

# Variational Attention for Sequence-to-Sequence Models

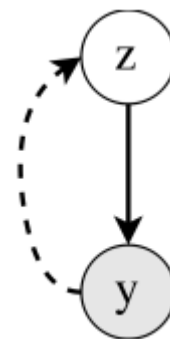
Hareesh Bahuleyan, Lili Mou, Olga Vechtomova, Pascal Poupart  
University of Waterloo, ON, Canada

# Variational Autoencoder

## □ Variational Autoencoder (VAE)

- Encode data  $Y$  as hidden random variables  $Z$ , then reconstruct  $Y$ .
- VAE models both  $q_{\phi}(z|y)$  and  $p_{\theta}(y|z)$  with neural networks.
- Training objective:

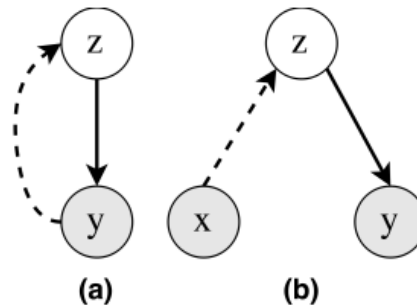
$$\underbrace{J_{\text{rec}}(\boldsymbol{\theta}, \boldsymbol{\phi}, \mathbf{y}^{(n)})}_{\text{reconstruction loss}} + \underbrace{\text{KL}\left(q_{\phi}(\mathbf{z}|\mathbf{y}^{(n)})\|p(\mathbf{z})\right)}_{\text{regularization}}$$



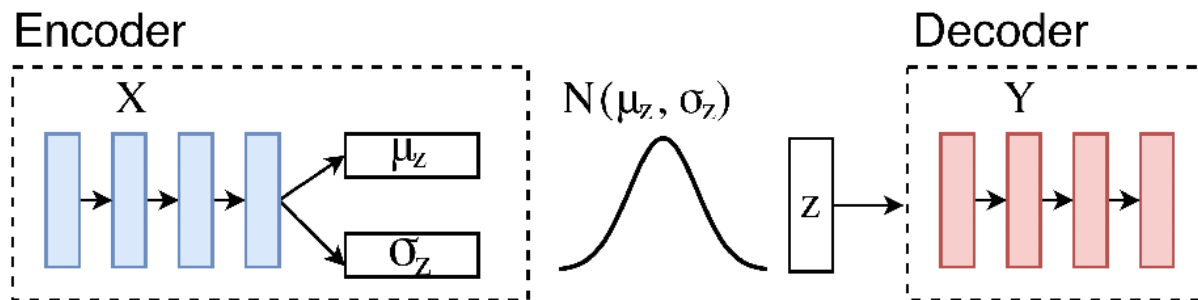
# Variational Encoder-Decoder

## □ Variational Encoder-Decoder (VED)

- Extend VAE to VED.



- Model:

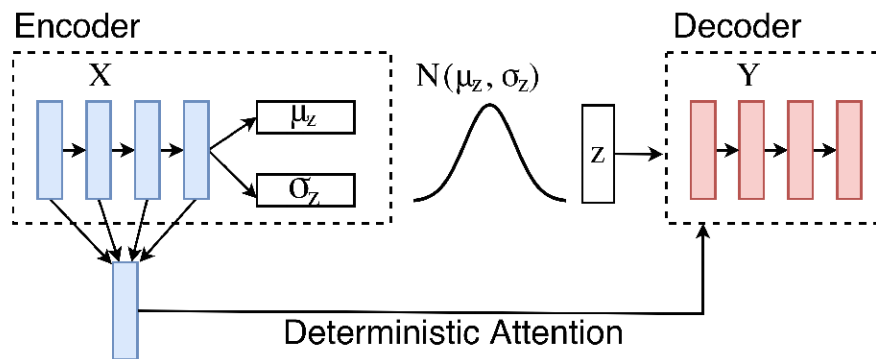


# Attention Mechanism

- ❑ Attention significantly improves Seq2Seq performance in translation, summarization, etc.

