

Sentence Generation

Lili Mou

l mou@ualberta.ca

lili-mou.github.io



UNIVERSITY OF
ALBERTA

Roapmap

- Motivation and examples
- Techniques
 - Generation from latent space
 - Generation from word space

Motivation

- Sentence generation
 - Dialogue systems
 - Paraphrase generation
 - Machine translation
- A seq2seq model may not suffice
 - No input
 - Constructing new information
 - Diversity needed
- Probabilistic sentence generation
 - Prior sampling, posterior sampling

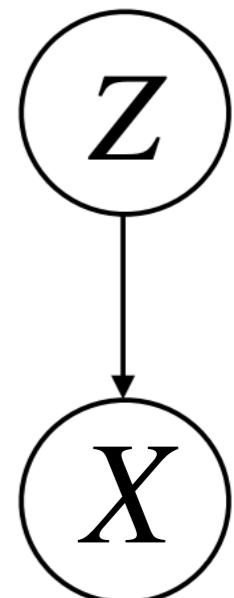


Latent Space Sampling

Variational Autoencoder

- Humans' sentence generation involves two steps
 - First, we have some “vague” idea of the sentence
 - Then, we flesh it out by words
- A sentence $x = (x_1, \dots, x_T)$ is subject to some latent representation z

$$p(z, x) = p(z)p(x | z)$$



Variational Autoencoder

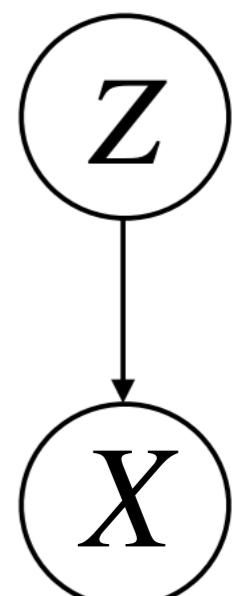
- Humans' sentence generation involves two steps
 - First, we have some “vague” idea of the sentence
 - Then, we flesh it out by words
- A sentence $x = (x_1, \dots, x_T)$ is subject to some latent representation z

$$p(z, x) = p(z)p(x | z)$$

How can we learn a model with latent variables?

$$\text{E-step: } p(z|x) = \frac{p(z)p(x|z)}{p(x)} = \frac{p(z)p(x|z)}{\int p(z')p(x|z')dz'}$$

$$\text{M-step: maximize } \mathbb{E}_{z \sim p(z|x)} \log p(z, x)$$



Variational Autoencoder

- Humans' sentence generation involves two steps
 - First, we have some “vague” idea of the sentence
 - Then, we flesh it out by words
- A sentence $x = (x_1, \dots, x_T)$ is subject to some latent representation z

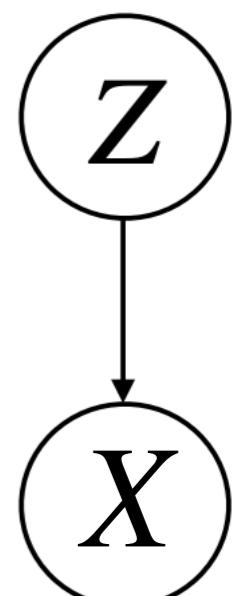
$$p(z, x) = p(z)p(x | z)$$

How can we learn a model with latent variables?

E-step: $p(z|x) = \frac{p(z)p(x|z)}{p(x)} = \frac{p(z)p(x|z)}{\int p(z')p(x|z')dz'}$

M-step: maximize $\mathbb{E}_{z \sim p(z|x)} \log p(z, x)$

Intractable



Variational Autoencoder

- Humans' sentence generation involves two steps
 - First, we have some “vague” idea of the sentence
 - Then, we flesh it out by words
- A sentence $x = (x_1, \dots, x_T)$ is subject to some latent representation z

$$p(z, x) = p(z)p(x | z)$$

How can we learn a model with latent variables?

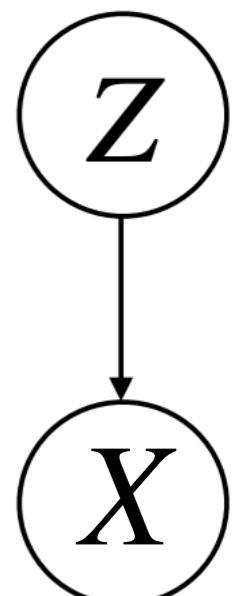
Recognition

$$\text{E-step: } p(z|x) = \frac{p(z)p(x|z)}{p(x)} = \frac{p(z)p(x|z)}{\int p(z')p(x|z')dz'}$$

Reconstruction

$$\text{M-step: maximize } \mathbb{E}_{z \sim p(z|x)} \log p(z, x)$$

(in a more general sense)



Variational Inference

$$\begin{aligned}
 \log p(\mathbf{x}; \boldsymbol{\theta}) &= \log \left(\int_z p(\mathbf{x}, z; \boldsymbol{\theta}) dz \right) \\
 &= \underbrace{\int q(z|\mathbf{x}) \log \frac{p(\mathbf{y}, z; \boldsymbol{\theta})}{q(z|\mathbf{x})} dz}_{L(q, \boldsymbol{\theta})} + \underbrace{\int q(z|\mathbf{x}) \log \frac{q(z|\mathbf{x})}{p(z|\mathbf{x}; \boldsymbol{\theta})} dz}_{\text{KL}(q(Z|\mathbf{x})||p(Z|\mathbf{x}))}
 \end{aligned}$$

Variational inference

vs

EM

Variational family

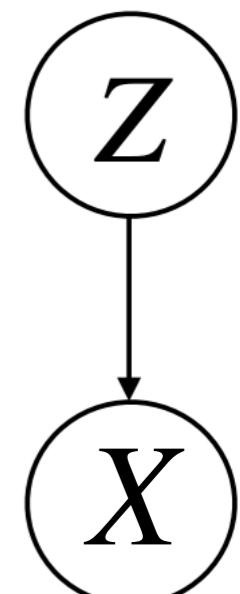
$$q \in Q$$

q can be any distribution

True posterior is the best

Ignore the KL-term

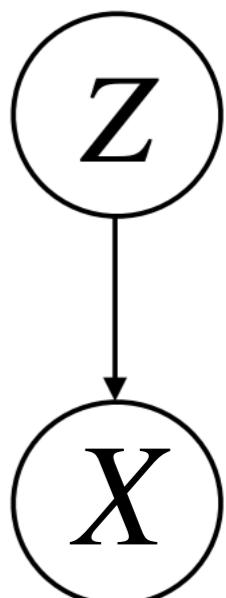
KL=0 after E step



Variational Inference

$$\begin{aligned}
 \log p(\mathbf{x}; \boldsymbol{\theta}) &= \log \left(\int_z p(\mathbf{x}, z; \boldsymbol{\theta}) dz \right) \\
 &= \underbrace{\int q(z|\mathbf{x}) \log \frac{p(\mathbf{y}, z; \boldsymbol{\theta})}{q(z|\mathbf{x})} dz}_{L(q, \boldsymbol{\theta})} + \underbrace{\int q(z|\mathbf{x}) \log \frac{q(z|\mathbf{x})}{p(z|\mathbf{x}; \boldsymbol{\theta})} dz}_{\text{KL}(q(Z|\mathbf{x})||p(Z|\mathbf{x}))}
 \end{aligned}$$

- Two extremes
 - $Q = \text{any function} \Rightarrow \text{EM}$
 \Rightarrow powerful model; optimization intractable
 - $Q = \{\text{a fixed distribution}\}$
 \Rightarrow degenerated model; optimization easy



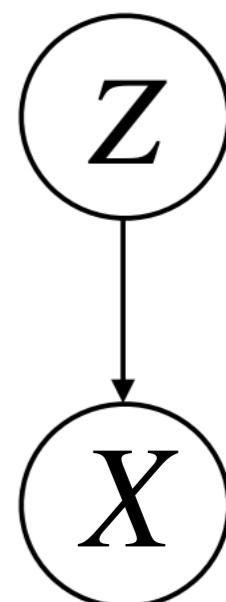
Variational Inference

$$\begin{aligned}
 \log p(\mathbf{x}; \boldsymbol{\theta}) &= \log \left(\int_z p(\mathbf{x}, z; \boldsymbol{\theta}) dz \right) \\
 &= \underbrace{\int q(z|\mathbf{x}) \log \frac{p(\mathbf{y}, z; \boldsymbol{\theta})}{q(z|\mathbf{x})} dz}_{L(q, \boldsymbol{\theta})} + \underbrace{\int q(z|\mathbf{x}) \log \frac{q(z|\mathbf{x})}{p(z|\mathbf{x}; \boldsymbol{\theta})} dz}_{\text{KL}(q(Z|\mathbf{x})||p(Z|\mathbf{x}))}
 \end{aligned}$$

- Two extremes
 - $Q = \text{any function} \Rightarrow \text{EM}$
 - ⇒ powerful model; optimization intractable
 - $Q = \{\text{a fixed distribution}\}$
 - ⇒ degenerated model; optimization easy

Trade-off, e.g.,

- Independent assumption
- Gaussian assumption



Example

[PRML]

$$p(\mathcal{D}|\mu, \tau) = \left(\frac{\tau}{2\pi}\right)^{N/2} \exp \left\{ -\frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2 \right\}$$

$$\begin{aligned} p(\mu|\tau) &= \mathcal{N}(\mu|\mu_0, (\lambda_0\tau)^{-1}) \\ p(\tau) &= \text{Gam}(\tau|a_0, b_0) \end{aligned}$$

