

# Towards Controlled Generation of Text

Zhiting Hu   Zichao Yang   Xiaodan Lang   **Ruslan Salakhutdinov**   Eric P.Xing  
Carnegie Mellon University

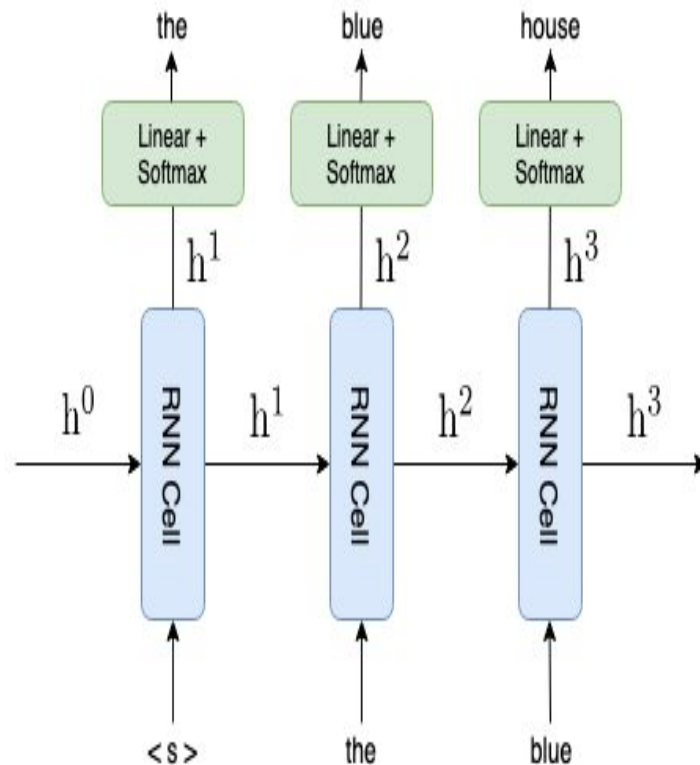
ICML 2017

**Presenter:** Anchit Bhattacharya  
Arizona State University

# Text Generation Methods

## RNN based text generation -

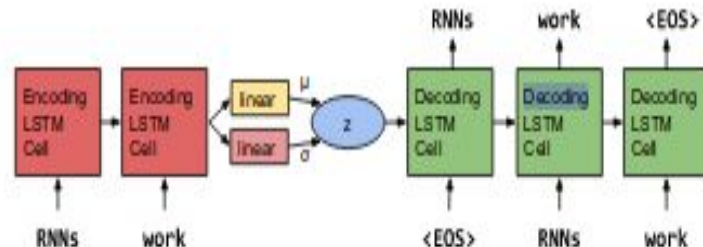
- Train a RNN based language model and use it to generate the next token based on all previous tokens.
- Can be character level or word level language model



# Text Generation Methods (Continued)

## VAE based text generation -

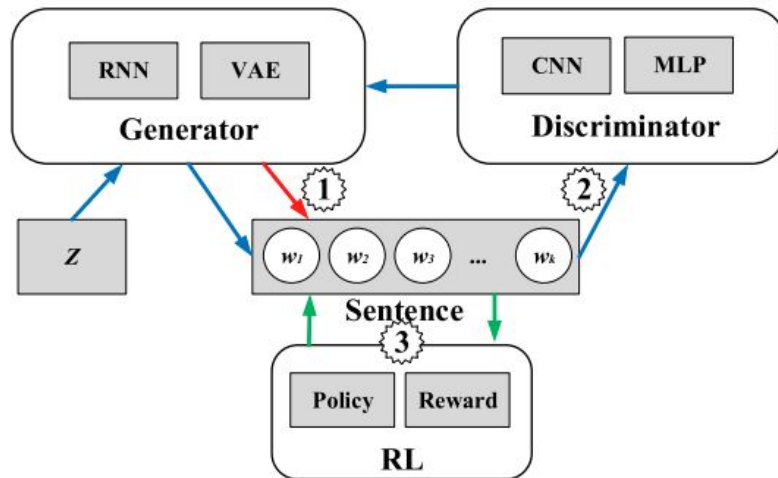
- Incorporates distributed latent representations of entire sentences.
- Explicitly model holistic properties of sentences.