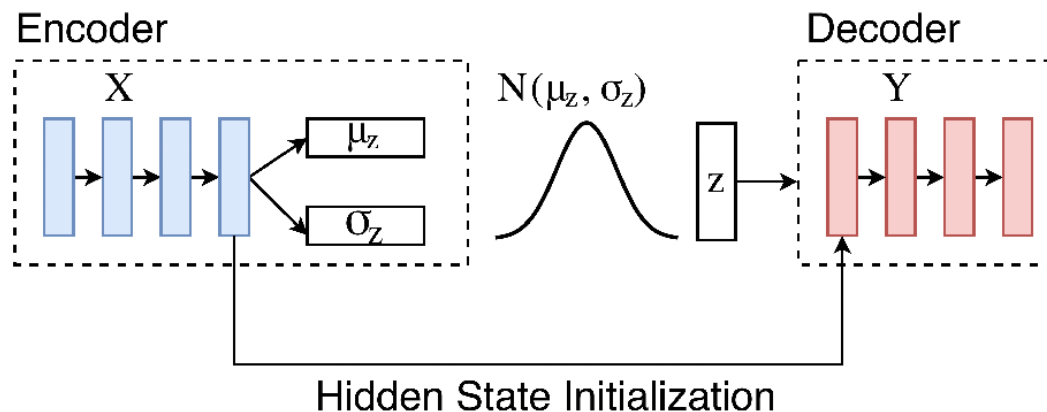
$$\alpha_{ji} = \frac{\exp\{\tilde{\alpha}_{ji}\}}{\sum_{i'=1}^{|\mathbf{x}|} \exp\{\tilde{\alpha}_{ji'}\}} \quad \mathbf{a}_j = \sum_{i=1}^{|\mathbf{x}|} \alpha_{ji} \mathbf{h}_i^{(\text{src})}$$

- ❑ Add attention to variational Seq2Seq model?



# “Bypassing” Phenomenon

- ❑ Observation: if the decoder has a direct, deterministic access to the encoder, the latent variables  $Z$  might not capture much information, so that VED does not learn much and play a role in the process.
- ❑ Example: variational Seq2Seq with hidden state initialization



# A Pilot Study

---

**Input:** *the men are playing musical instruments*

---

**(a) VAE w/o hidden state init. (Avg entropy: 2.52)**

---

*the men are playing musical instruments*

*the men are playing video games*

*the musicians are playing musical instruments*

*the women are playing musical instruments*

---

**(b) VAE w/ hidden state init. (Avg entropy: 2.01)**

---

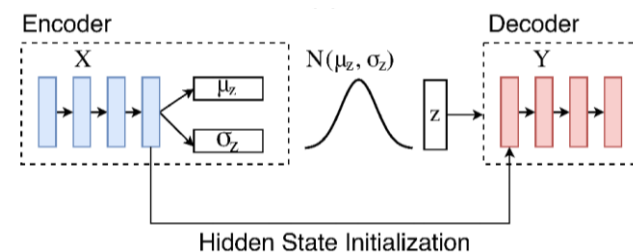
*the men are playing musical instruments*

*the men are playing musical instruments*

*the men are playing musical instruments*

*the man is playing musical instruments*

---

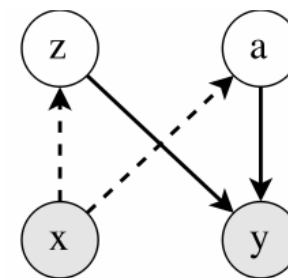
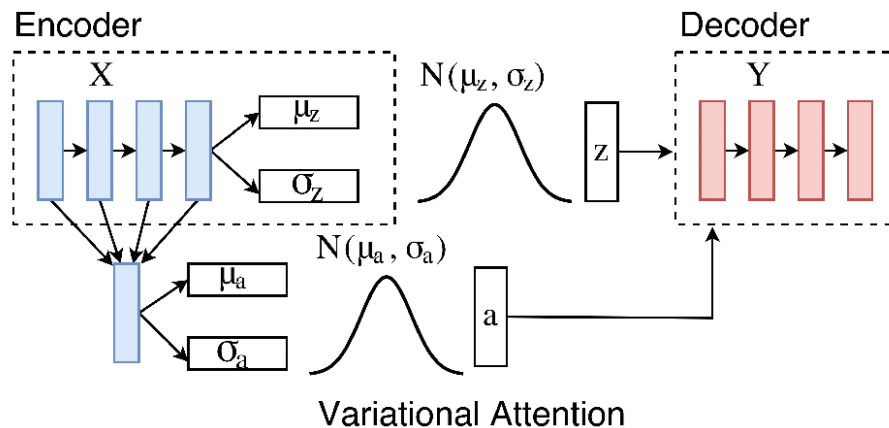


*Less diversified!*

Table 1: Sentences obtained by sampling from the VAE latent space. (a) VAE without hidden state initialization. (b) VAE with hidden state initialization.