Variational family: factorized distribution

$$q(\mu, \tau) = q_\mu(\mu)q_\tau(\tau)$$

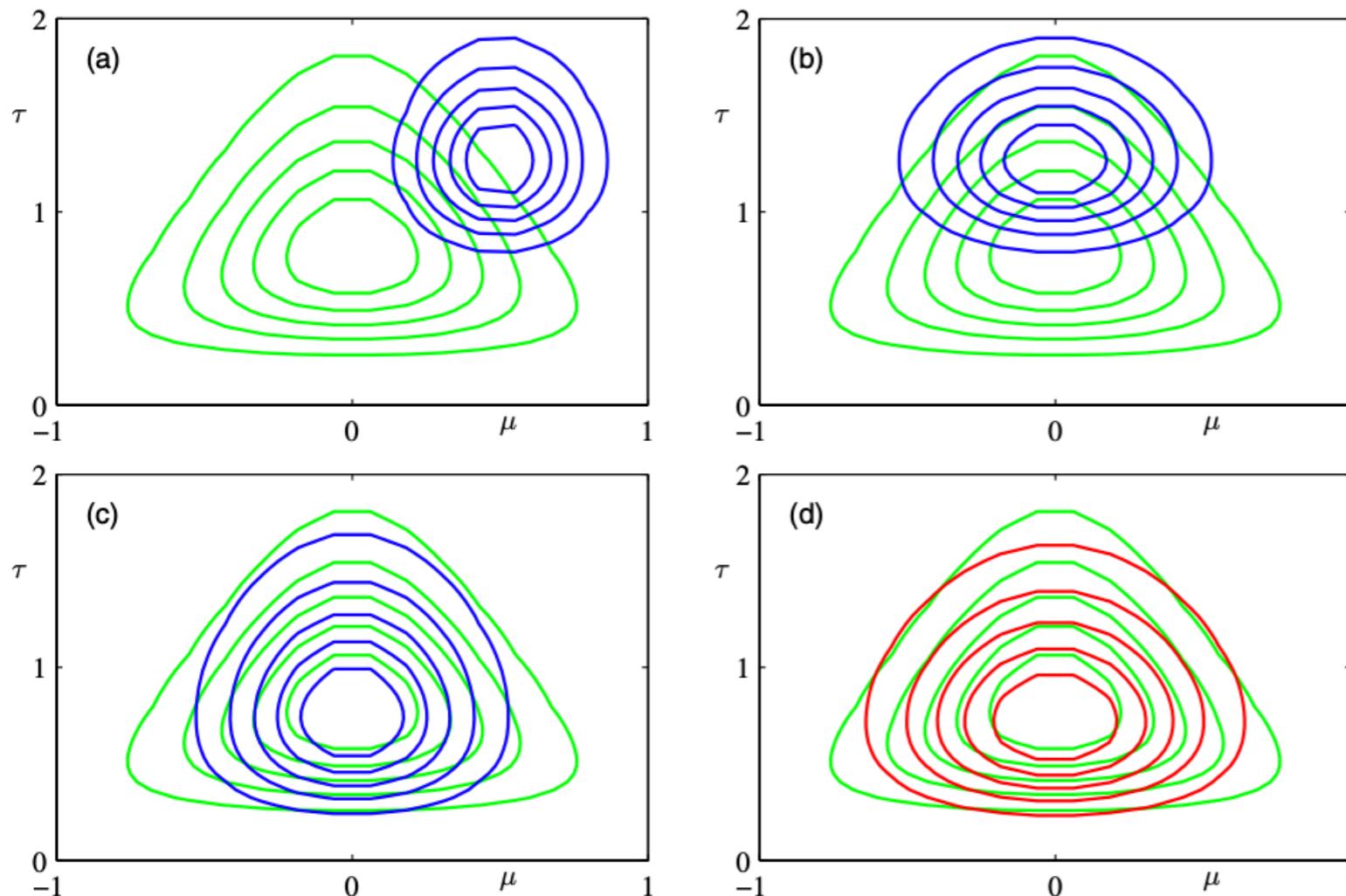
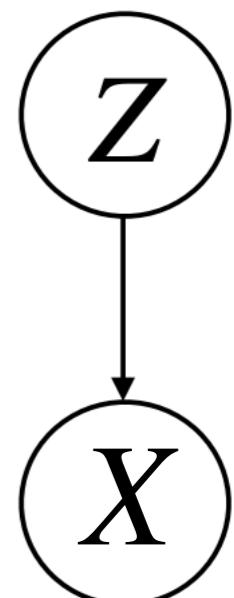


Figure 10.4 Illustration of variational inference for the mean μ and precision τ of a univariate Gaussian distribution. Contours of the true posterior distribution $p(\mu, \tau | \mathcal{D})$ are shown in green. (a) Contours of the initial factorized approximation $q_\mu(\mu)q_\tau(\tau)$ are shown in blue. (b) After re-estimating the factor $q_\mu(\mu)$. (c) After re-estimating factor $q_\tau(\tau)$. (d) Contours of the optimal factorized approximation, to which the iterative scheme converges, are shown in red.



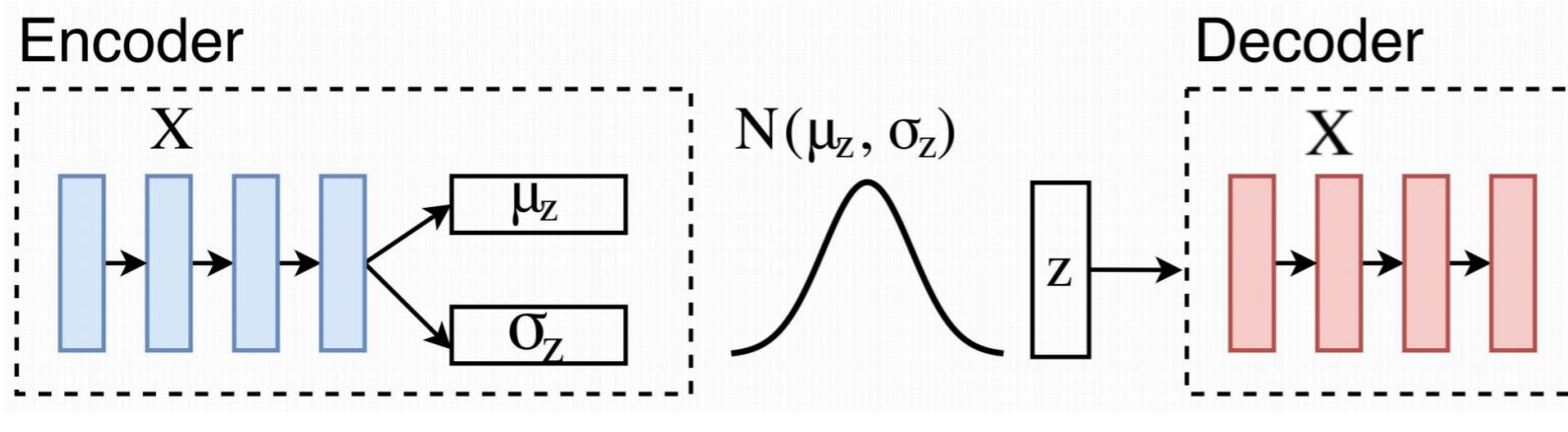
Variational Autoencoder

- Variational autoencoder
 - Variational family: $\mathcal{Q} = \{\mathcal{N}(\boldsymbol{\mu}, \text{diag } \boldsymbol{\sigma}^2) : \boldsymbol{\mu} \in \mathbb{R}^d, \boldsymbol{\sigma} \in \mathbb{R}_{++}^d\}$
 - Recognizing $\boldsymbol{\mu}, \boldsymbol{\sigma}$ by NN
 - Modeling x also by NN (need a little bit more efforts)



Variational Autoencoder

$$\begin{aligned}
 \log p_\theta(x) &= \log \left(\int_z p_\theta(x, z) dz \right) \\
 &= \int q_\phi(z|x) \log \frac{p_\theta(y, z)}{q_\phi(z|x)} dz + \int q_\phi(z|x) \log \frac{q_\phi(z|x)}{p_\theta(z|x)} dz \\
 &\geq \int q_\phi(z|x) \log \frac{p_\theta(y, z; \theta)}{q_\phi(z|x)} dz \\
 &= \int q_\phi(z|x) \log p_\theta(y|z) dz + \int q_\phi(z|x) \log \frac{p_\theta(z)}{q_\phi(z|x)} dz \\
 &= \mathbb{E}_{z \sim q_\phi(z|x)} \log p_\theta(y|z) - \text{KL}(q_\phi(z|x) \| p_\theta(z))
 \end{aligned}$$