# Text Generation Methods (Continued)

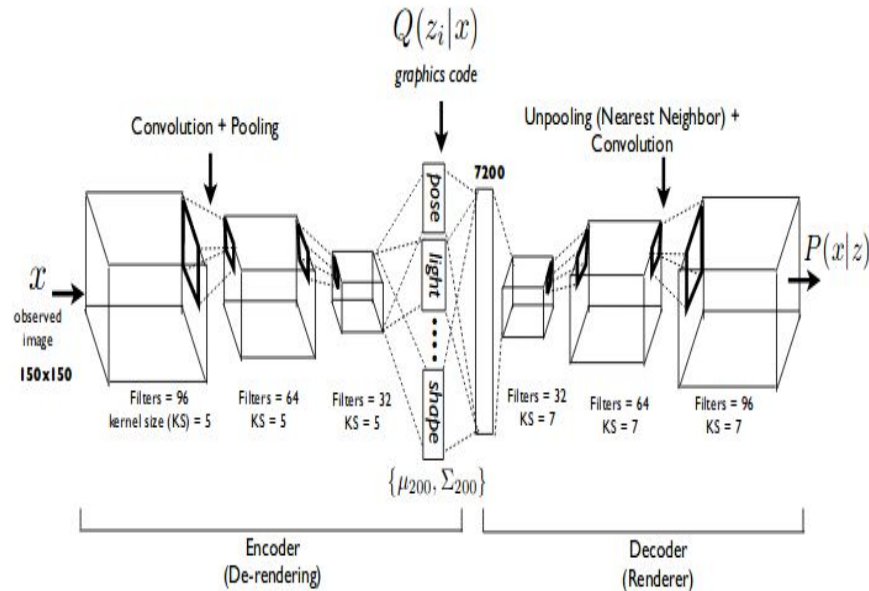
## GAN based text generation -

- Adversarial training to improve generated text.
- Can be used as a training method on top of previous models.



# Disentangled Representations

- Neurons in the neural network are somehow learning complete concepts alone.
- We want the latent representation of the texts to be disentangled, so that we can control the generation effectively.



# Challenges for Controllable Text Generation

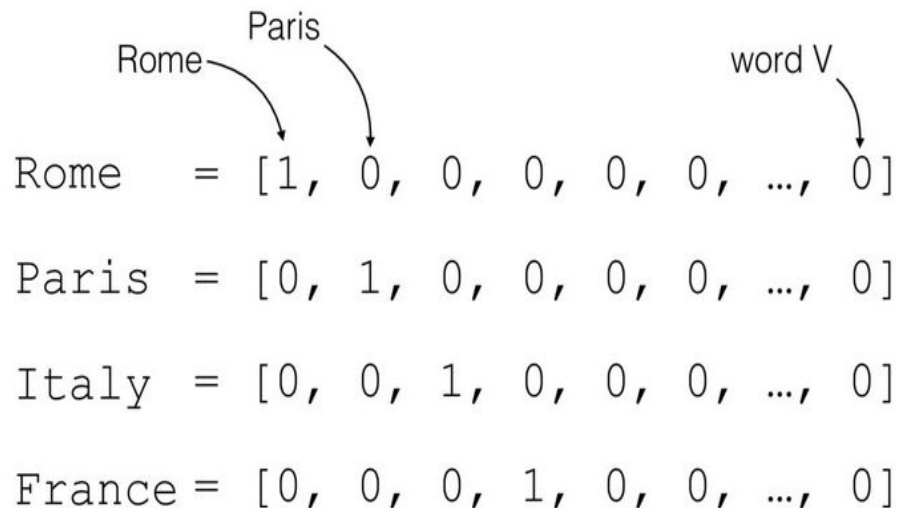
## 1> **Discrete Nature of Text Samples** -

Difficult to do backpropagation and train the generator in adversarial based methods.

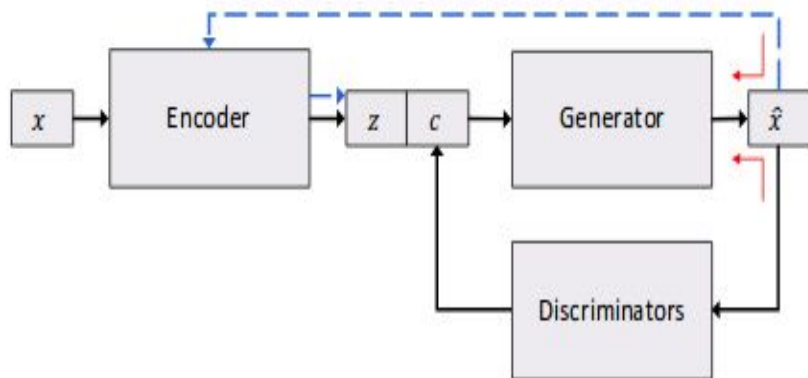
## 2> **Learning Disentangled representations** -

Lacks explicit enforcement of independence property on full latent representation.