# Variational Attention

- Treat latent space  $Z$  and attention vector  $a_j$  as random variables.



- Variational lower bound

$$\begin{aligned} \mathcal{L}_j^{(n)}(\theta, \phi) &= \mathbb{E}_{\mathbf{z}, \mathbf{a} \sim q_\phi(\mathbf{z}, \mathbf{a} | \mathbf{x}^{(n)})} \left[ \log p_\theta(\mathbf{y}^{(n)} | \mathbf{z}, \mathbf{a}) \right] \\ &\quad - \text{KL} \left( q_\phi(\mathbf{z}, \mathbf{a} | \mathbf{y}^{(n)}) \| p(\mathbf{z}, \mathbf{a}) \right) \end{aligned} \quad (6)$$

$$\begin{aligned} &= \mathbb{E}_{\mathbf{z} \sim q_\phi^{(z)}(\mathbf{z} | \mathbf{x}^{(n)}), \mathbf{a} \sim q_\phi^{(a)}(\mathbf{a} | \mathbf{x}^{(n)})} \left[ \log p_\theta(\mathbf{y}^{(n)} | \mathbf{z}, \mathbf{a}) \right] \\ &\quad - \text{KL} \left( q_\phi^{(z)}(\mathbf{z} | \mathbf{y}^{(n)}) \| p(\mathbf{z}) \right) \\ &\quad - \text{KL} \left( q_\phi^{(a)}(\mathbf{a} | \mathbf{y}^{(n)}) \| p(\mathbf{a}) \right) \end{aligned} \quad (7)$$

$z$  and  $a$  are conditional independent given  $x$ , also marginally independent.

# Variational Attention

## □ Prior $p(a_j)$

- The same as  $z$ , set  $p(a_j) = N(0, I)$ .
- $p(a_j) = N(\bar{h}^{(src)}, I)$ , where  $\bar{h}^{(src)} = \frac{1}{|x|} \sum_{i=1}^{|x|} h_i^{(src)}$ .

## □ Training Objective

$$\begin{aligned} J^{(n)}(\boldsymbol{\theta}, \boldsymbol{\phi}) = & J_{\text{rec}}(\boldsymbol{\theta}, \boldsymbol{\phi}, \mathbf{y}^{(n)}) \\ & + \lambda_{\text{KL}} \left[ \text{KL} \left( q_{\boldsymbol{\phi}}^{(z)}(\mathbf{z}) \| p(\mathbf{z}) \right) \right. \\ & \left. + \gamma_a \sum_{j=1}^{|y|} \text{KL} \left( q_{\boldsymbol{\phi}}^{(a)}(\mathbf{a}_j) \| p(\mathbf{a}_j) \right) \right] \end{aligned}$$



# Experiments

- ❑ Task: generate questions based on a sentence in a paragraph.
- ❑ Generated questions do need some variety.
- ❑ Dataset: follow Du et al. (2017) and use SQuAD.
- ❑ Metrics:
  - Measure accuracy: BLEU-1 to BLEU-4.
  - Measure diversity: entropy and distinct metrics.

The percentage of distinct unigrams or bigrams.

