with a keyboard



a man is standing next to a  
brown horse

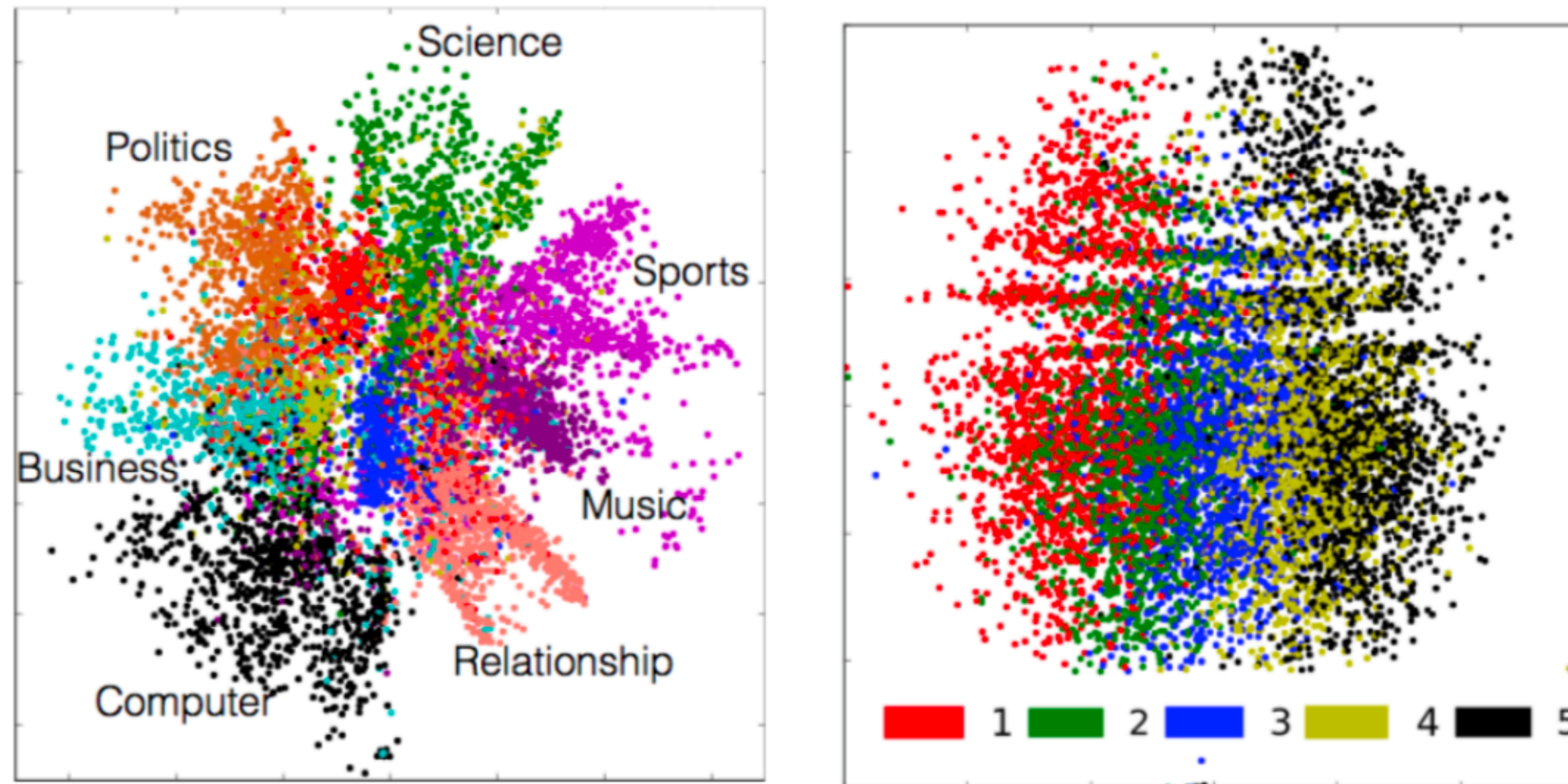


a box full of apples and  
oranges

**Figure 2: Examples of generated caption from unseen images on the validation dataset of ImageNet.**

Y. Pu, Z. Gan, R. Henao, X. Yuan, C. Li, A. Stevens, and L. Carin. Variational autoencoder for deep learning of images, labels and captions. In NIPS, 2016.

# Applications: document clustering



(a) Yahoo

(b) Yelp

Figure 3: Visualizations of learned latent representations.

Z. Yang, Z. Hu, R. Salakhutdinov, and T. Berg-Kirkpatrick. Improved variational autoencoders for text modeling using dilated convolutions. In *Proceedings of The 34th International Conference on Machine Learning*, 2017.