**3> Controlled Generation of Text** - Control the generation of text based on some properties generically. Byproduct of point 2.



# Model



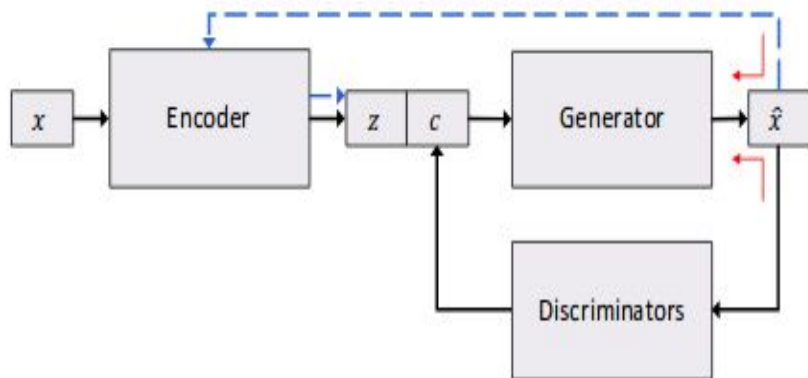
## Discrete Nature(Solution)

- Continuous Approximation based on softmax with a decreasing temperature( $T$ ).

$$x = \text{softmax}(o \mid T)$$

- Low variance and fast convergence as opposed to Policy Gradient based method.

# Model



## Controlled Generation Solution

- Add a structured code to the latent representation for every attribute that needs to be controlled.

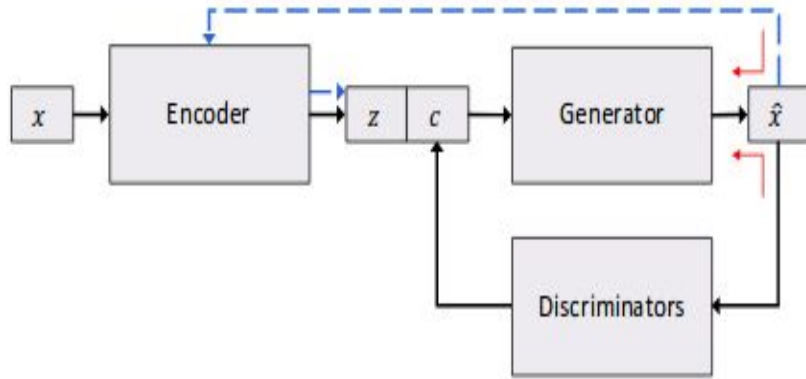
***$z$  - Gaussian prior***

***$c$  - Continuous/Discrete***

- Different discriminator for each attribute



# Model



## Generator Training(Part1)

- Train parameters of Encoder( $\theta_E$ ) and Generator( $\theta_G$ ) based on the reconstruction Error of real sentences.
- Also, train the encoder so that the latent representation is close to a Gaussian prior.

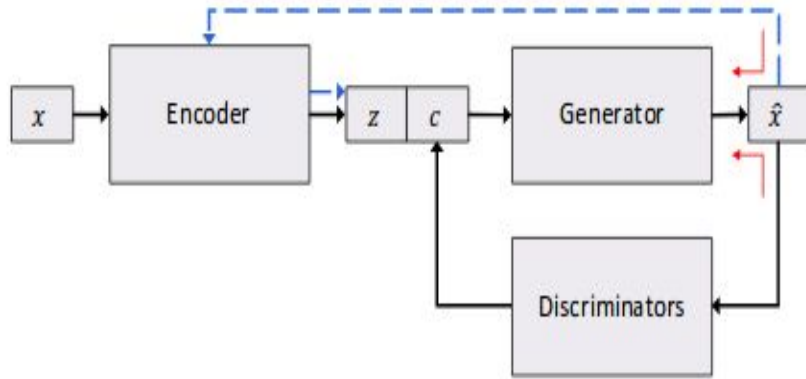
Mathematically:-

$$L_{\text{VAE}}(\theta_G, \theta_E; x) = \text{KL}(q_E(z|x) \| p(z)) - \mathbb{E}_{q_E(z|x)q_D(c|x)}[\log p_G(x|z, c)],$$

Loss function

Forcing latent code towards Gaussian prior

# Model



## Generator Training(Part 2)

- Discriminator produces extra learning signals enforcing the generator to produce text with certain attributes conditioned on code.
- It is still possible that other attributes not explicitly modelled may also entangle with the code.

