

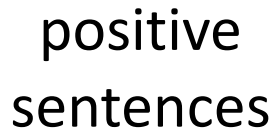
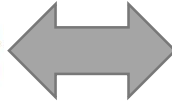
Text *Style* Transfer

Hung-yi Lee 李宏毅

Audio Style



female



negative
sentences

Text Style Transfer

Text Style Transfer

你真笨
(negative)



Seq2seq

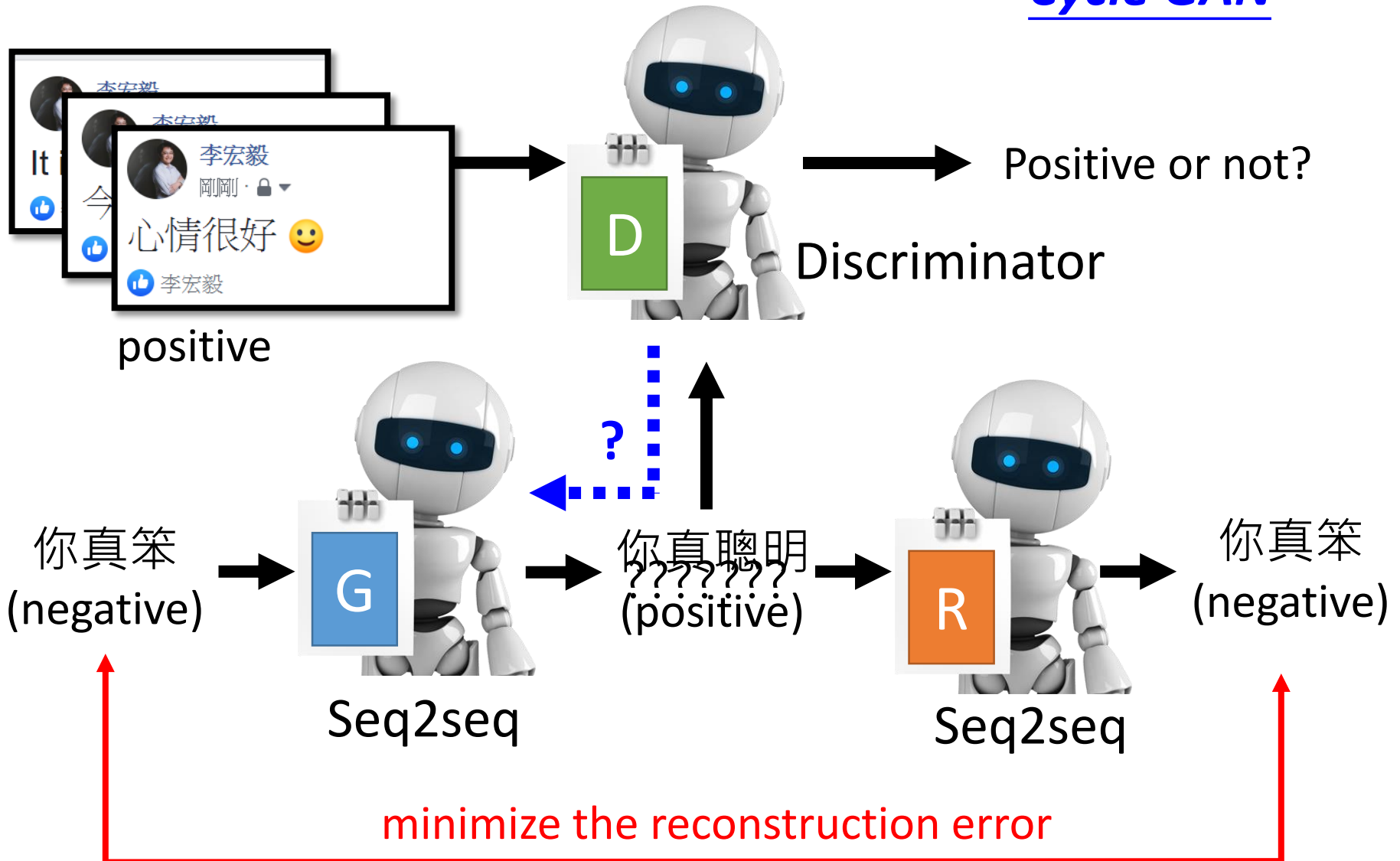


你真聰明
??????
(positive)



Text Style Transfer

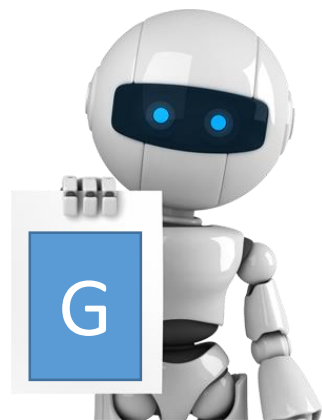
Cycle GAN



Can we use
gradient ascent?


NO!

scalar 



Seq2seq

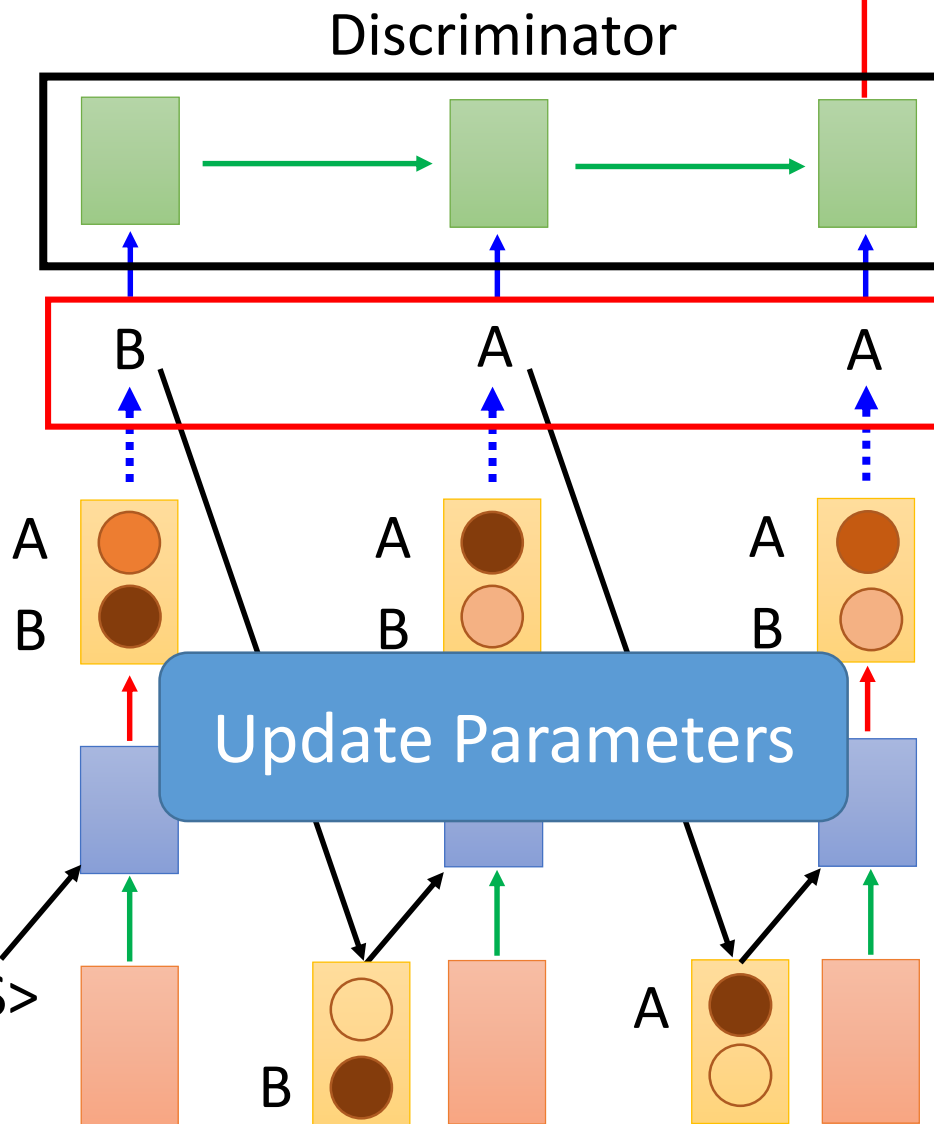
Generator

 : obtained
by attention

<BOS>

B

A



Three Categories of Solutions

Gumbel-softmax

- [Matt J. Kusner, et al., arXiv, 2016][Weili Nie, et al. ICLR, 2019]