# Applications: sign clustering

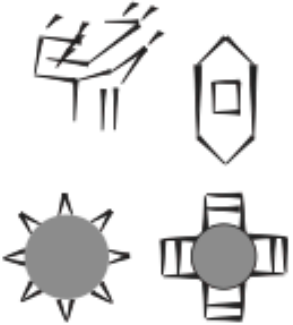
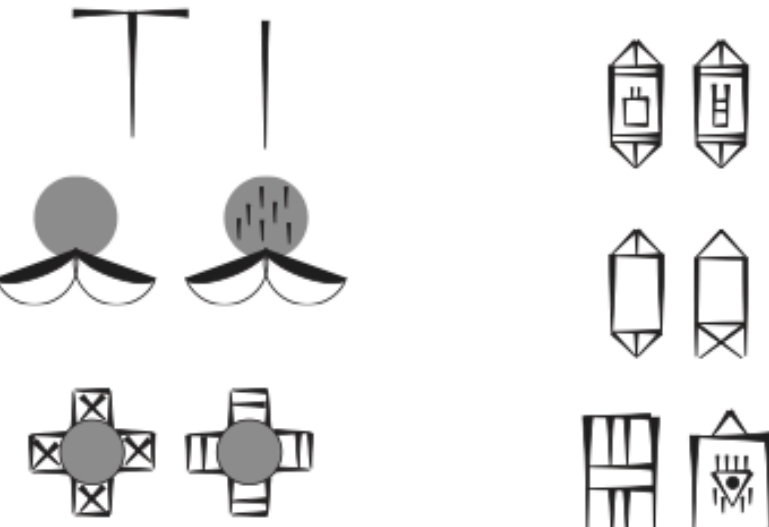
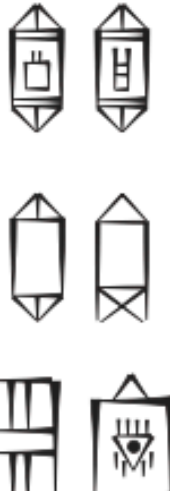
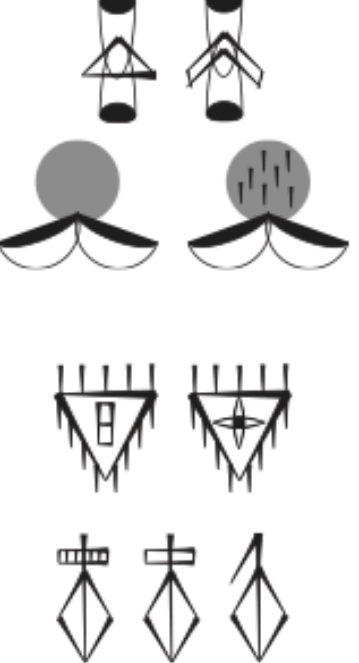
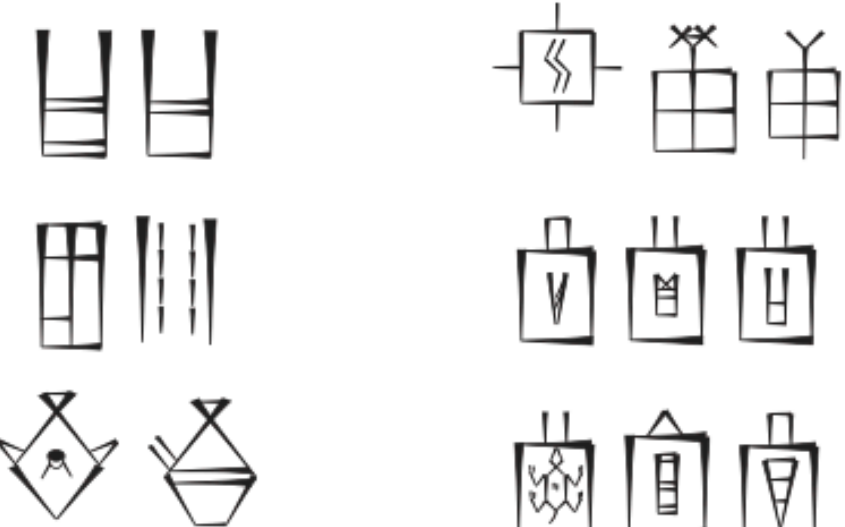

VAE+Neighbor	VAE+LSTM	VAE+Transformer			
					

Table 3: Pairs/triplets of character images which have distinct labels in the working signlist, but which our models merge into single clusters.