Formula Zoo

$$\begin{aligned}
 & \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log p(x)] \geq \underbrace{\mathbb{E}_{x \sim p_{\text{data}}(x)} [\mathbb{E}_{q(z|x)} [\log p(x|z)]]}_{-\text{VAE Reconstruction}} - \underbrace{\mathbb{E}_{x \sim p_{\text{data}}(x)} [\text{KL}(q(z|x) || p(z))]}_{\text{VAE Regularization}} \quad (1) \\
 &= \underbrace{\mathbb{E}_{x \sim p_{\text{data}}(x)} [\mathbb{E}_{q(z|x)} [\log p(x|z)]]}_{-\text{AVB Reconstruction}} - \underbrace{\text{KL}(q(z,x) || p(z)p_{\text{data}}(x))}_{\text{AVB Regularization}} \quad (2) \\
 &- \underbrace{\mathbb{E}_{x \sim p_{\text{data}}(x)} [\mathbb{E}_{q(z|x)} [\log p(x|z)]]}_{-\text{AAE Reconstruction}} - \underbrace{\text{KL}(q(z) || p(z))}_{\text{AAE Regularization}} - \underbrace{\mathcal{I}(z;x)}_{\text{Mutual info.}} \quad (3) \\
 &= -\underbrace{\mathbb{E}_{z \sim q(z)} [\text{KL}(q(x|z) || p(x|z))]}_{\text{IAE Reconstruction}} - \underbrace{\text{KL}(q(z) || p(z))}_{\text{IAE Regularization}} - \underbrace{H_{\text{data}}(x)}_{\text{Entropy of data}} \quad (4) \\
 &= -\underbrace{\text{KL}(q(x,z) || r(x,z))}_{\text{IAE Reconstruction}} - \underbrace{\text{KL}(q(z) || p(z))}_{\text{IAE Regularization}} - \underbrace{H_{\text{data}}(x)}_{\text{Entropy of data}} \quad (5) \\
 &= -\underbrace{\text{KL}(q(x,z) || p(x,z))}_{\text{ALI/BiGAN Cost.}} - \underbrace{H_{\text{data}}(x)}_{\text{Entropy of data}} \quad (6)
 \end{aligned}$$

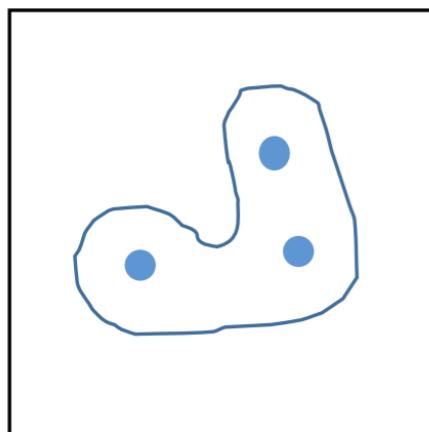


Adversarial/Wasserstein Autoencoder

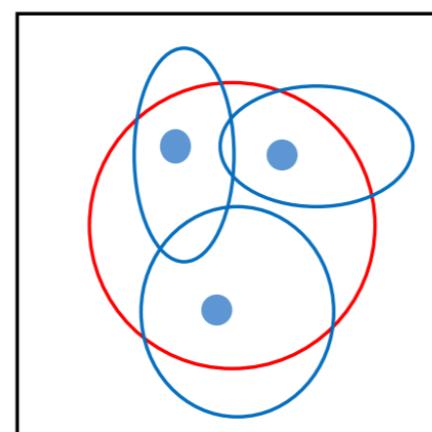
- VAE: $q(z|x) \rightarrow p(z)$

- WAE: $q(z) = \int p_{\mathcal{D}}(x)q(z|x)dx \xrightarrow{\text{close}} p(z)$

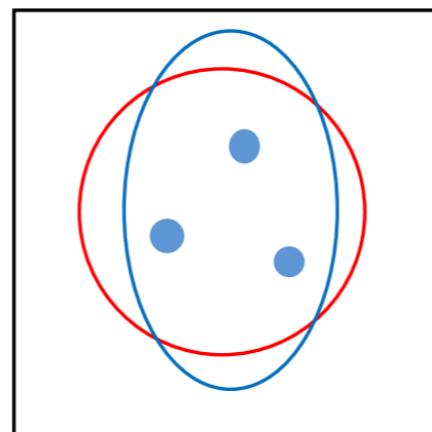
$$J = \mathbb{E}_{x \in p_{\mathcal{D}}(x)} \mathbb{E}_{z \in q(z|x)} \log p(z|x) + \mathbb{D}(q(z), p(z))$$



(a) DAE



(b) VAE



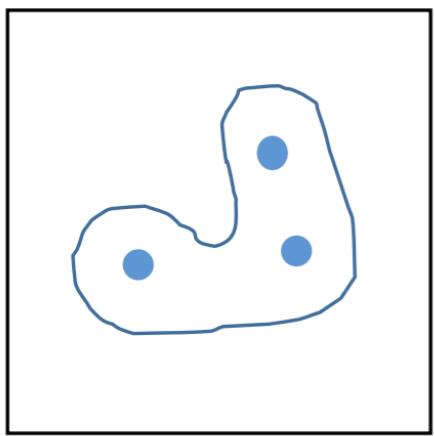
(c) WAE

Implicit Distributions

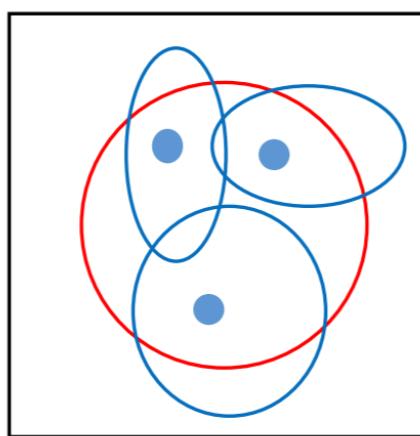
- We penalize some distance between $q(z)$ and $p(x)$
- We do not have an explicit form or $q(z)$

- But samples from $q(z) := \int p_{\mathcal{D}}(x)q(z|x)dx$

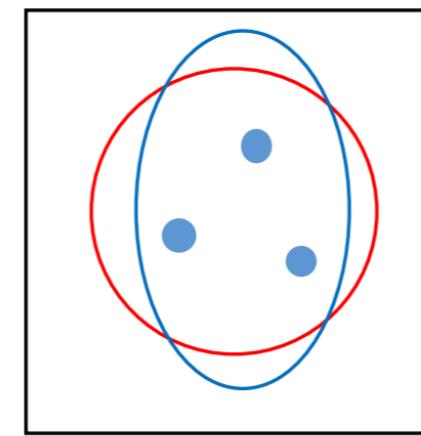
$$x^{(i)} \sim p_{\mathcal{D}}(x), \quad z \sim q(z|x^{(i)})$$



(a) DAE



(b) VAE



(c) WAE



Adversarial Training

- We penalize some distance between $q(z)$ and $p(z)$
- We deliberately introduce a classifier (discriminator/adversary) to distinguish samples from $q(z)$ and $p(z)$
- We train the model to fool the discriminator
 - Flipping gradient
 - Maximizing predicted entropy