Continuous Input for Discriminator

- [Sai Rajeswar, et al., arXiv, 2017][Ofir Press, et al., ICML workshop, 2017][Zhen Xu, et al., EMNLP, 2017][Alex Lamb, et al., NIPS, 2016][Yizhe Zhang, et al., ICML, 2017]

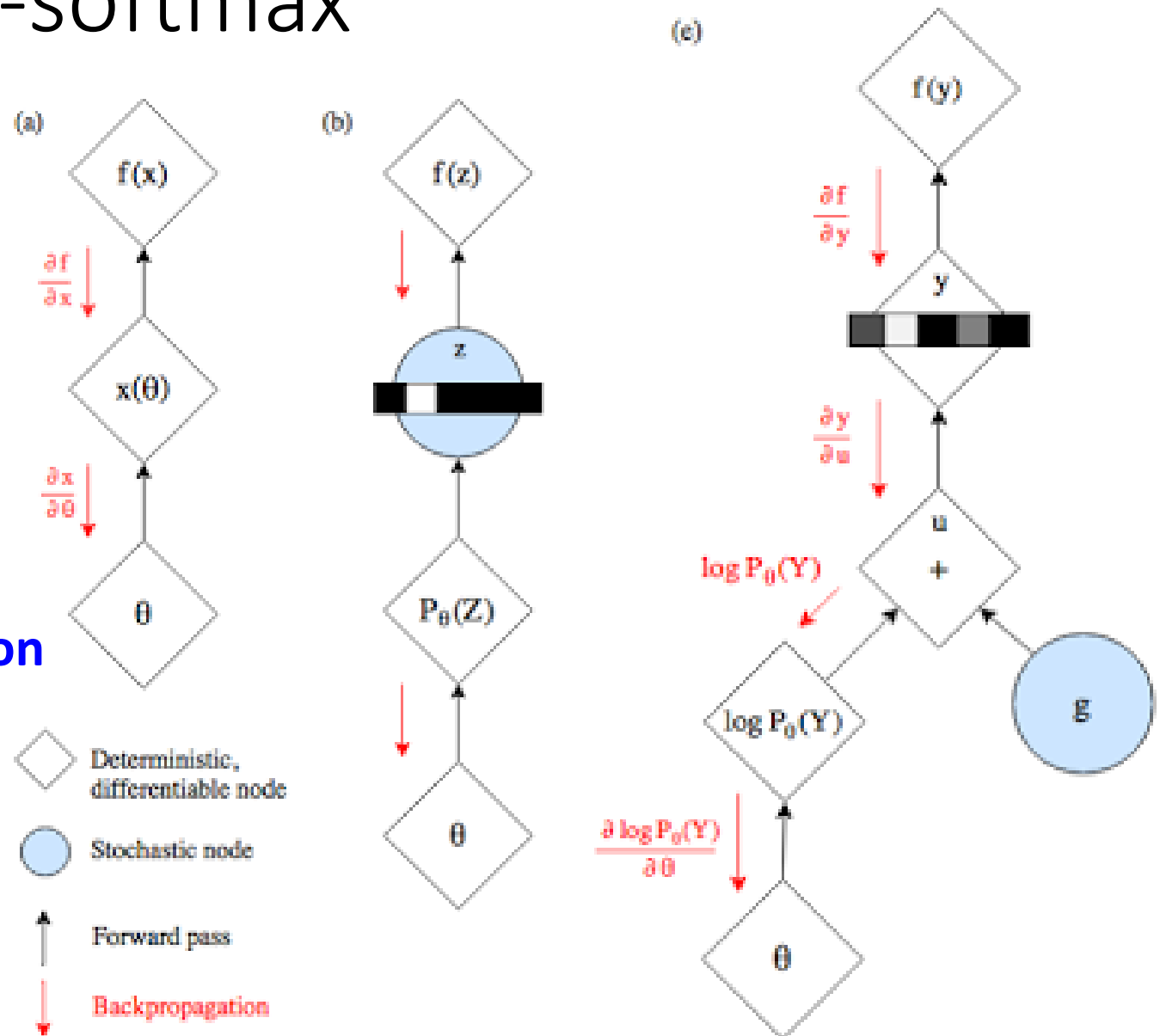
Reinforcement Learning

- [Yu, et al., AAAI, 2017][Li, et al., EMNLP, 2017][Tong Che, et al, arXiv, 2017][Jiaxian Guo, et al., AAAI, 2018][Kevin Lin, et al, NIPS, 2017][William Fedus, et al., ICLR, 2018]

Gumbel-softmax

Using the
reparameterization
trick

As what people
do for training
VAE



Three Categories of Solutions

Gumbel-softmax

- [Matt J. Kusner, et al., arXiv, 2016][Weili Nie, et al. ICLR, 2019]

Continuous Input for Discriminator

- [Sai Rajeswar, et al., arXiv, 2017][Ofir Press, et al., ICML workshop, 2017][Zhen Xu, et al., EMNLP, 2017][Alex Lamb, et al., NIPS, 2016][Yizhe Zhang, et al., ICML, 2017]

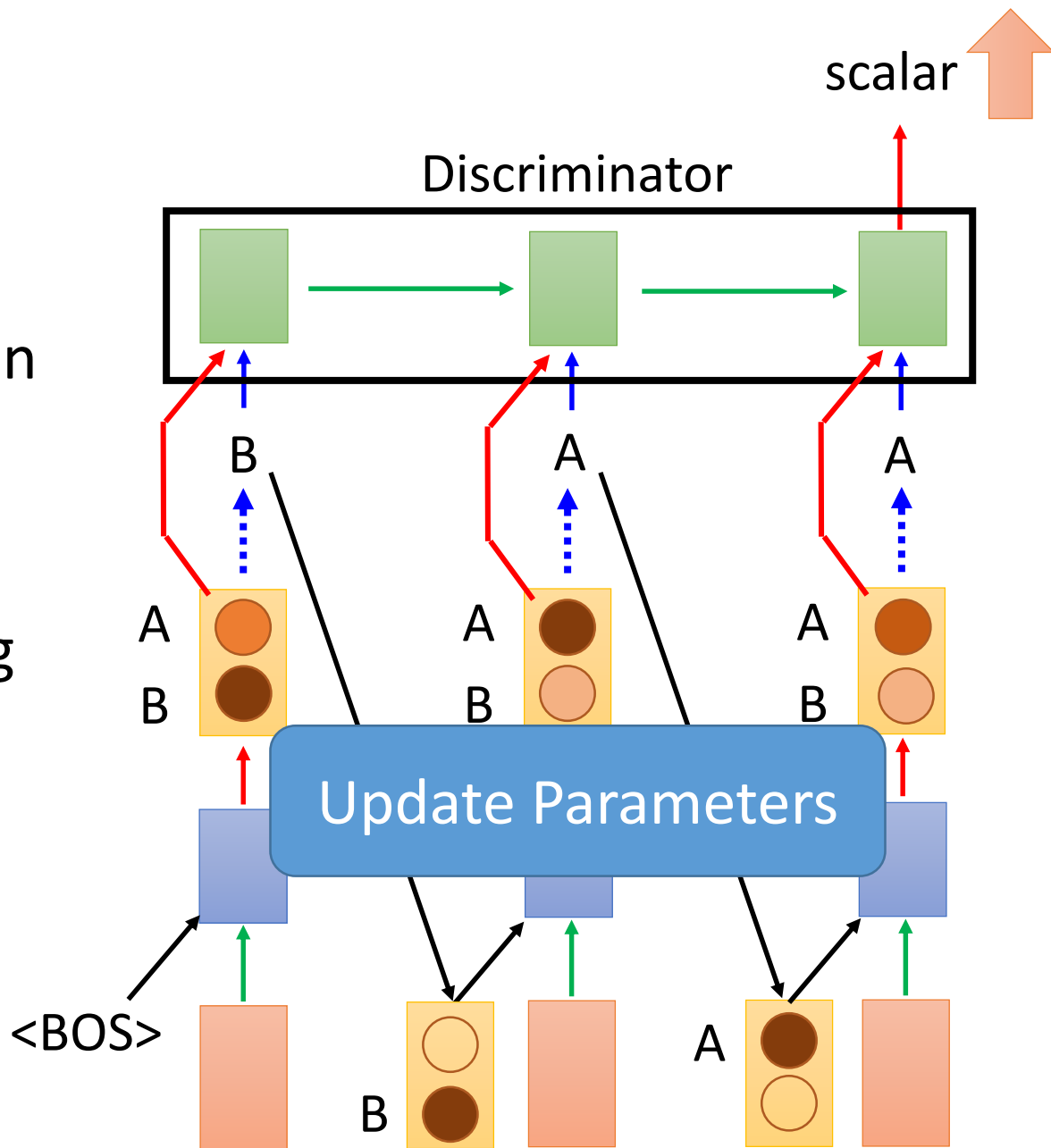
Reinforcement Learning

- [Yu, et al., AAAI, 2017][Li, et al., EMNLP, 2017][Tong Che, et al, arXiv, 2017][Jiaxian Guo, et al., AAAI, 2018][Kevin Lin, et al, NIPS, 2017][William Fedus, et al., ICLR, 2018]

Use the distribution
as the input of
discriminator

Avoid the sampling
process

We can do
backpropagation
now.