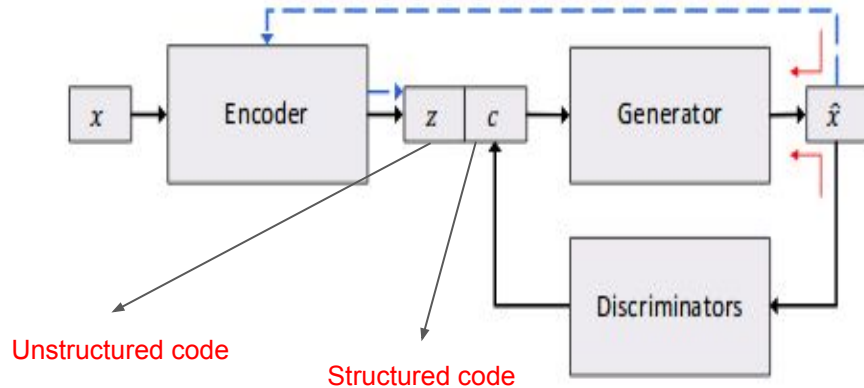
Mathematically:-

$$L_{\text{Attr},c}(\theta_G) = -\mathbb{E}_{p(z)p(c)}[\log_{q_D}(c | G_T(z,c))]$$

Loss function

Train the generator  
based on discriminator  
O/P of generated text

# Model



## Generator Training(Part 3)

- We reuse the Encoder as a discriminator, and match it with the unstructured code(z)

Mathematically:-

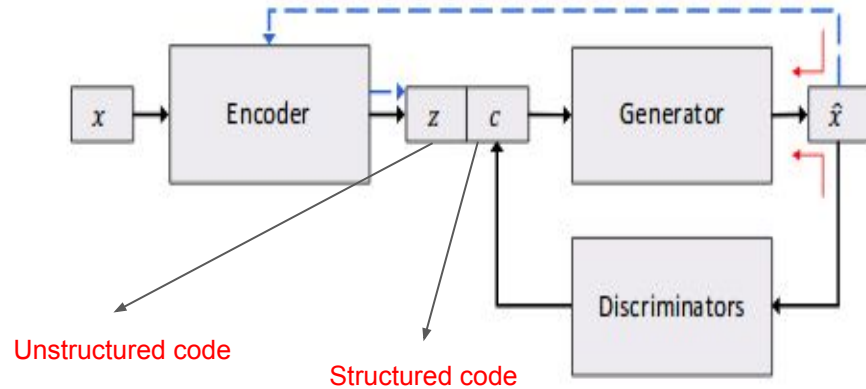
$$L_{Attr,z}(\theta_G) = -\mathbb{E}_{p(z)p(c)} [\log_{qE}(\mathbf{z} | G_T(z,c))]$$

Loss function

Train the generator  
based on discriminator  
O/P of generated text

O/P is z instead of c

# Model



Overall Loss function:-

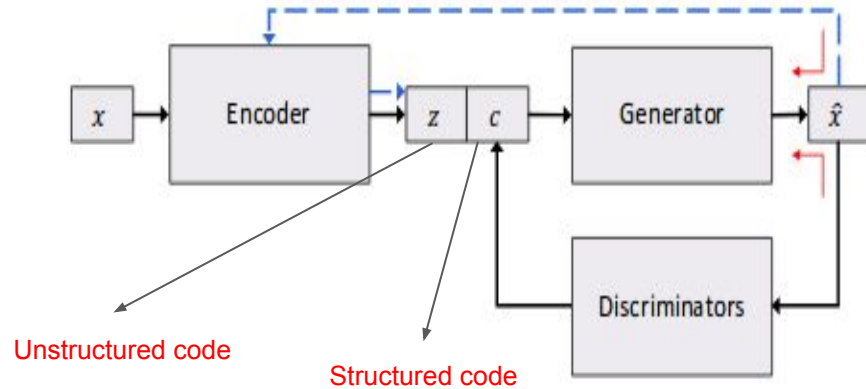
$$\min_{\theta_G} L_G = L_{VAE} + \lambda_c L_{Attr,c} + \lambda_z L_{Attr,z}$$