Algorithm

foreach *mini-batch* **do**

minimize $J_{\text{dis}}(\theta_{\text{dis}})$ w.r.t. θ_{dis}

minimize J_{ovr} w.r.t. θ_E, θ_D

end

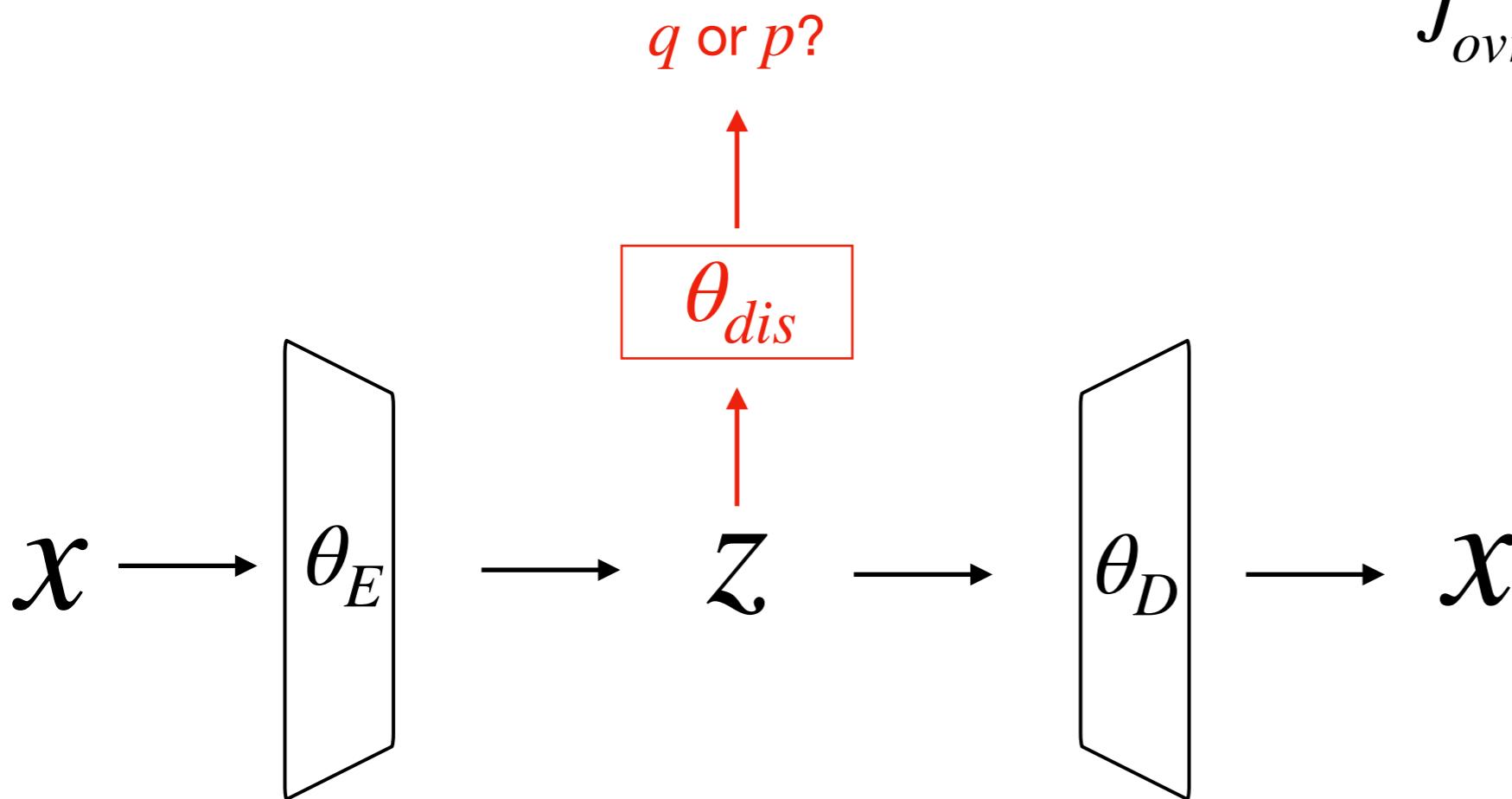
- Modeling J_{ovr}

- Flipping gradient

$$J_{\text{ovr}} = J_{\text{rec}} - J_{\text{dis}}$$

- Maximizing predicted entropy

$$J_{\text{ovr}} = J_{\text{rec}} - \mathcal{H}(y_{\text{dis}})$$



Applications of Adv Training

- Adversarial/Wasserstein autoencoder

$$q(z) \quad vs \quad p(z)$$

- Generative adversarial network

$$p_{gen}(z) \quad vs \quad p_{\mathcal{D}}(z)$$

- Domain adaptation

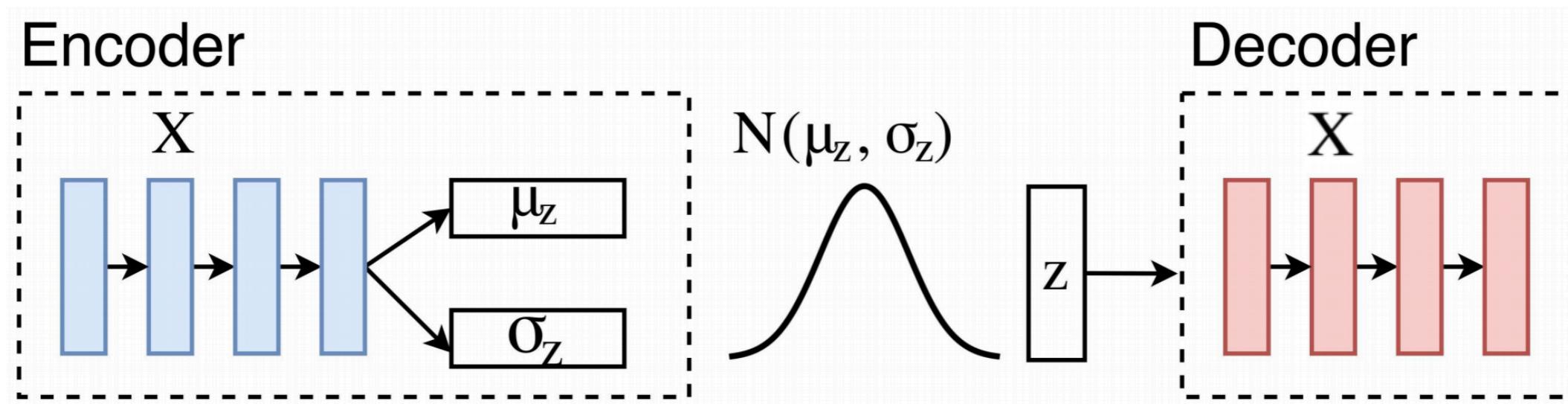
$$p_{D1}(z) \quad vs \quad p_{D2}(z)$$

A Few Fundamental Questions

- VAE: KL collapse

$$J = \mathbb{E}_{z \sim q(z|x)}[\log p(x|z)] + \text{KL}(q(z|x) \| p(z))$$

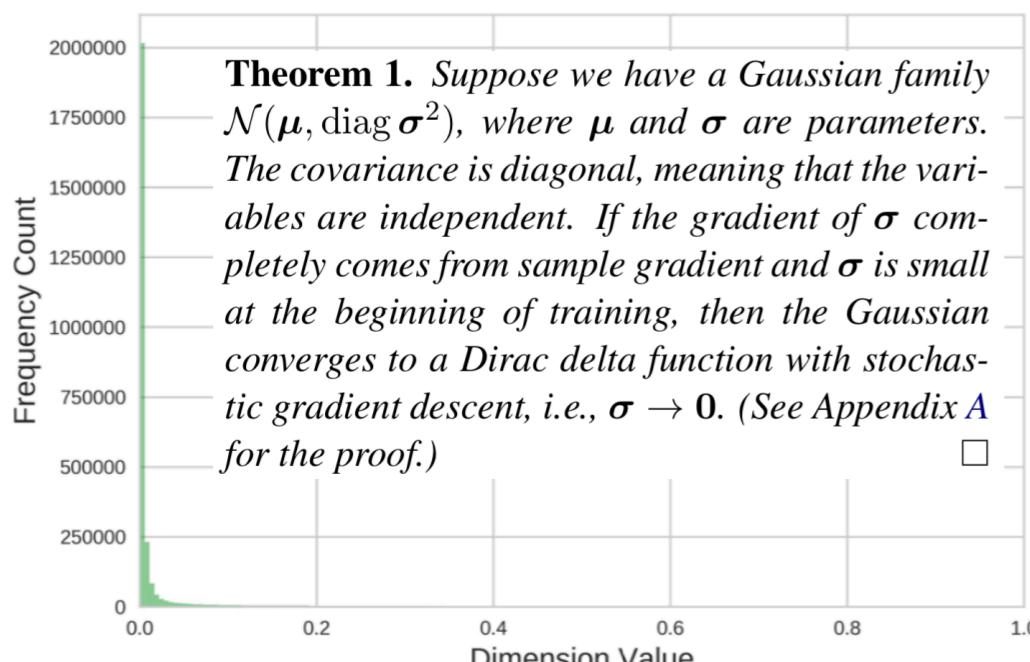
- WAE alleviates this problem



A Few Fundamental Questions

- WAE: Stochasticity collapse

$$J = \mathbb{E}_{z \sim q(z|x)}[\log p(x|z)] + \mathbb{D}(q(z), p(z))$$



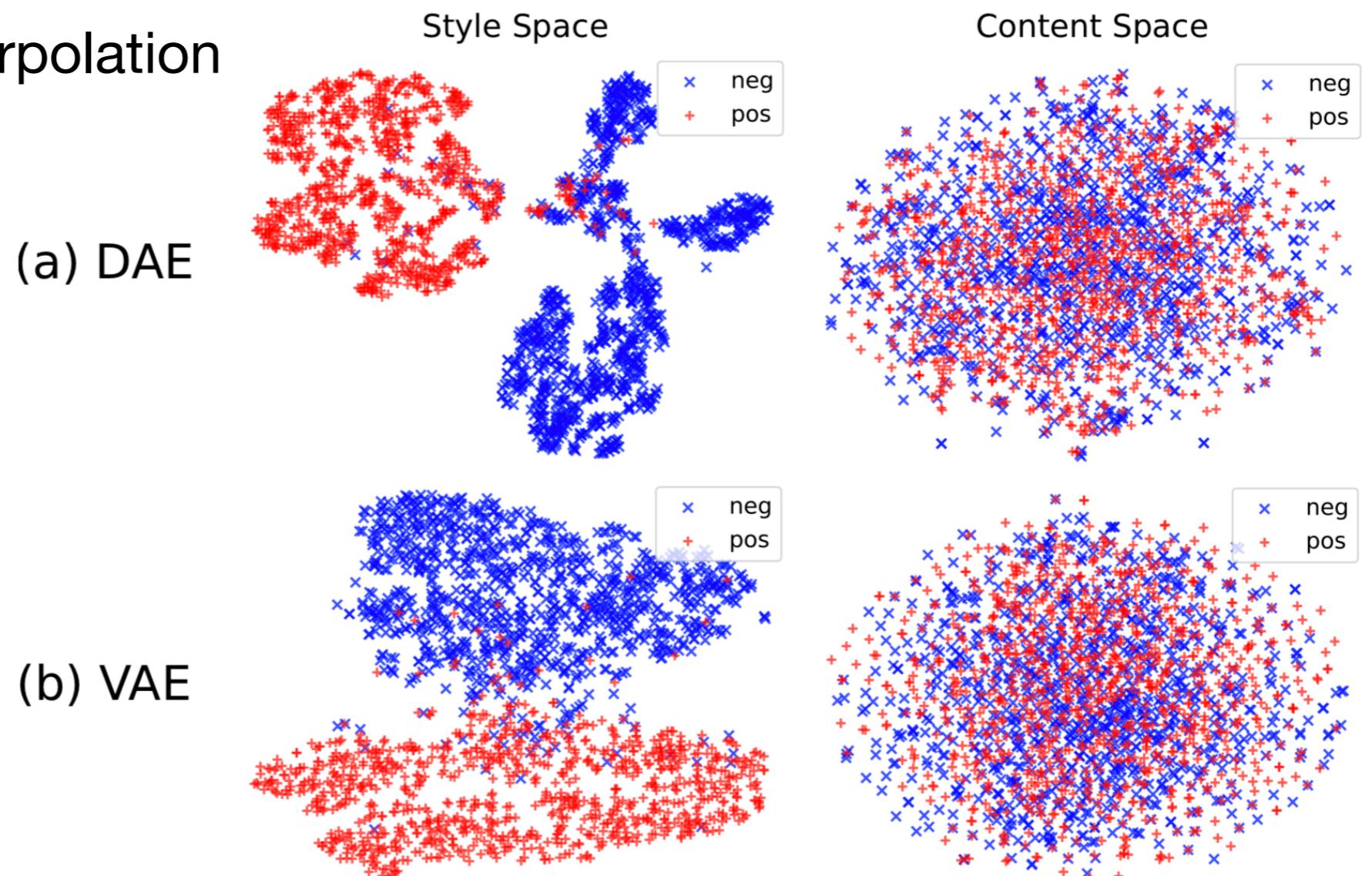
(a) $\lambda_{\text{KL}} = 0$

- If
 - Gaussian encoder
 - Gradient comes from samples
 - Sampling var \ll Data var
- Then
 - $\sigma^2 \rightarrow 0$ by SGD