What is the problem?

Discriminator with constraint
(e.g. WGAN) can be helpful.

- Real sentence

1	0	0	0	0
0	1	0	0	0
0	0	1	0	0
0	0	0	1	0
0	0	0	0	1

Discriminator can immediately find the difference.

- Generated

0.9	0.1	0.1	0	0
0.1	0.9	0.1	0	0
0	0	0.7	0.1	0
0	0	0.1	0.8	0.1
0	0	0	0.1	0.9

Can never
be 1-hot

Three Categories of Solutions

Gumbel-softmax

- [Matt J. Kusner, et al., arXiv, 2016][Weili Nie, et al. ICLR, 2019]

Continuous Input for Discriminator

- [Sai Rajeswar, et al., arXiv, 2017][Ofir Press, et al., ICML workshop, 2017][Zhen Xu, et al., EMNLP, 2017][Alex Lamb, et al., NIPS, 2016][Yizhe Zhang, et al., ICML, 2017]

Reinforcement Learning

- [Yu, et al., AAAI, 2017][Li, et al., EMNLP, 2017][Tong Che, et al, arXiv, 2017][Jiaxian Guo, et al., AAAI, 2018][Kevin Lin, et al, NIPS, 2017][William Fedus, et al., ICLR, 2018]

The reward function
may change

→ Different from typical RL

Reward ← scalar ↑

Discriminator

Environment

Actions taken

Generator
= Agent in RL

Trained by RL algorithm
(e.g. Policy Gradient)