Loss function

## Generator Training(Overall Loss)

- Overall Loss is the sum of all the losses.
- $\lambda_c$  and  $\lambda_z$  are balancing parameters for the loss

# Model



## Discriminator Training

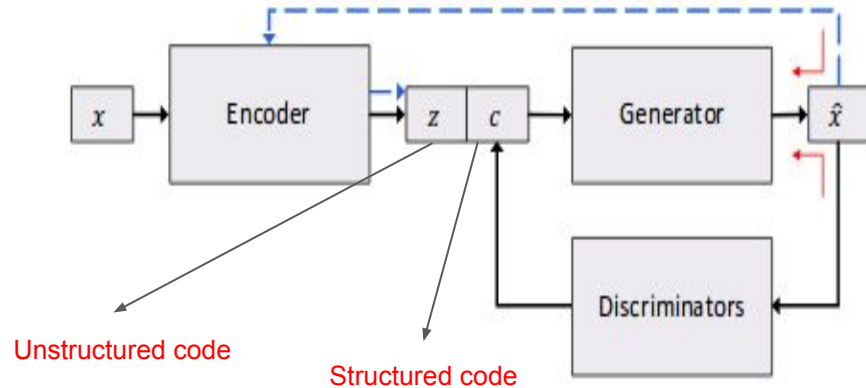
- Can be formulated as a sentence classifier for categorical code, or a probabilistic regressor for continuous code.
- Use labelled examples to train the discriminator, to embed the text characteristics into the code.

Mathematically:-

$$L_s(\theta_D) = - \mathbb{E}_{X_L} [\log_{q_D}(c_L|x_L)]$$

Training on labelled data

# Model



## Discriminator Training

- Use samples generated from generator as a augmented data for training the discriminator
- Use labelled examples to train the discriminator, to embed the text characteristics into the code.

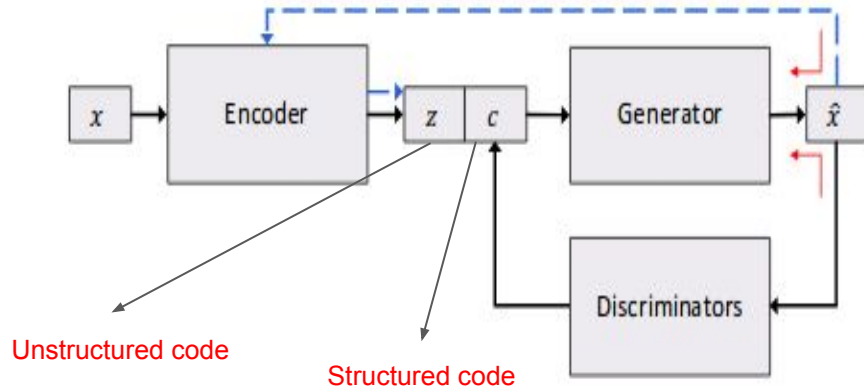
Mathematically:-

$$L_u(\theta_D) = - \mathbb{E}_{p_G(\tilde{x}|z,c)p(z)p(c)} [ \log q_D(c|\tilde{x}) + \beta H(q_D(c|\tilde{x})) ],$$

Training discriminator  
on generator data

Shannon entropy on generated data to  
compensate for noisy data from generator

# Model



Overall loss function :-

$$\min_{\theta_D} L_D = L_s + \lambda_u L_u$$

## Discriminator Training(Overall Loss)

- Use samples generated from generator as a augmented data for training the discriminator
- Use labelled examples to train the discriminator, to embed the text characteristics into the code.

# Training

---

**Algorithm 1** Controlled Generation of Text

---

**Input:** A large corpus of unlabeled sentences  $\mathcal{X} = \{\mathbf{x}\}$

A few sentence attribute labels  $\mathcal{X}_L = \{(\mathbf{x}_L, \mathbf{c}_L)\}$

Parameters:  $\lambda_c, \lambda_z, \lambda_u, \beta$  – balancing parameters

1: Initialize the base VAE by minimizing Eq.(4) on  $\mathcal{X}$  with  $\mathbf{c}$  sampled from prior  $p(\mathbf{c})$

2: **repeat**

3:   Train the discriminator  $D$  by Eq.(11)

4:   Train the generator  $G$  and the encoder  $E$  by Eq.(8) and minimizing Eq.(4), respectively.