Applications of *AEs

- Regularization

Especially good for interpolation





Applications of *AEs

- Prior sampling

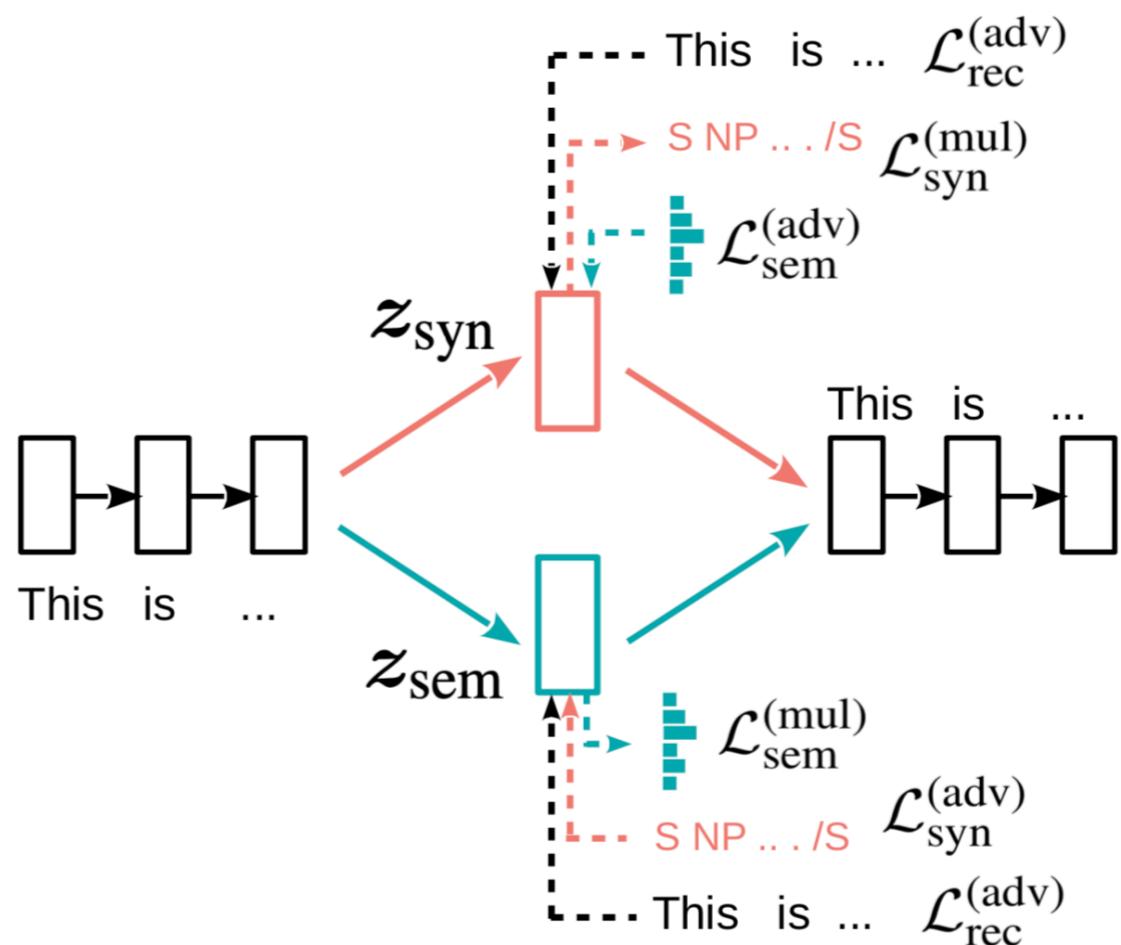
WAE-D ($\lambda_{\text{WAE}} = 10$)
<i>the lone man is working .</i>
<i>the group of men is using ice at the sunset .</i>
<i>a family is outside in the background .</i>
<i>two women are standing on a busy street outside a fair</i>
<i>a tourists is having fun on a sunny day</i>

WAE-S ($\lambda_{\text{WAE}} = 10, \lambda_{\text{KL}} = 0.01$)
<i>an asian man is dancing in a highland house .</i>
<i>a person wearing a purple snowsuit jumps over the tree .</i>
<i>the vocalist is at the music and dancing with a microphone .</i>
<i>a young man is dressed in a white shirt cleaning clothes .</i>
<i>three children lie together and a woman falls in a plane .</i>

Training Samples
<i>a mother and her child are outdoors.</i>
<i>the people are opening presents.</i>
<i>the girls are looking toward the water.</i>
<i>a small boy walks down a wooden path in the woods.</i>
<i>a person in a green jacket it surfing while holding on to a line.</i>
DAE
<i>two families walking in a towel down alaska sands a cot .</i>
<i>a blade is rolling its nose furiously paper .</i>
<i>a woman in blue shirts is passing by a some beach</i>
<i>transporting his child are wearing overalls .</i>
<i>a guys are blowing on professional thinks the horse .</i>
VAE without Annealing
<i>a man is playing a guitar .</i>
<i>a man is playing with a dog .</i>
<i>a man is playing with a dog .</i>
<i>a man is playing a guitar .</i>
<i>a man is playing with a dog .</i>
VAE with Annealing
<i>the band is sitting on the main street .</i>
<i>couple dance on stage in a crowded room .</i>
<i>two people run alone in an empty field .</i>
<i>the group of people have gathered in a picture .</i>
<i>a cruise ship is docking a boat ship .</i>

Applications of *AEs

- Posterior sampling
 - Paraphrase generation
 - Style-transfer generation



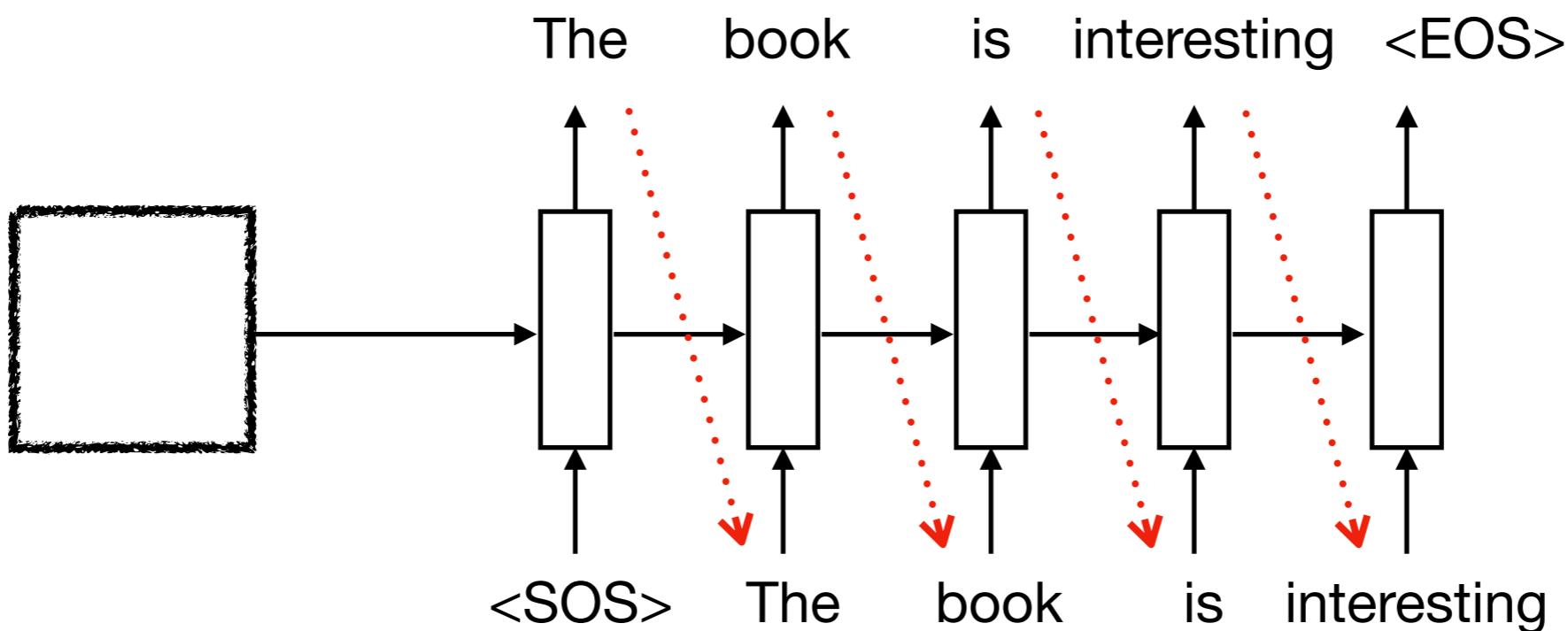
Semantic and Syntactic Providers	Syntax-Transfer Output
Ref_{syn}: There is an apple on the table.	VAE: The man is in the kitchen.
Ref_{sem}: The airplane is in the sky.	DSS-VAE: There is a airplane in the sky.
Ref_{syn}: The shellfish was cooked in a wok.	VAE: The man was filled with people.
Ref_{sem}: The stadium was packed with people.	DSS-VAE: The stadium was packed with people.
Ref_{syn}: The child is playing in the garden.	VAE: There is a person in the garden.
Ref_{sem}: There is a dog behind the door.	DSS-VAE: A dog is walking behind the door.