<BOS>

B

A

A

B

A

B

A

B

B

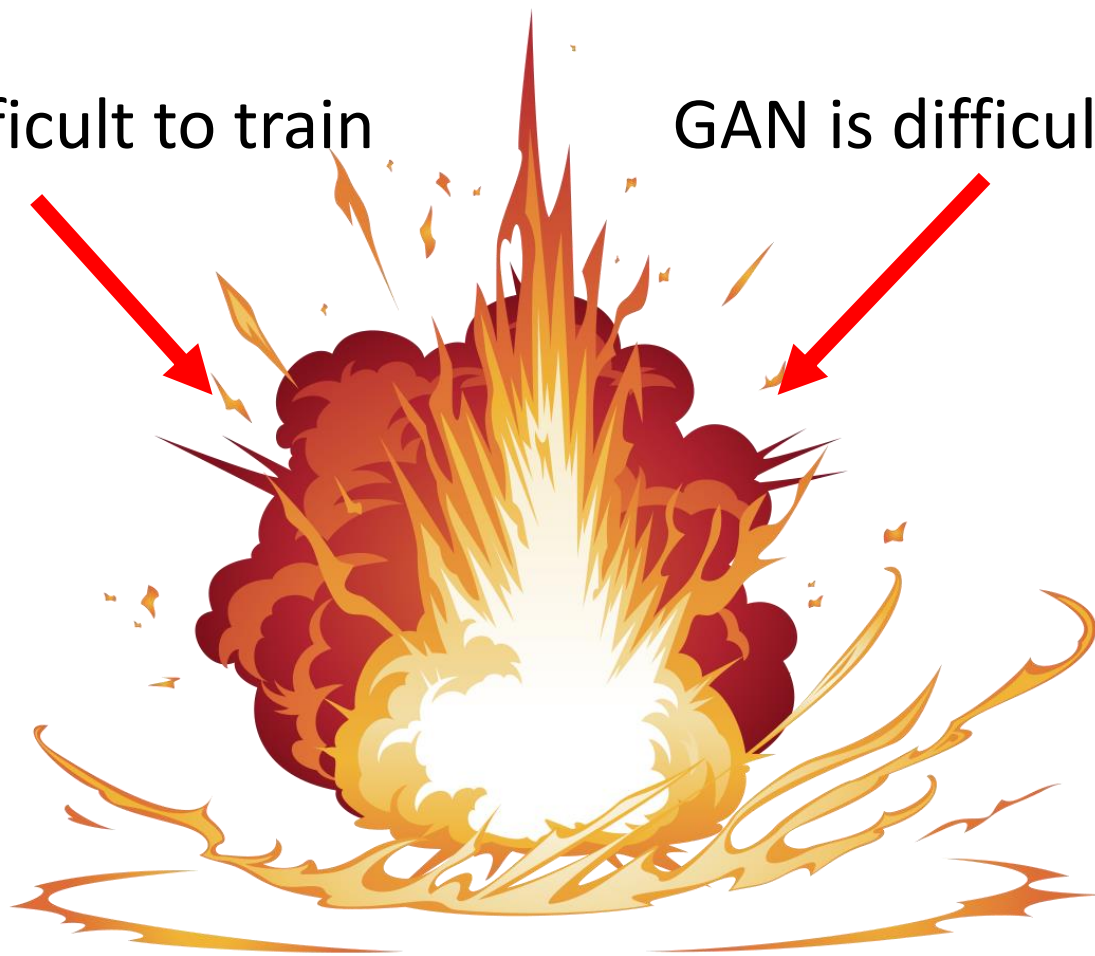
A

A

Disaster

· RL is difficult to train

GAN is difficult to train

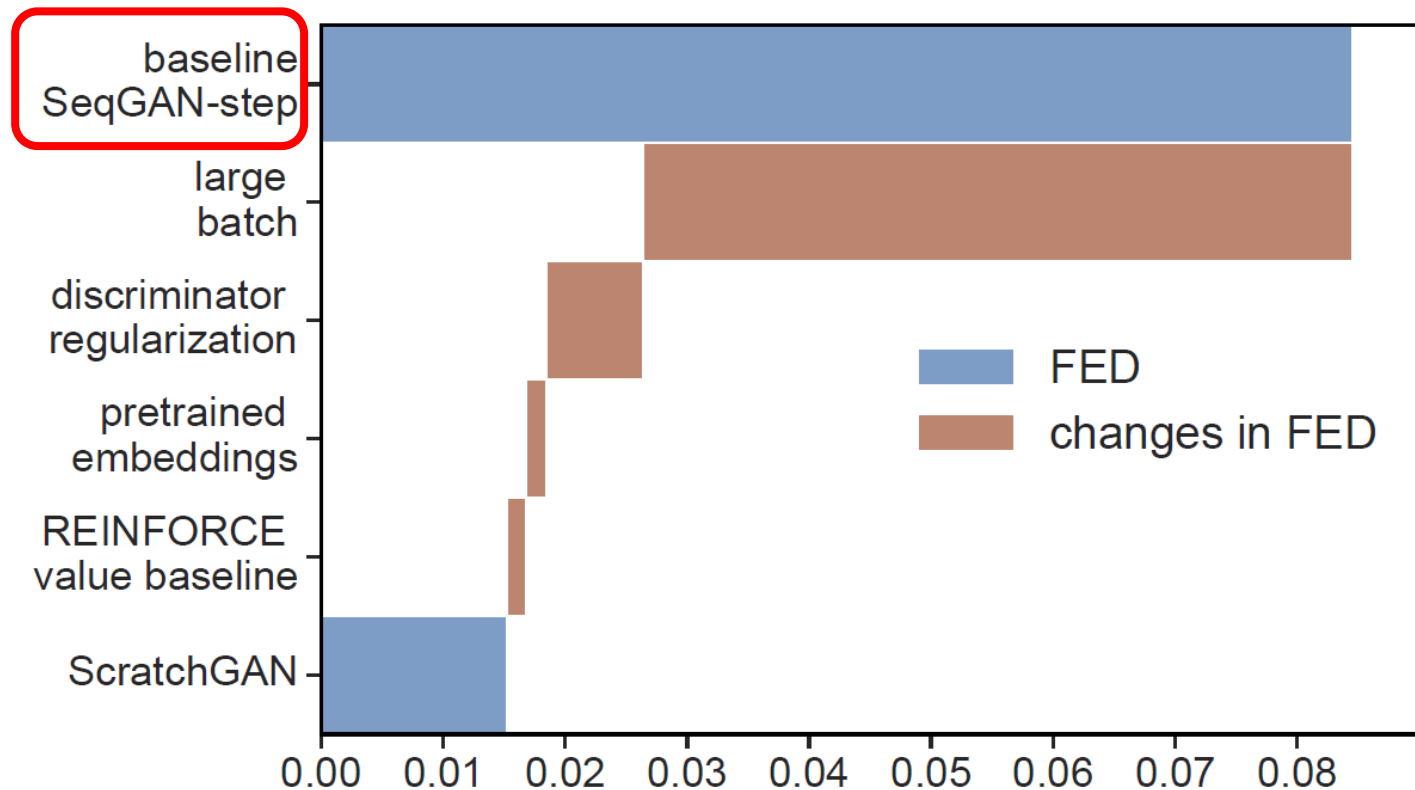


RL+GAN

Tips?

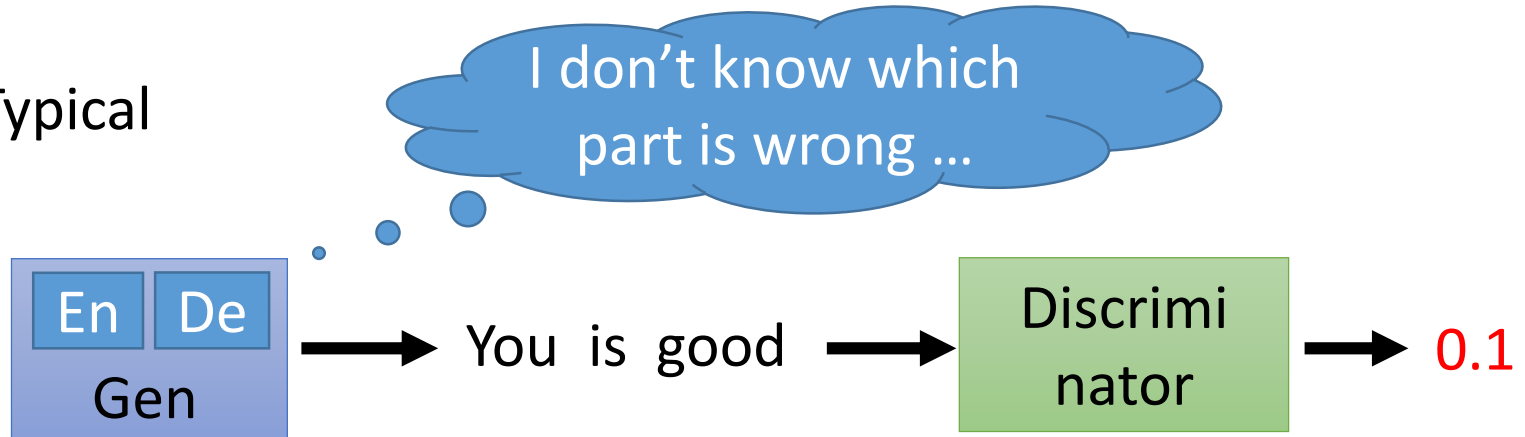
[Cyprien de Masson d'Autume, et al., arXiv 2019]

- ScarchGAN

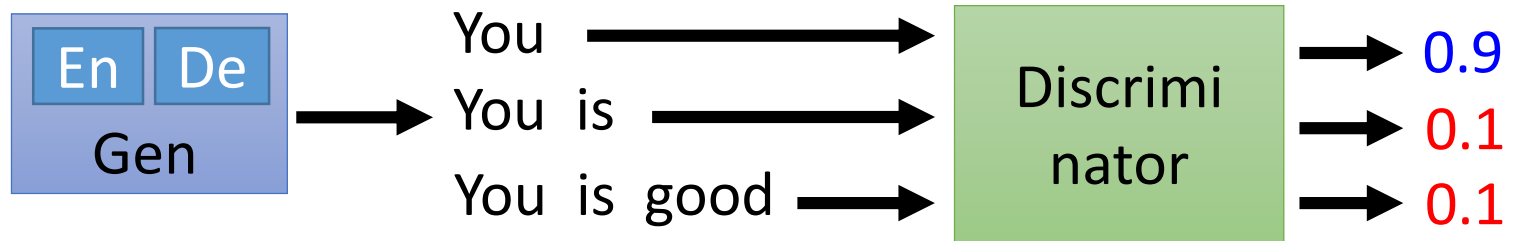


Tips?

- Typical

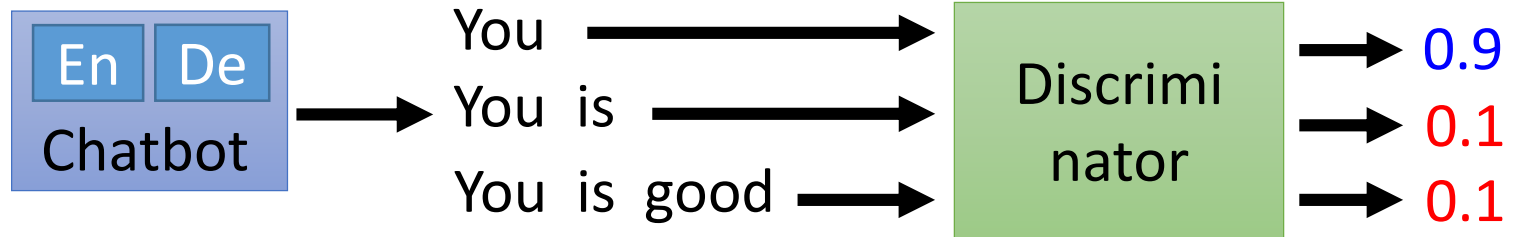


- Reward for Every Generation Step



Tips?

- Reward for Every Generation Step



Method 1. Monte Carlo (MC) Search [\[Yu, et al., AAAI, 2017\]](#)

Method 2. Discriminator For Partially Decoded Sequences

[\[Li, et al., EMNLP, 2017\]](#)

Method 3. Step-wise evaluation [\[Tual, Lee, TASLP, 2019\]](#)[\[Xu, et al., EMNLP, 2018\]](#)[\[William Fedus, et al., ICLR, 2018\]](#)

Text Style Transfer

[Lee, et al., ICASSP'18]



- From **negative** sentence to **positive** one

胃疼, 沒睡醒, 各種不舒服

我都想去上班了, 真夠賤的!

暈死了, 吃燒烤、竟然遇到個變態狂

我肚子痛的厲害

Relaxed ↔ Annoyed

Relaxed	Sitting by the Christmas tree and watching Star Wars after cooking dinner. What a nice night 🍷🌲💎
Annoyed	Sitting by the computer and watching The Voice for the second time tonight. What a horrible way to start the weekend 😡😡😡
Annoyed	Getting a speeding ticket 50 feet in front of work is not how I wanted to start this month 🙄
Relaxed	Getting a haircut followed by a cold foot massage in the morning is how I wanted to start this month 😊

Male ↔ Female

Male	Gotta say that beard makes you look like a Viking...
Female	Gotta say that hair makes you look like a Mermaid...
Female	Awww he's so gorgeous 😍 can't wait for a cuddle. Well done 🙄 xxx
Male	Bro he's so f***ing dope can't wait for a cuddle. Well done bro