5: **until** convergence

**Output:** Sentence generator  $G$  conditioned on disentangled representation  $(\mathbf{z}, \mathbf{c})$

---



# Experiments

Generate short sentences(length <= 15)

**Sentence corpus** - IMDB text corpus(1.4m)

## Sentiment -

- **SST dataset** - 6920/872/1821
- **IMDB review datasets** - 5K/1K/10K
- **SST small** - 250 labelled samples
- **Lexicon** - word level sentiment labels

## Tense -

- Dataset of labelled words
- **5250 words** and phrases labeled with one of {"past", "present", "future"}

# Results (Disentangled Representation)

w/ independency constraint	w/o independency constraint
the film is strictly routine ! the film is full of imagination .	the acting is bad . the movie is so much fun .
after watching this movie , i felt that disappointed . after seeing this film , i 'm a fan .	none of this is very original . highly recommended viewing for its courage , and ideas .
the acting is uniformly bad either . the performances are uniformly good .	too bland highly watchable
this is just awful . this is pure genius .	i can analyze this movie without more than three words . i highly recommend this film to anyone who appreciates music .

# Results (contd..)

---

## Varying the code of tense

---

i thought the movie was too bland and too much  
i guess the movie is too bland and too much  
i guess the film will have been too bland

this was one of the outstanding thrillers of the last decade  
this is one of the outstanding thrillers of the all time  
this will be one of the great thrillers of the all time

---

---

## Varying the unstructured code $z$

---

(*“negative”, “past”*)

the acting was also kind of hit or miss .  
i wish i 'd never seen it  
by the end i was so lost i just did n't care anymore

(*“negative”, “present”*)

the movie is very close to the show in plot and characters  
the era seems impossibly distant  
i think by the end of the film , it has confused itself

(*“negative”, “future”*)

i wo n't watch the movie  
and that would be devastating !  
i wo n't get into the story because there really is n't one

(*“positive”, “past”*)

his acting was impeccable  
this was spectacular , i saw it in theaters twice  
it was a lot of fun

(*“positive”, “present”*)

this is one of the better dance films  
i 've always been a big fan of the smart dialogue .  
i recommend you go see this, especially if you hurt

(*“positive”, “future”*)

i hope he 'll make more movies in the future  
i will definitely be buying this on dvd  
you will be thinking about it afterwards, i promise you

---

# Results (contd...)

---

## Failure cases

---

the plot is not so original

the plot weaves us into <unk>

it does n't get any better the other dance movies

it does n't reach them , but the stories look

he is a horrible actor 's most part

he 's a better actor than a standup

i just think so

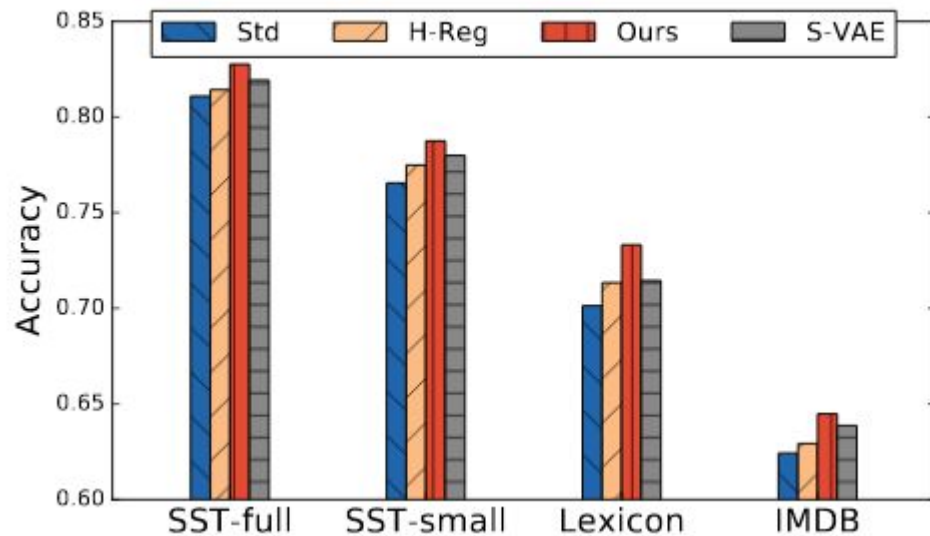
i just think !

---

## Results(Sentiment Accuracy)

Model	Dataset		
	SST-full	SST-small	Lexicon
S-VAE	0.822	0.679	0.660
Ours	<b>0.851</b>	<b>0.707</b>	<b>0.701</b>

## Results(Data Augmentation)



# Questions