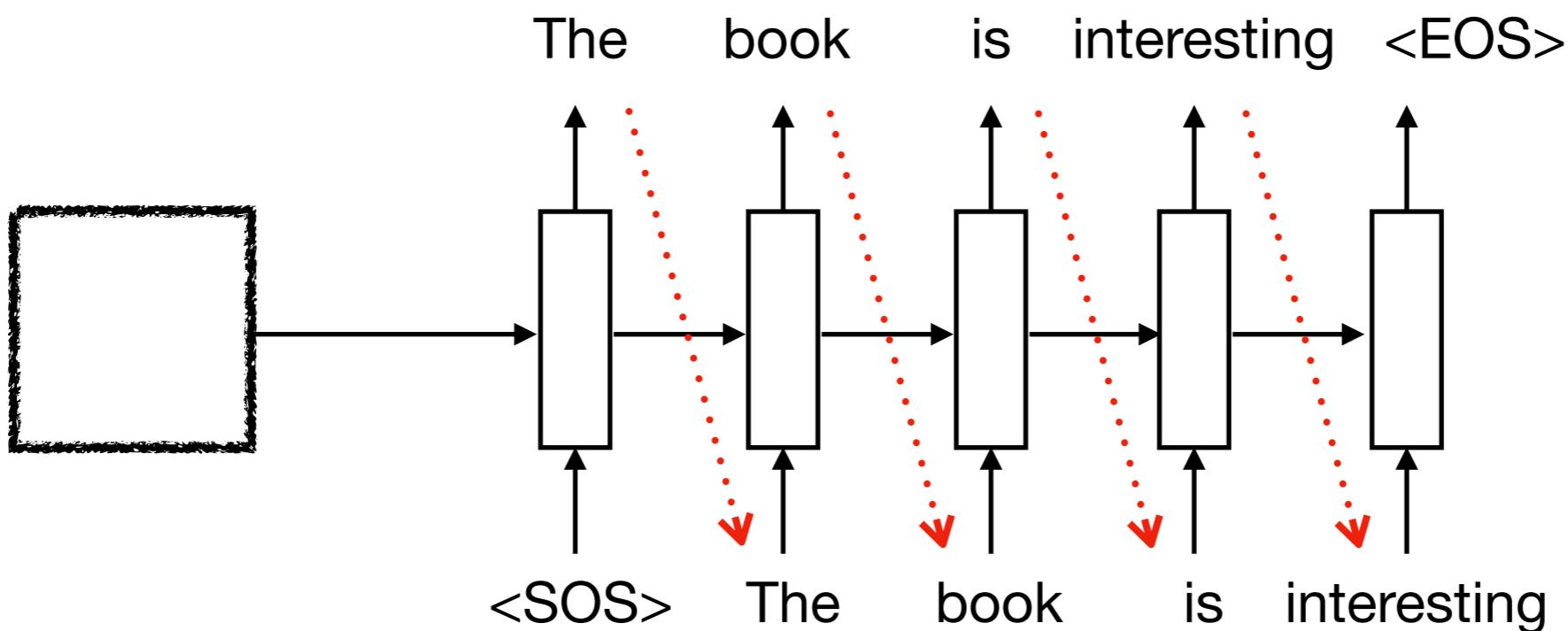
Word Space Sampling

Miao, Zhou, Mou, Yan, Li. CGMH: Constrained sentence generation by Metropolis-Hastings sampling. In AAAI, 2019.

RNN Generation



RNN Generation



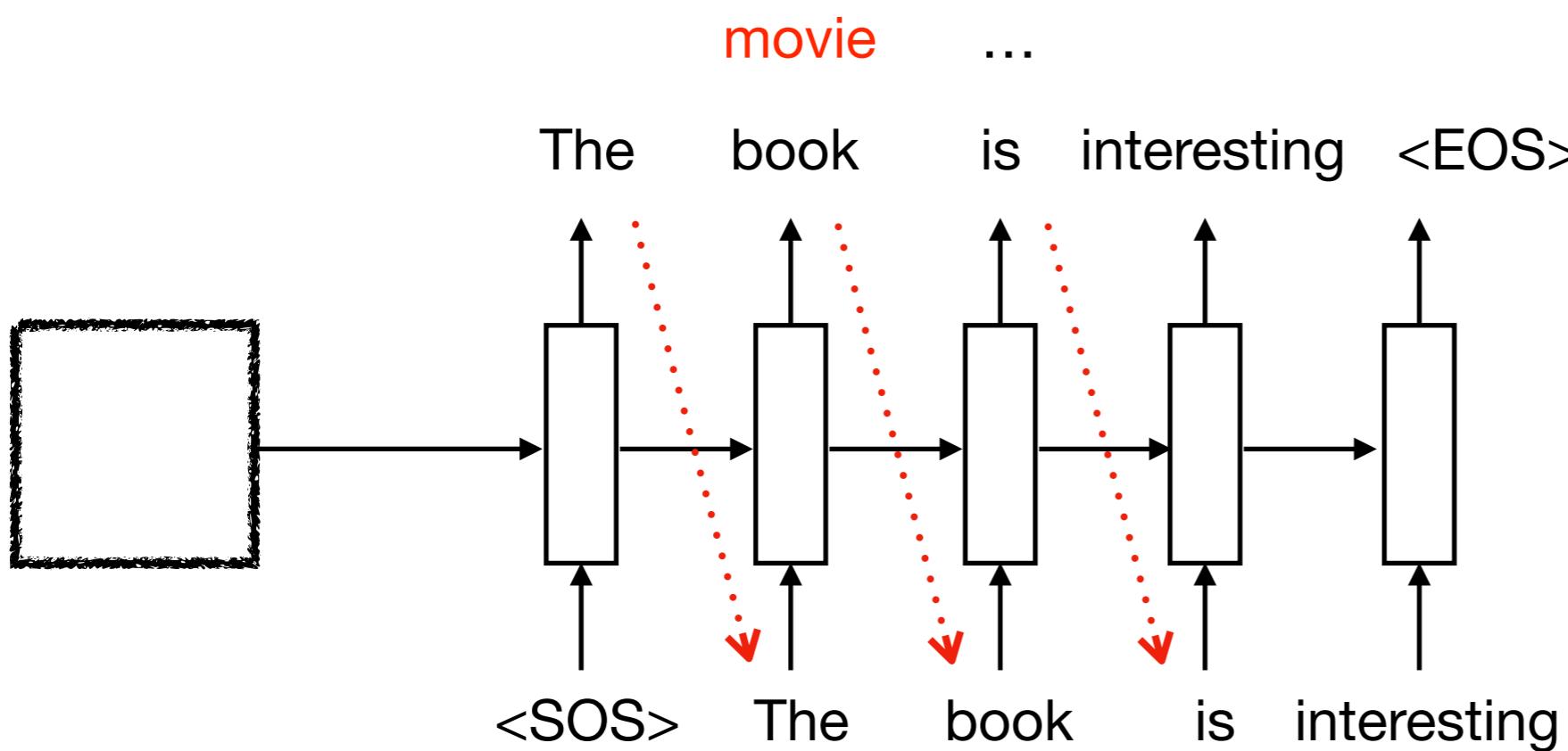
Question: Can we generate a sentence right-to-left?



UNIVERSITY OF
ALBERTA

Issues with Single Directional Generation

- Information bottleneck
- Error cumulation
 - Due to sampling or incompetency of the RNN



Generation by Local Changes

- Suppose we have a blueprint

The book is interesting <EOS>

Generation by Local Changes

- Suppose we have a blueprint

The book is interesting <EOS>

~~The~~ book is interesting <EOS>
This

Generation by Local Changes

- Suppose we have a blueprint

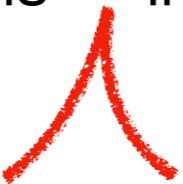
The book is interesting <EOS>

~~The~~ book is interesting <EOS>
This
quite

Generation by Local Changes

- Suppose we have a blueprint

The book is interesting <EOS>

~~The~~ book is ~~interesting~~ <EOS>
This 
quite

Applications

- Paraphrase generation
 - “Sample” a sentence with similar semantics but different wordings
- Summarization
 - “Sample” a sentence with similar semantics
- Grammatical error correction
 - “Sample” a more likely sentence with the same semantics

Sampling Methods



**UNIVERSITY OF
ALBERTA**

Independent Sampling

- Sampling from CDF

- Probabilistic density function (PDF)

$$\Pr[a \leq x \leq b] = \int_a^b f(x) dx$$

- Cumulative density function (CDF)

$$F(x) = \int_{-\infty}^x f(u) du = \Pr[u \leq x]$$

- Sampling procedure

$$u \sim U[0,1]; \quad x = \text{CDF}^{-1}(u)$$

- Problems

- CDF not analytic
 - Especially, the conditional CDF in multivariate cases

Independent Sampling

- Rejection sampling

- To sample from

$$p(x) = \frac{1}{Z} \tilde{p}(x)$$

- We instead sample

$$x \sim q(x)$$

- Accept the sample x with probability

$$\frac{\tilde{p}(x)}{k \cdot q(x)}$$

where k is a constant s.t. $kq(x) \geq \tilde{p}(x), \forall x$

- Reject x w.p. $1 - \frac{\tilde{p}(x)}{k \cdot q(x)}$

- Many other sampling methods

Dependent Sampling

- Goal: Sample from $p(x)$
- MCMC sampling
 - Start from an arbitrary initial sample $x^{(0)}$
 - Sample $x^{(1)} \sim p(x^{(1)} | x^{(0)})$, $x^{(2)} \sim p(x^{(2)} | x^{(1)})$, ...
 - Hope $p(x^{(n)}) \rightarrow p(x)$ as $n \rightarrow \infty$

Markov Chain

- States: $S = \{s_1, s_2, \dots\}$
- Initial distribution $\pi^{(0)}$
- Transition probability: $\mathcal{T}_{i \rightarrow j} = p(x^{(t+1)} = s_j | x^{(t)} = s_i)$
 - $x^{(t+1)}$ is independent of $x^{(t-1)}$, given $x^{(t)}$
 - $\mathcal{T}_{i \rightarrow j}$ works for all time steps t
- **Thm:** Starting from an arbitrary initial distribution, a Markov Chain converges to a **unique** stationary distribution (under mild assumptions).