Age 18-24 ↔ 65+

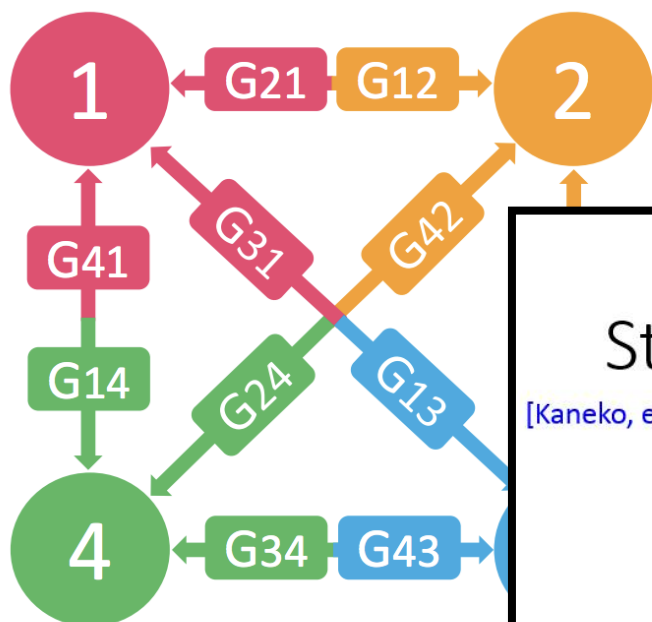
18-24	You cheated on me but now I know nothing about loyalty 😂 ok
65+	You cheated on America but now I know nothing about patriotism. So ok.
65+	Ah! Sweet photo of the sisters. So happy to see them together today .
18-24	Ah 😂 Thankyou ❤️ #sisters ❤️ happy to see them together today



¹Note that using “gender” (or any other attribute for that matter) as a differentiating attribute between several bodies of text implies that there are indeed signatures of gender in the data. These signatures could be as innocuous as some first names like Mary being usually associated with women, or disheartening like biases and stereotypes exposed by statistical methods, (e.g., “man is to computer programmer as woman is to home-maker” (Bolukbasi et al., 2016)). We certainly do not condone those stereotypes, and on the contrary, we hope that showing that our models can uncover these biases might down the line turn them into powerful tools for researchers who study fairness and debiasing (Reddy & Knight, 2016).

Source of image: <https://openreview.net/forum?id=H1g2NhC5KQ>

(a) Cross-domain models



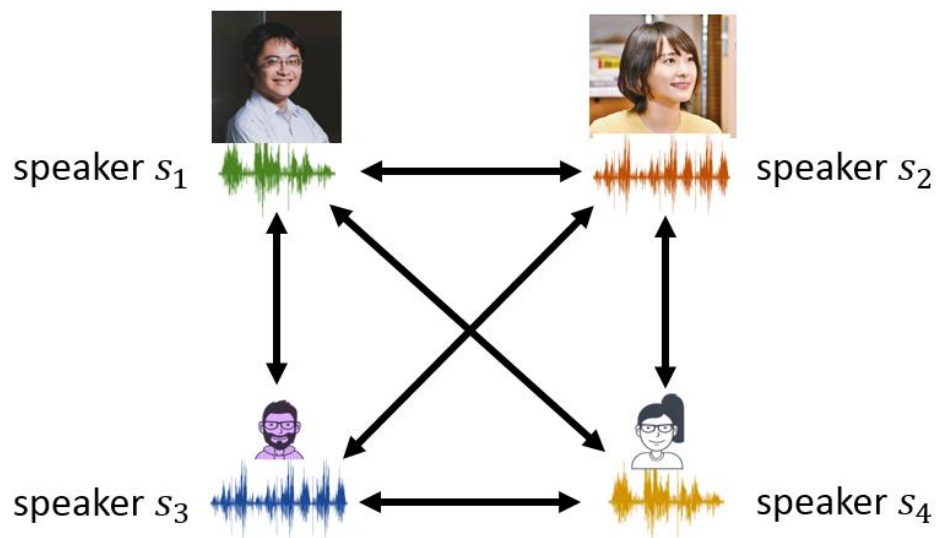
(b) StarGAN

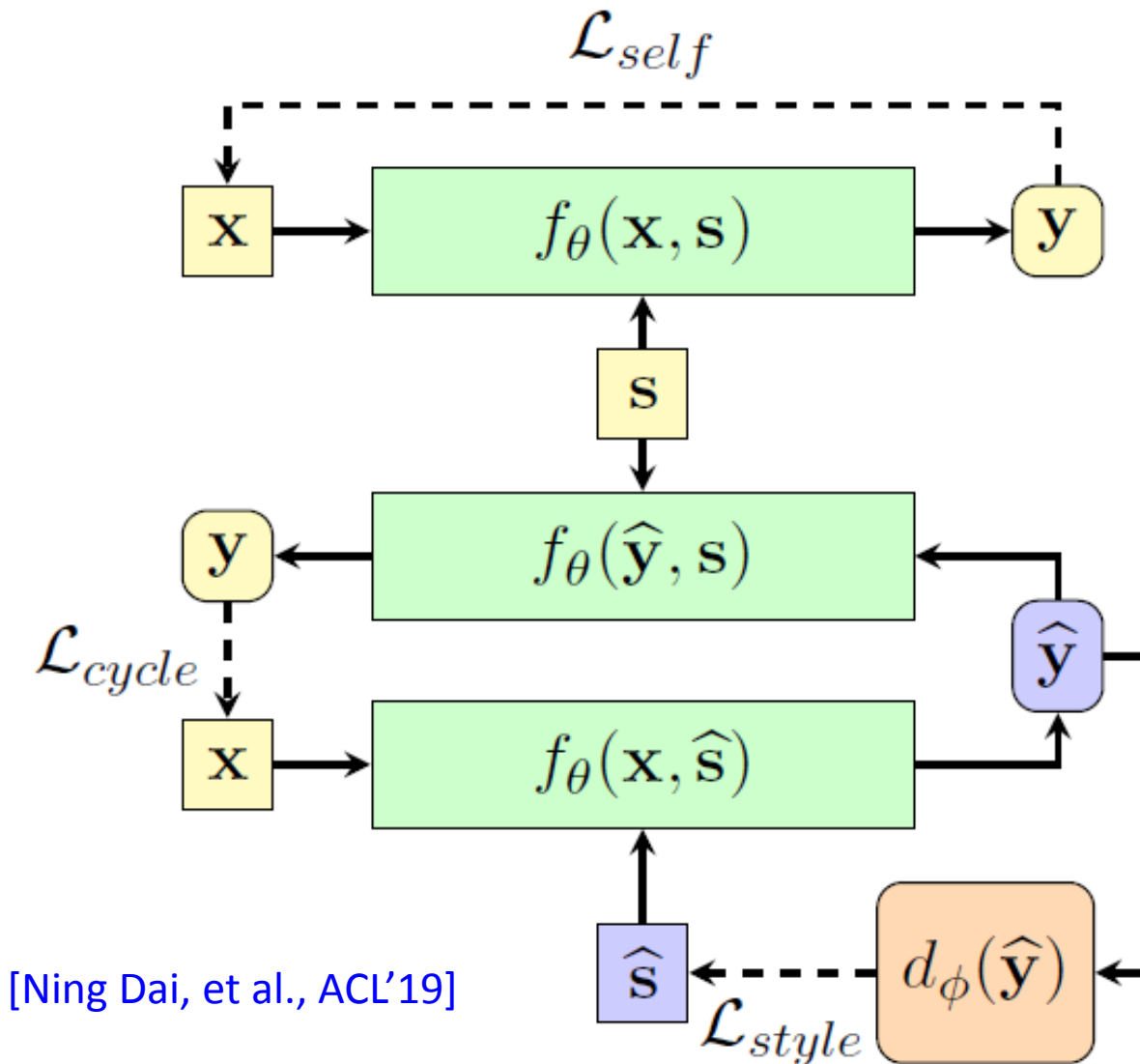
StarGAN

[Kaneko, et al., INTERSPEECH'19]

For CycleGAN:

If there are N speakers, you need $N \times (N-1)$ generators.

