

# A Hierarchical Reinforced Sequence Operation Method for Unsupervised Text Style Transfer

---

Chen Wu<sup>1</sup>, Xuancheng Ren<sup>2</sup>, Fuli Luo<sup>2</sup>, Xu Sun<sup>2,3</sup>

Presented at ACL-2019

<sup>1</sup>Tsinghua University

<sup>2</sup>MOE Key Laboratory of Computational Linguistics, School of EECS, Peking University

<sup>3</sup>Center for Data Science, Beijing Institute of Big Data Research, Peking University

# Outline

Background

Our Approach

Main Experiment

Take-Home Messages

# Outline

Background

Our Approach

Main Experiment

Take-Home Messages

# Unsupervised Text Style Transfer

😊 I will be going back and enjoying this great place !

→ 😞 I will not be going back and avoid this horrible place !

## Training data (Unsupervised)

- $\mathcal{X}_1 = \{x_1^{(1)}, \dots, x_1^{(n)}\}$  of style  $s_1$
- $\mathcal{X}_2 = \{x_2^{(1)}, \dots, x_2^{(m)}\}$  of style  $s_2$
- Non-aligned!

## Goal

- $p(x_{1 \rightarrow 2} | x_1)$  that transfers style  $s_1$  into style  $s_2$
- $p(x_{2 \rightarrow 1} | x_2)$  that transfers style  $s_2$  into style  $s_1$

# Prior Work

## Disentanglement Approach

Disentangling latent style and content<sup>1,2,3</sup>

## Two-step Approach

Neutralization (deletion) + stylization (reconstruction)<sup>4,5</sup>

---

<sup>1</sup>Hu et al. “Toward Controlled Generation of Text”. *ICML-17*.

<sup>2</sup>Fu et al. “Style Transfer in Text: Exploration and Evaluation”. *AAAI-18*.

<sup>3</sup>John et al. “Disentangled Representation Learning for Non-Parallel Text Style Transfer”. *ACL-19*.

<sup>4</sup>Li et al. “Delete, Retrieve, Generate: A Simple Approach to Sentiment and Style Transfer”. *NAACL-HLT-18*.

<sup>5</sup>Xu et al. “Unpaired Sentiment-to-Sentiment Translation: A Cycled Reinforcement Learning Approach”. *ACL-18*.

# Challenges

## Poor Content Preservation

---

**Original.**    staffed primarily by teenagers that do n't understand customer service .

---

**System 1.**    staffed , the best and sterile by flies , how fantastic customer service .

**System 2.**    staffed established each tech feel when great customer service professional .

**System 3.**    staffed distance that love customer service .

---

## Lack of Interpretability

# Outline

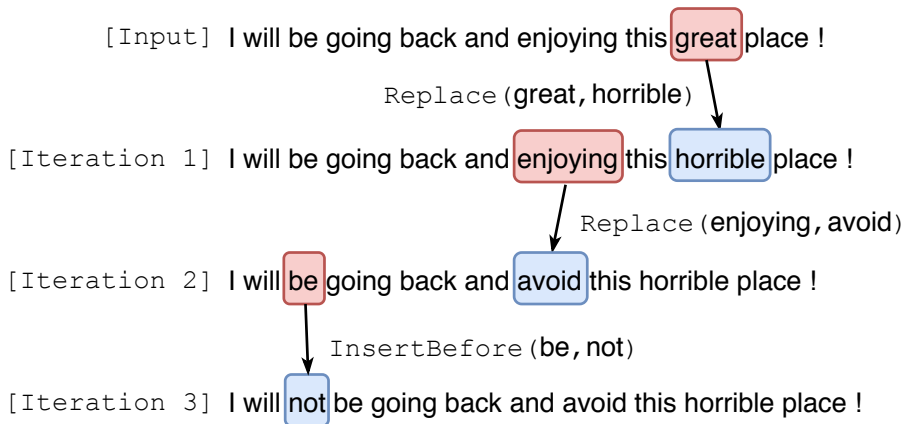
Background

**Our Approach**

Main Experiment

Take-Home Messages

## An Example Case





# A Hierarchy of Agents

## The Options Framework

An HRL framework proposed by  
Sutton et al. (1999)

### High-Level Agent

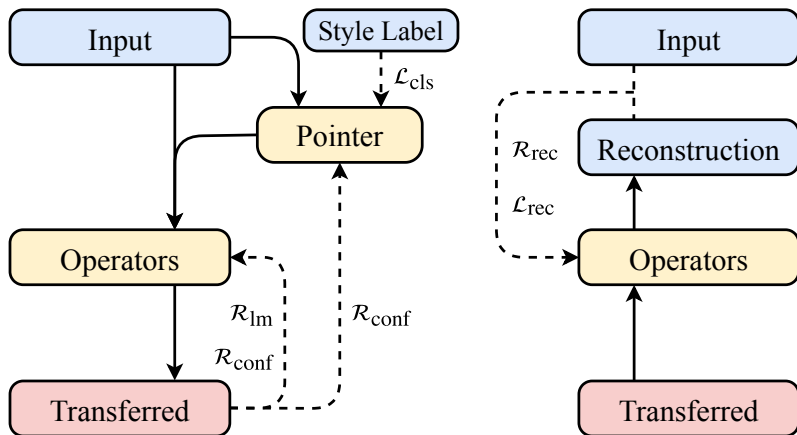
Propose a position to be operated  
around

### Low-Level Agent

Select an operator from the table  
and generate a word  $\hat{w}$  (*optional*)