# Results

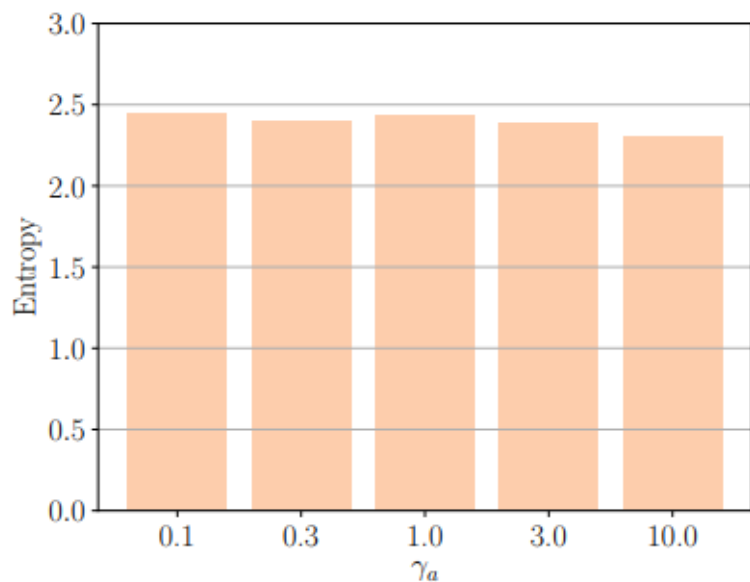
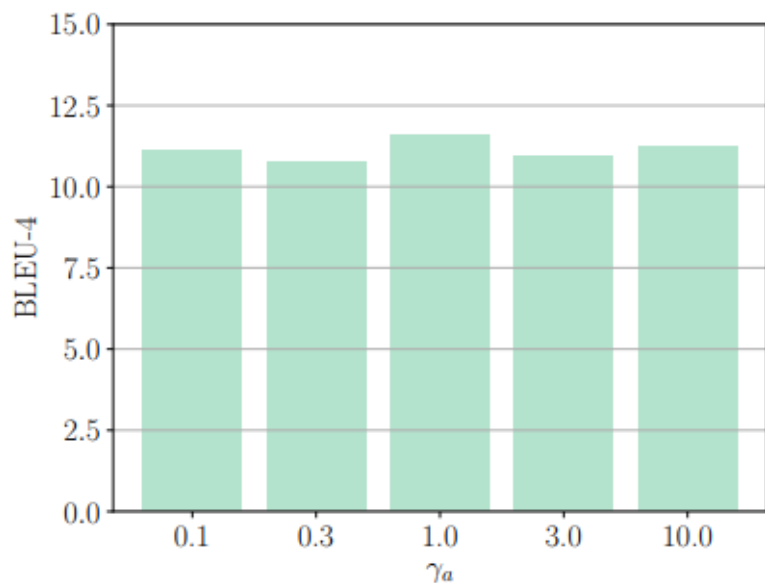
Model	Inference	BLEU-1	BLEU-2	BLEU-3	BLEU-4	Entropy	Dist-1	Dist-2
Previous work (Du et al., 2017)	MAP	43.09	25.96	17.50	12.28	-	-	-
DED (w/o Attn)	MAP	39.46	28.49	20.74	8.10	-	-	-
DED+DAttn	MAP	42.34	30.86	22.74	11.60	-	-	-
VED+DAttn	MAP	<b>42.50</b>	<b>31.13</b>	<b>23.09</b>	<b>12.38</b>	-	-	-
	Sampling	42.48	31.10	23.08	12.30	2.37	0.18	0.26
VED+DAttn (2-stage training)	MAP	42.17	30.96	22.95	11.98	-	-	-
	Sampling	41.98	30.82	22.81	11.78	2.41	0.19	0.27
VED+VAttn-0	MAP	41.77	30.54	22.53	11.37	-	-	-
	Sampling	41.73	30.51	22.49	11.27	<b>2.44</b>	<b>0.20</b>	0.28
VED+VAttn- $\bar{h}$	MAP	42.10	30.71	22.70	11.55	-	-	-
	Sampling	42.03	30.62	22.66	11.50	<b>2.44</b>	<b>0.20</b>	<b>0.29</b>

Table 2: BLEU, entropy, and distinct scores. We compare the deterministic encoder-decoder (DED) and variational encoder-decoders (VEDs). For VED, we have several variates: deterministic attention (DAttn) and the proposed variational attention (VAttn). We evaluate the sentences obtained by both max *a posteriori* (MAP) inference and sampling.

Sampling: draw 10 samples each time.

2-stage training: train VED without attention for first 20 epochs.

# Strength of Attention's KL Loss



$$\begin{aligned} J^{(n)}(\boldsymbol{\theta}, \boldsymbol{\phi}) = & J_{\text{rec}}(\boldsymbol{\theta}, \boldsymbol{\phi}, \mathbf{y}^{(n)}) \\ & + \lambda_{\text{KL}} \left[ \text{KL} \left( q_{\boldsymbol{\phi}}^{(z)}(\mathbf{z}) \| p(\mathbf{z}) \right) \right. \\ & \left. + \gamma_a \sum_{j=1}^{|\mathbf{y}|} \text{KL} \left( q_{\boldsymbol{\phi}}^{(a)}(\mathbf{a}_j) \| p(\mathbf{a}_j) \right) \right] \end{aligned}$$

$\gamma_a$  has little effects.