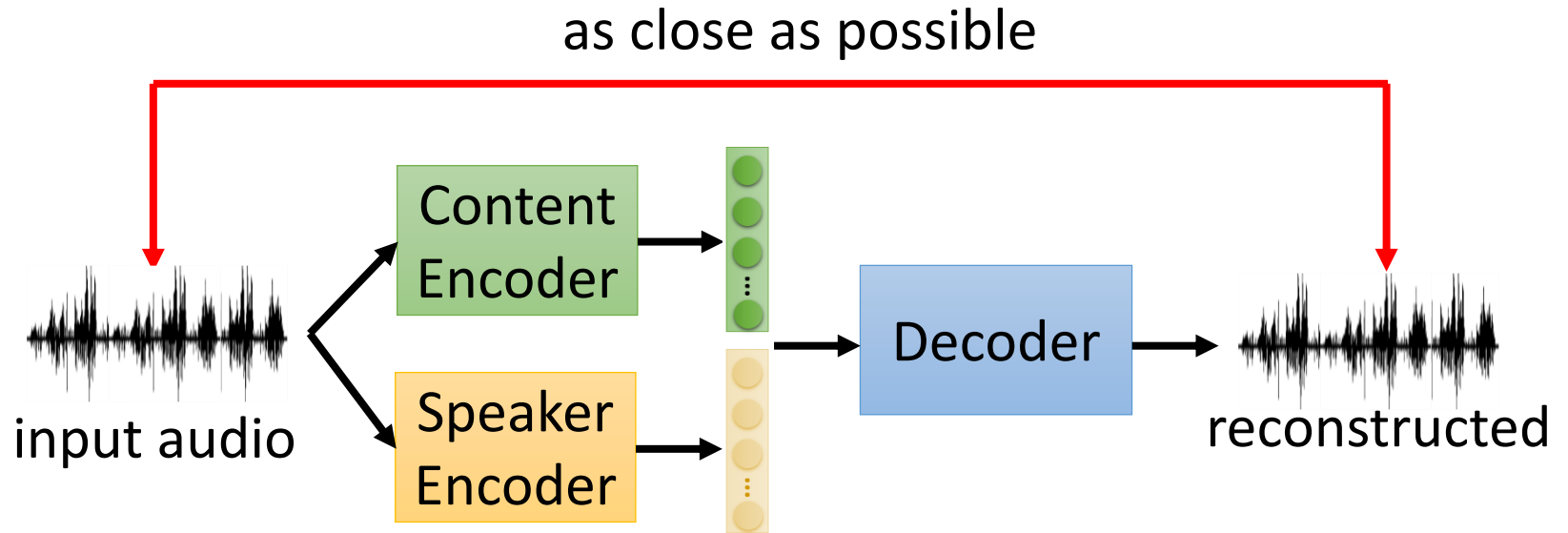




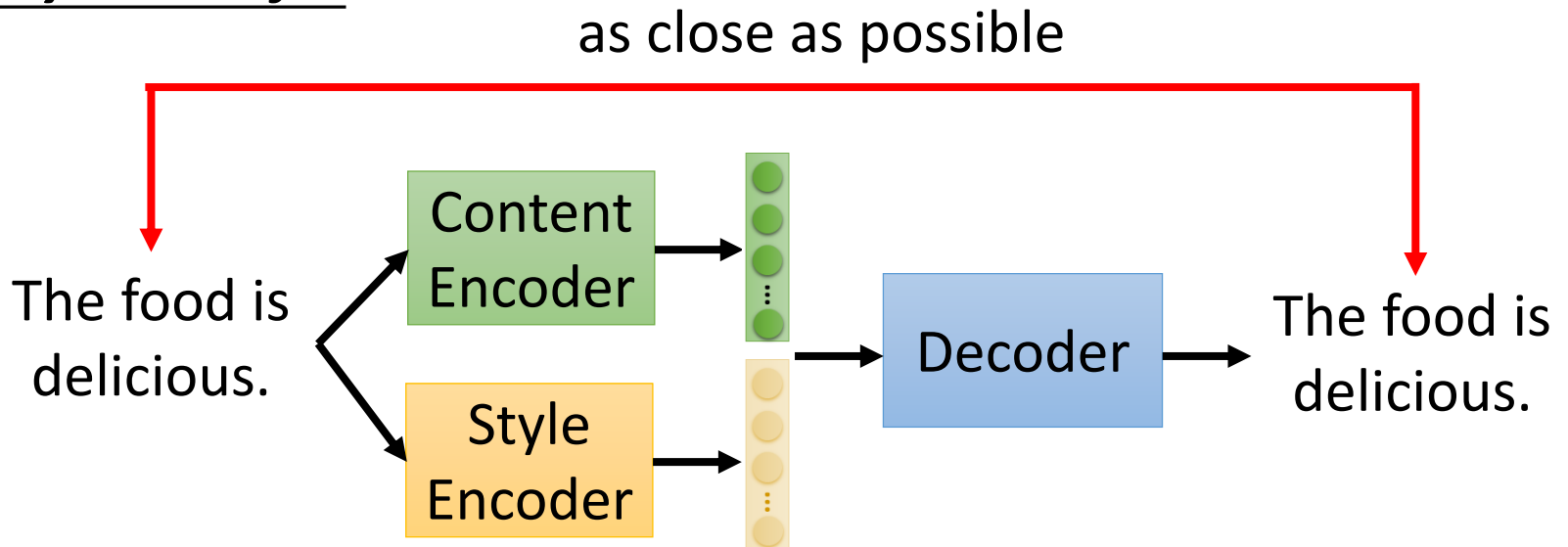
Style Transformer (Text version of StarGAN)

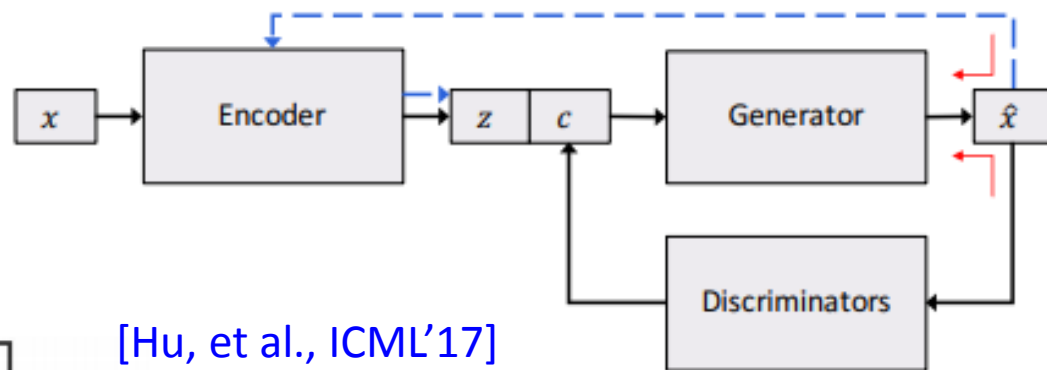
Source of image: <https://arxiv.org/abs/1905.05621>