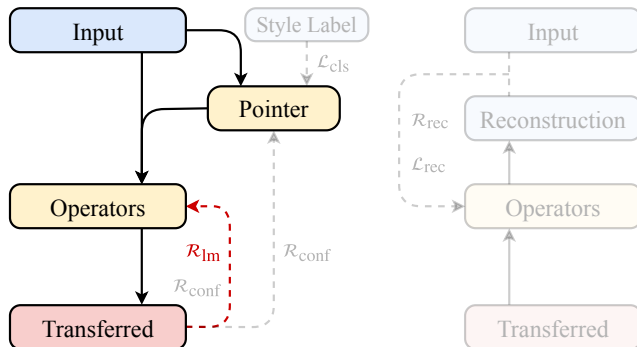
Operator	Operation
IF <sub><math>\phi_1</math></sub>	Insert a word $\hat{w}$ in the <b>F</b> ront
IB <sub><math>\phi_2</math></sub>	Insert a word $\hat{w}$ <b>B</b> ehind
Rep <sub><math>\phi_3</math></sub>	<b>R</b> eplace the word with $\hat{w}$
DC	<b>D</b> elete the <b>C</b> urrent word
DF	<b>D</b> elete the word in the <b>F</b> ront
DB	<b>D</b> elete the word <b>B</b> ehind
Skip	Do not change anything

## Graphical Overview for Training



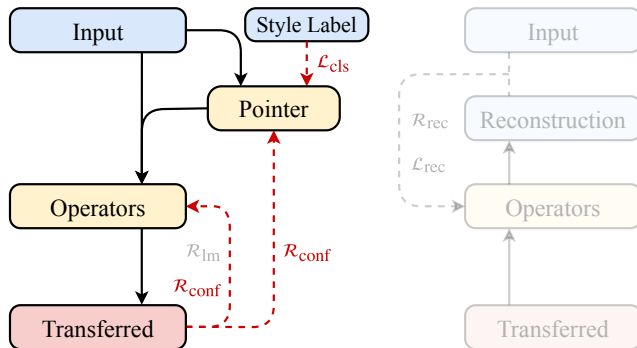
# Hierarchical Policy Learning

- Language model reward



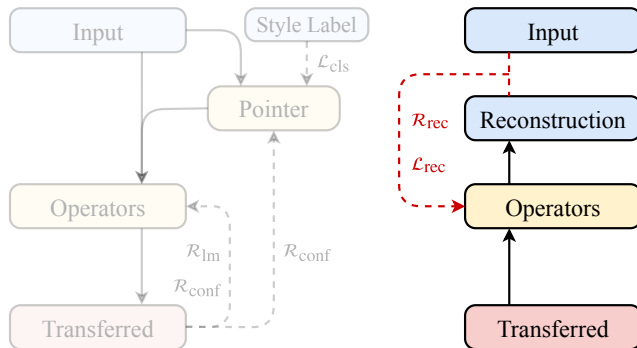
# Hierarchical Policy Learning

- Classification confidence reward
- Auxiliary task: style classification



# Hierarchical Policy Learning

- Self-supervised reconstruction loss
- Reconstruction reward



# Training with Single-Option Trajectory

## Motivations for Single-Step Training

- High variance of policy gradients
- Ambiguity in self-supervised reconstruction

# Iterative and Dynamic Inference

## Basic Ideas

- Enumerate the operators and select the *best* one in each iteration
- Greedy modifications
- **Mask**-mechanisms
- Until the sentence does not show the original style anymore  
(or beyond maximum iterations)

# Outline

Background

Our Approach

**Main Experiment**

Take-Home Messages



# Automatic Evaluation

	Yelp		Amazon	
	Acc	BLEU	Acc	BLEU
CrossAligned	74.7	9.06	75.1	1.90
MultiDecoder	50.6	14.54	69.9	9.07
StyleEmbedding	8.4	21.06	38.2	15.07
TemplateBased	81.2	22.57	64.3	34.79
DeleteOnly	86.0	14.64	47.0	33.00
Del-Ret-Gen	88.6	15.96	51.0	30.09
BackTranslate	94.6	2.46	<b>76.7</b>	1.04
UnpairedRL	57.5	18.81	56.3	15.93
UnsuperMT	<b>97.8</b>	22.75	72.4	33.95
Human	74.7	-	43.2	-
Point-Then-Operate	91.5	<b>29.86</b>	40.2	<b>41.86</b>

- Classification accuracy is **low** for human references
- BLEU of our method outperforms baselines by a large margin

## Human Evaluation<sup>6</sup>

	Yelp				Amazon			
	Flu.	Cont.	Sty.	Suc	Flu.	Cont.	Sty.	Suc
TemplateBased	3.47	3.76	3.25	68.0 %	3.46	4.08	2.15	9.0 %
Del-Ret-Gen	3.82	3.73	3.52	70.3 %	4.02	4.31	2.69	21.0 %
UnpairedRL	3.54	3.59	2.90	53.8 %	2.58	2.55	2.44	4.5 %
UnsuperMT	4.26	4.24	<b>4.03</b>	<b>82.5 %</b>	4.24	4.13	3.05	35.5 %
Point-Then-Operate	<b>4.39</b>	<b>4.56</b>	3.78	81.5 %	<b>4.28</b>	<b>4.47</b>	<b>3.31</b>	<b>47.0 %</b>

- Better overall performance
- Sacrificed style polarity on Yelp

---

<sup>6</sup>Baselines for human evaluation are selected based on automatic evaluation

# Outline

Background

Our Approach

Main Experiment

Take-Home Messages

## Take-Home Messages

1. A **sequence operation** method with hierarchical reinforcement learning (**HRL**) for unsupervised text style transfer
2. Address two challenges
  - **Content preservation**
  - **Interpretability**
3. Provide an iterative and dynamic **mask**-based inference algorithm that allows for single-option trajectory training



We make our code public.