Markov Chain Monte Carlo

- Goal: Sample from $p(x)$
- MCMC sampling
 - Start from an arbitrary initial sample $x^{(0)}$
 - Sample $x^{(1)} \sim p(x^{(1)} | x^{(0)})$, $x^{(2)} \sim p(x^{(2)} | x^{(1)})$, ...
 - Hope $p(x^{(n)}) \rightarrow p(x)$ as $n \rightarrow \infty$

Markov Chain Monte Carlo

- Goal: Sample from $p(x)$
- MCMC sampling
 - Start from an arbitrary initial sample $x^{(0)}$
 - Sample $x^{(1)} \sim p(x^{(1)} | x^{(0)})$, $x^{(2)} \sim p(x^{(2)} | x^{(1)})$, ...
by following a carefully designed Markov chain
 - ~~Hope~~ $p(x^{(n)}) \rightarrow p(x)$ as $n \rightarrow \infty$

Guaranteed that



UNIVERSITY OF
ALBERTA

Metropolis–Hastings Sampler

- **Input**
 - An **arbitrary** desired distribution $p(x)$
- **Output**
 - An unbiased sample $x \sim p(x)$
- **Algorithm**
 - Start from an **arbitrary** initial state $x^{(0)}$
 - For every step t
 - Propose a new state $x' \sim g(x' | x^{(t)})$
 - Accept x' w.p. $A(x'|x) = \min \left\{ 1, \frac{p(x')g(x^{(t)}|x')}{p(x)g(x'|x^{(t)})} \right\}$, i.e., $x^{(t+1)} = x'$
 - Reject x' otherwise, i.e., $x^{(t+1)} = x^{(t)}$
 - Return $x^{(t)}$ with a large t



UNIVERSITY OF
ALBERTA

Proof Sketch

- Detailed balance property = > Stationary distribution

If

$$\forall x, y, \quad \pi(x) \cdot \mathcal{T}_{x \rightarrow y} = \pi(y) \cdot \mathcal{T}_{y \rightarrow x}$$

Then

$\pi(x)$ is a stationary distribution

Because

$$\forall x, \quad \pi(x) = \sum_y \pi(y) \mathcal{T}_{y \rightarrow x} = \sum_y \pi(x) \mathcal{T}_{x \rightarrow y} = \pi(x)$$



UNIVERSITY OF
ALBERTA

Proof Sketch (Cont.)

- MH Sampler satisfies detailed balance

- $\forall x, y, \text{ if } x \neq y, p(x) \cdot \mathcal{T}_{x \rightarrow y} = p(x) \cdot g(y|x) \cdot \min \left\{ 1, \frac{p(y)g(x|y)}{p(x)g(y|x)} \right\}$ (1)

$$p(y) \cdot \mathcal{T}_{y \rightarrow x} = p(y) \cdot g(x|y) \cdot \min \left\{ 1, \frac{p(x)g(y|x)}{p(y)g(x|y)} \right\} \quad (2)$$

- W.L.O.G., we assume $p(x)g(y|x) \geq p(y)g(x|y)$

$$(1) = p(y) \cdot g(x|y)$$

$$(2) = p(y) \cdot g(x|y)$$

- $\forall x, y, \text{ if } x = y, p(x)\mathcal{T}_{x \rightarrow y} = p(y)\mathcal{T}_{y \rightarrow x}$ also holds

Gibbs Sampler

- Suppose $\mathbf{x} = (x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$
 - If the proposal distribution is $x'_i \sim p(x_i | \mathbf{x}_{-i})$
 - Then, the acceptance rate is $A(\mathbf{x}' | \mathbf{x}) = \min \left\{ 1, \frac{p(\mathbf{x}')g(\mathbf{x} | \mathbf{x}')}{p(\mathbf{x})g(\mathbf{x}' | \mathbf{x})} \right\}$
 - Notice that $\mathbf{x}' = (x_1, x_2, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)$
 - Thus, $\frac{p(\mathbf{x}')g(\mathbf{x} | \mathbf{x}')}{p(\mathbf{x})g(\mathbf{x}' | \mathbf{x})} = \frac{p(\mathbf{x}_{-i})p(x'_i | \mathbf{x}_{-i}) \cdot p(x_i | \mathbf{x}_{-i})}{p(\mathbf{x}_{-i})p(x_i | \mathbf{x}_{-i}) \cdot p(x'_i | \mathbf{x}_{-i})} = 1$
- => Gibbs step is a special case of an MH step, with AC rate = 1

Applying MH to Sentence Generation



UNIVERSITY OF
ALBERTA

MH Components

- State: Every sentence
- Target distribution: Depend on the task
- Proposal distribution
 - Task agnostic, or task specific
- Compute acceptance rate
 - We can't do anything here

Target distribution

- General formula
 - $p(\mathbf{x}) \propto p_{\text{LM}}(\mathbf{x}) \cdot s_1(\mathbf{x}) \cdots s_n(\mathbf{x})$
 - $s_i(\mathbf{x})$: scoring functions specific to the task

Target distribution

- General formula

- $p(\mathbf{x}) \propto p_{\text{LM}}(\mathbf{x}) \cdot s_1(\mathbf{x}) \cdots s_n(\mathbf{x})$

- $s_i(\mathbf{x})$: scoring functions specific to the task

- Keywords-to-sentence generation

$$s(\mathbf{x}) = \begin{cases} 1, & \text{if keywords in } \mathbf{x} \\ 0, & \text{otherwise} \end{cases}$$

- Paraphrase generation/Grammatical error correction

- $s(\mathbf{x}) = \text{sim}_{\text{semantic}}(\mathbf{x}, \mathbf{x}_0) + \text{diff}_{\text{word}}(\mathbf{x}, \mathbf{x}_0)$



Proposal Distribution

- Replace

$$g_{\text{replace}}(\mathbf{x}' | \mathbf{x}) = \pi(w_m^* = w^c | \mathbf{x}_{-m}) = \frac{\pi(w_1, \dots, w_{m-1}, w^c, w_{m+1}, \dots, w_n)}{\sum_{w \in \mathcal{V}} \pi(w_1, \dots, w_{m-1}, w, w_{m+1}, \dots, w_n)}$$

- Delete
- Insert
 - Also sample from posterior

Examples: Keywords-to-Sentence

Keyword(s)	Generated Sentences
friends	My good friends were in danger .
project	The first project of the scheme .
have, trip	But many people have never made the trip .
lottery, scholarships	But the lottery has provided scholarships .
decision, build, home	The decision is to build a new home .
attempt, copy, painting, denounced	The first attempt to copy the painting was denounced .



Examples: Paraphrase Generation

Model	BLEU-ref	BLEU-ori	NLL
Origin Sentence	30.49	100.00	7.73
VAE-SVG (100k)	22.50	-	-
VAE-SVG-eq (100k)	22.90	-	-
VAE-SVG (50k)	17.10	-	-
VAE-SVG-eq (50k)	17.40	-	-
Seq2seq (100k)	22.79	33.83	6.37
Seq2seq (50k)	20.18	27.59	6.71
Seq2seq (20k)	16.77	22.44	6.67
VAE (unsupervised)	9.25	27.23	7.74
CGMH <i>w/o matching</i>	18.85	50.28	7.52
<i>w/ KW</i>	20.17	53.15	7.57
<i>w/ KW + WVA</i>	20.41	53.64	7.57
<i>w/ KW + WVM</i>	20.89	54.96	7.46
<i>w/ KW + ST</i>	20.70	54.50	7.78

Type	Examples
Ori	what 's the best plan to lose weight
Ref	what is a good diet to lose weight
Gen	what 's the best way to slim down quickly
Ori	how should i control my emotion
Ref	how do i control anger and impulsive emotions
Gen	how do i control my anger
Ori	why do my dogs love to eat tuna fish
Ref	why do my dogs love eating tuna fish
Gen	why do some dogs like to eat raw tuna and raw fish

Examples: Paraphrase Generation

Model	#parallel data	GLEU
AMU	2.3M	44.85
CAMB-14	155k	46.04
MLE	720k	52.75
NRL	720k	53.98
CGMH	0	45.5

Ori	Even if we are failed , We have to try to get a new things .
Ref	Even if we all failed , we have to try to get new things .
Gen	Even if we are failing , We have to try to get some new things .
Ori	In the world oil price very high right now .
Ref	In today 's world , oil prices are very high right now .
Gen	In the world , oil prices are very high right now .