Voice Conversion

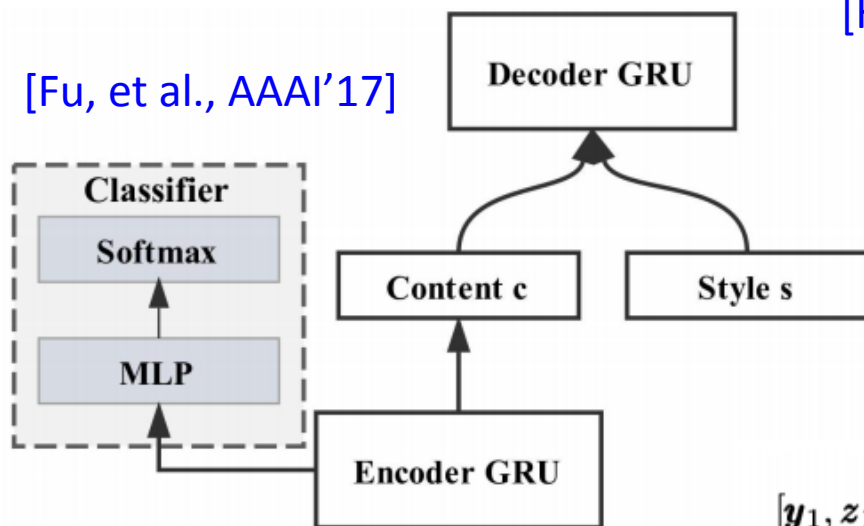


Text Style Transfer

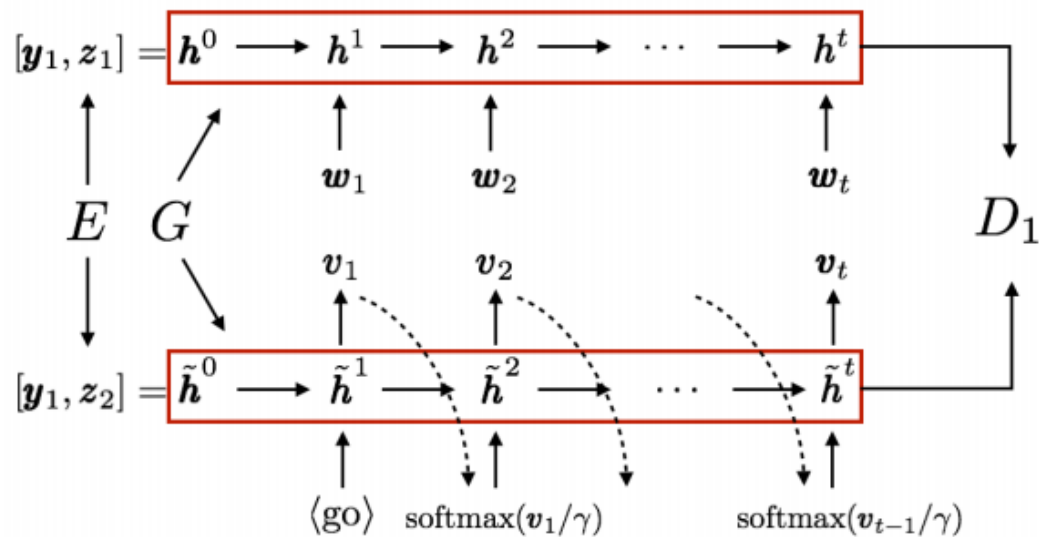




[Hu, et al., ICML'17]

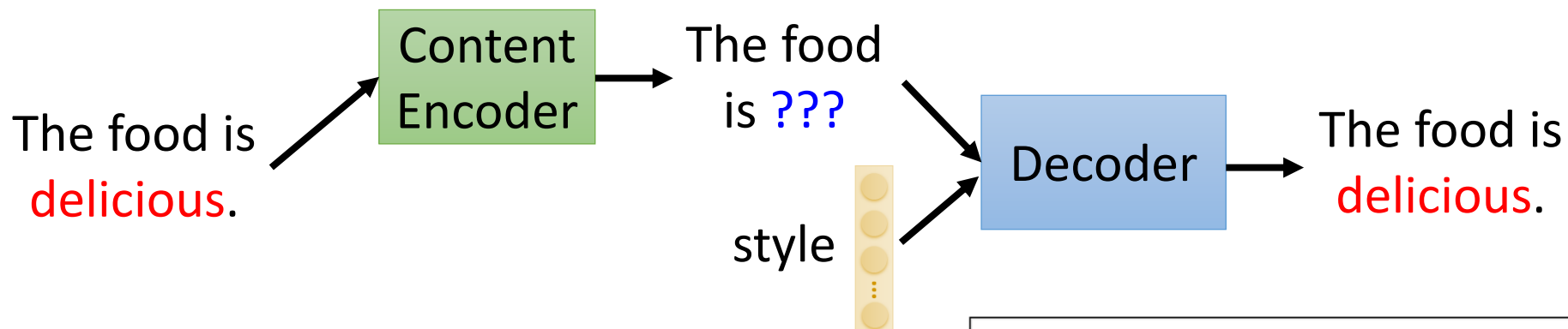


[Fu, et al., AAAI'17]



[Shen, et al., NIPS'17]

Text Style Transfer



great food but horrible staff and very very rude workers !

Delete attribute markers

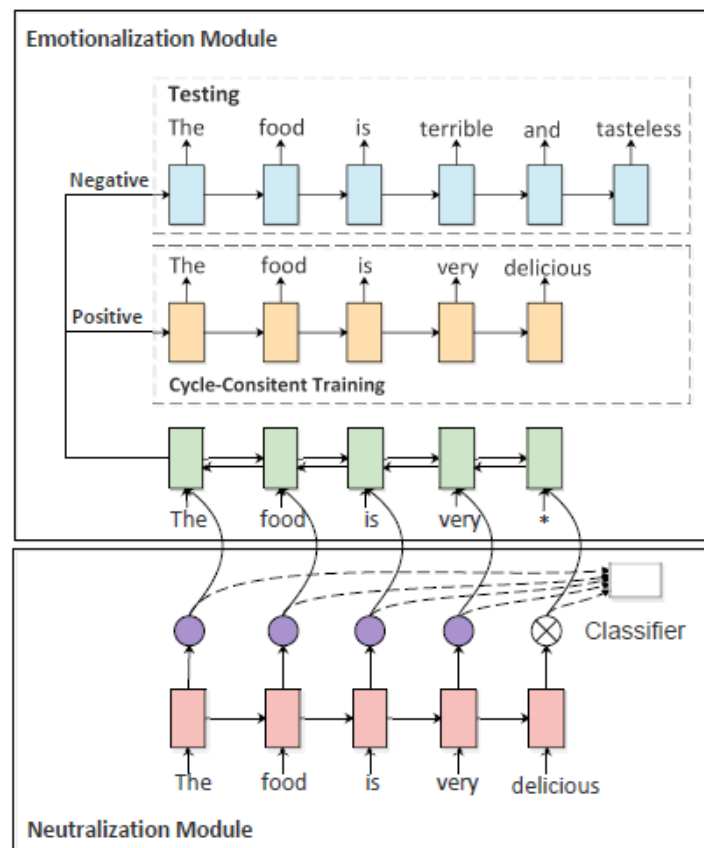
great food staff and very workers ! target=positive

Run system

*great food , awesome staff , very personable
and very efficient atmosphere !*

[Li, et al., NAACL'18]

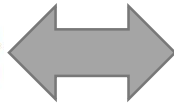
[Xu, et al., ACL'18]