Analysis

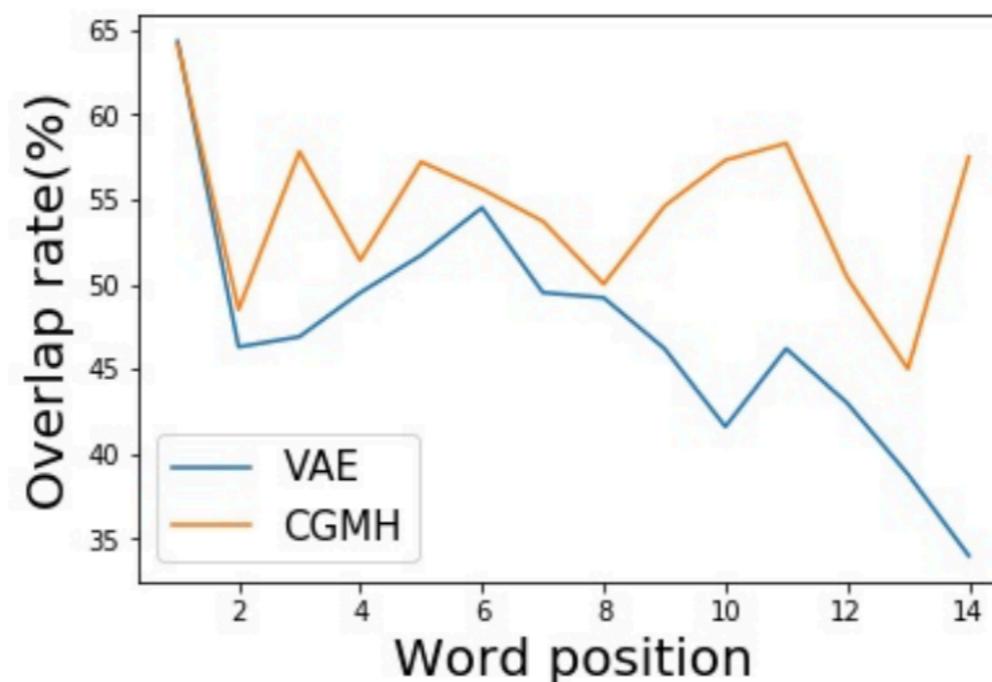
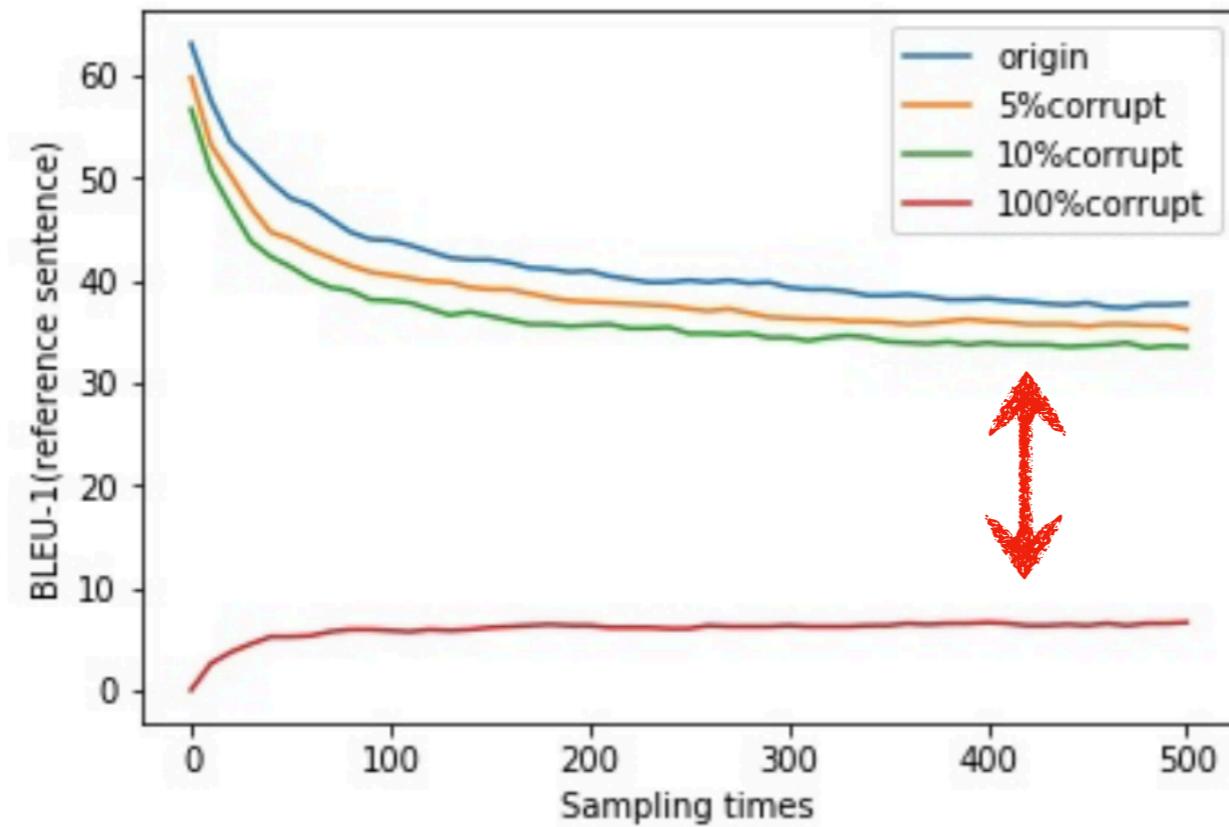


Figure 3: Overlap rates of CGMH and VAE for each word position of sentences.

Analysis (Cont.)



**The Markov Chain never mixes.
We mainly use MH as SA.**

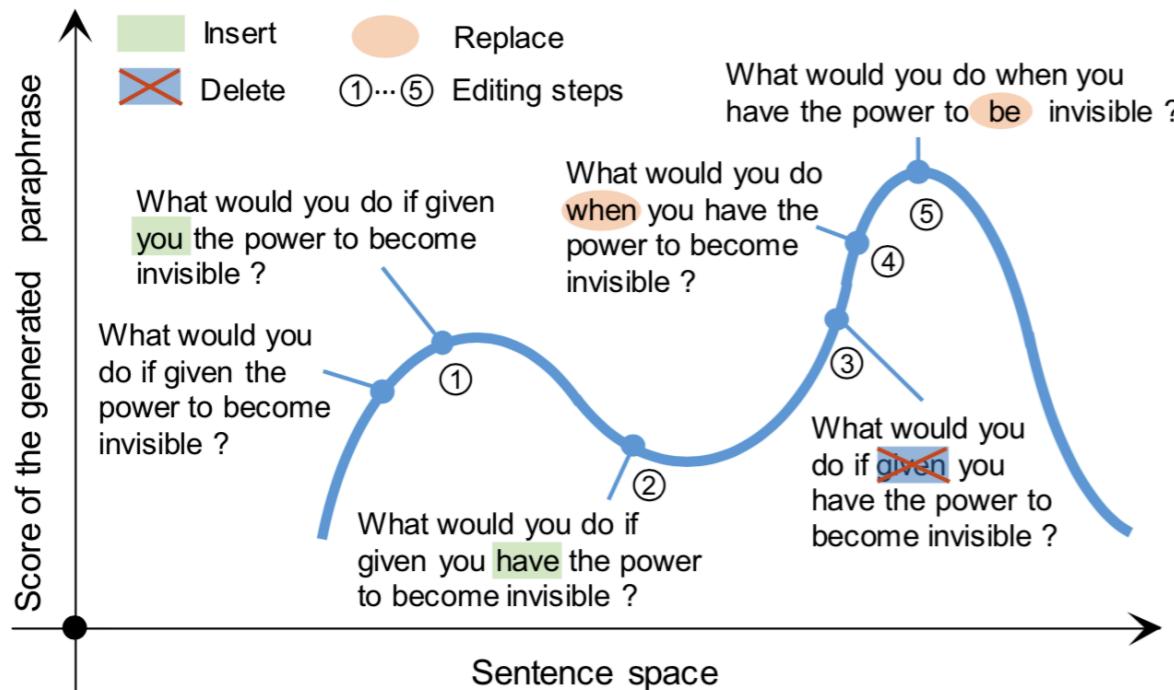
Figure 2: Generation quality with corrupted initial states. At each situation, 0/5%/10%/100% of the words in initial sentences are randomly replaced with other words.





Simulated Annealing

- MH: Sampling from $\propto \exp\{s(x)\}$
- SA: Searching the optimum of $s(x)$
 - Define $p_\tau(x) \propto \exp\{s(x)/\tau\}$
 - Start from high temperature, but cool it down gradually
 - With $\tau \rightarrow 0$,
$$p_\tau(x) = 1 \text{ if } x = \operatorname{argmax} s(x), \text{ or } 0 \text{ otherwise}$$



$$p(\text{accept} | x_*, x_t, T) = \min(1, e^{\frac{f(x_*) - f(x_t)}{T}})$$

$$T = \max(0, T_{\text{init}} - C \cdot t)$$

References

- Kingma, D. P., & Welling, M. Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114, 2013.
- Makhzani A, Shlens J, Jaitly N, Goodfellow I, Frey B. Adversarial autoencoders. arXiv preprint arXiv:1511.05644. 2015.
- Bowman SR, Vilnis L, Vinyals O, Dai AM, Jozefowicz R, Bengio S. Generating sentences from a continuous space. In *CoNLL*, 2016.
- Bahuleyan, Mou, Zhou, and Vechtomova. Stochastic Wasserstein autoencoder for probabilistic sentence generation. In *NAACL-HLT*, 2019.
- Miao, Zhou, Mou, Yan, Li. CGMH: Constrained sentence generation by Metropolis-Hastings sampling. In *AAAI*, 2019.
- Liu, Mou, Meng, Zhou, Zhou, Song. Unsupervised Paraphrasing by Simulated Annealing. arXiv preprint arXiv:1909.03588, 2019.

Thank you!

Q&A



**UNIVERSITY OF
ALBERTA**