Audio Style



female



negative
sentences

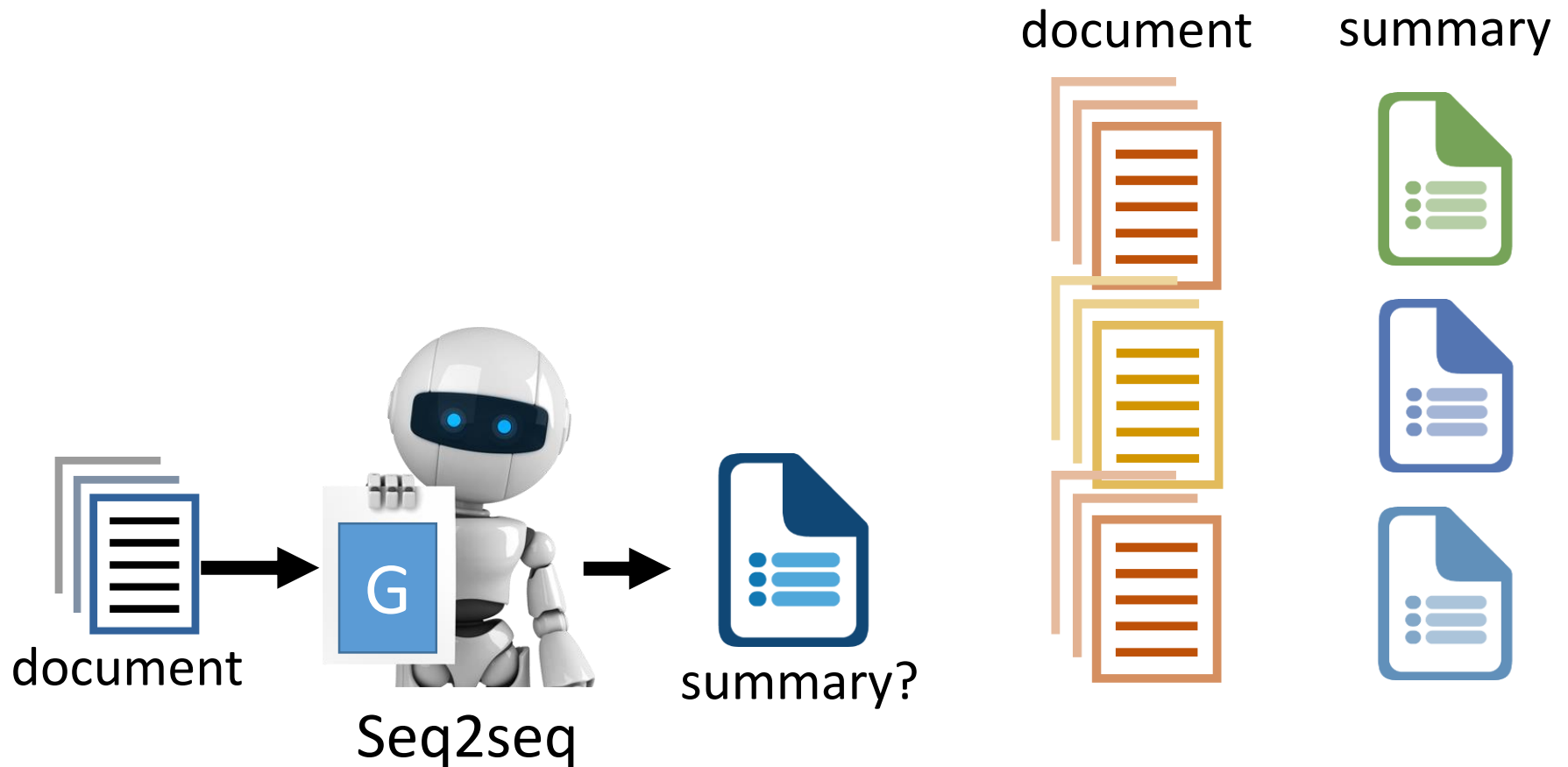
Text Style Transfer



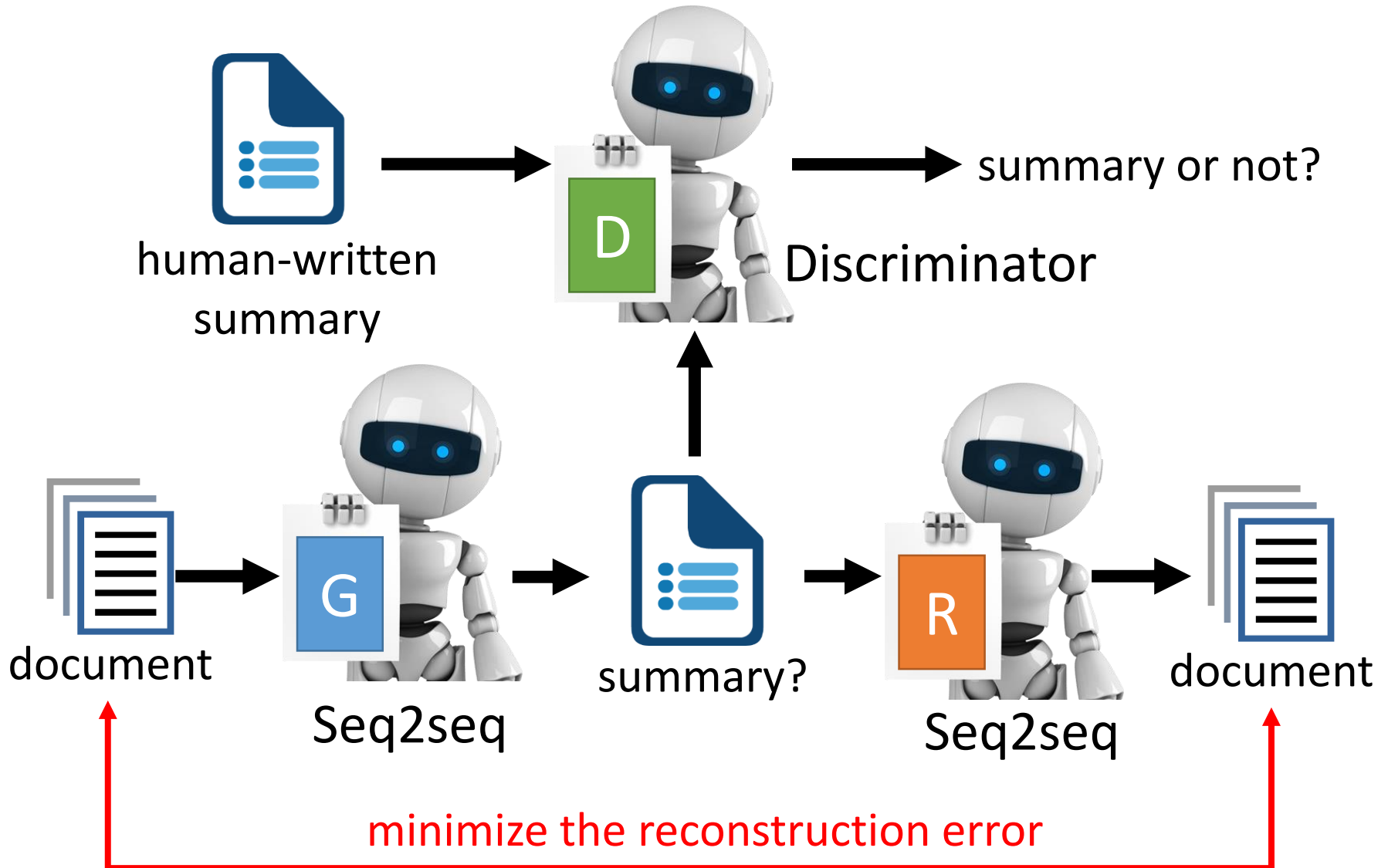
summary

Unsupervised Abstractive Summarization

Unsupervised Abstractive Summarization



Unsupervised Abstractive Summarization



Summarization

[Wang, Lee,
EMNLP 2018]

English Gigaword (Document title as summary)

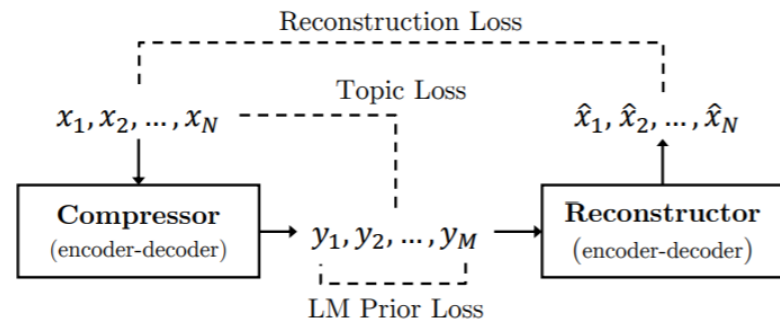
	ROUGE-1	ROUGE-2	ROUGE-L
Supervised	33.2	14.2	30.5
Trivial	21.9	7.7	20.5
Unsupervised (matched data)	28.1	10.0	25.4
Unsupervised (no matched data)	27.2	9.1	24.1

- Matched data: using the title of English Gigaword to train Discriminator
- No matched data: using the title of CNN/Diary Mail to train Discriminator

More Unsupervised Summarization

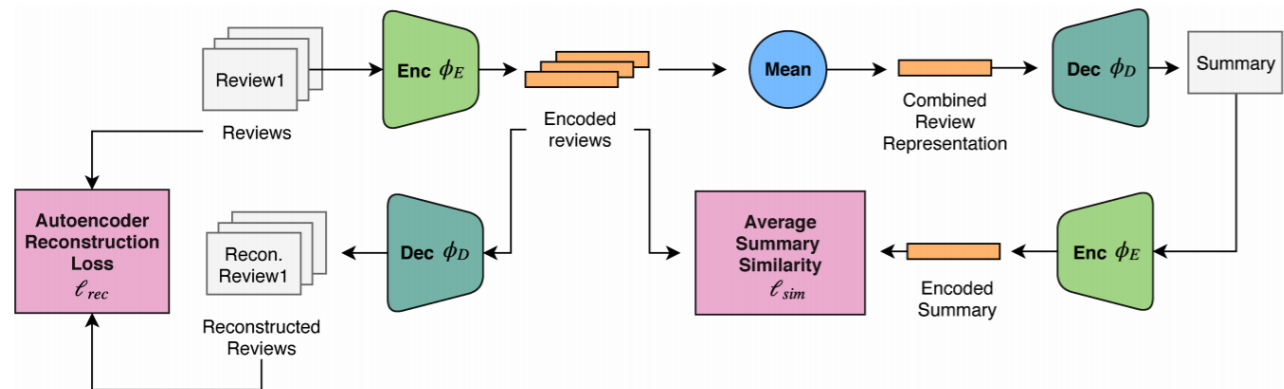
- Unsupervised summarization with language prior

[Baziotis, et al.,
NAACL 2019]



- Unsupervised multi-document summarization

[Chu, et al.,
ICML 2019]



Audio Style



male



female



positive
sentences



negative
sentences

Text Style Transfer



document



summary

Unsupervised Abstractive Summarization



Language 1

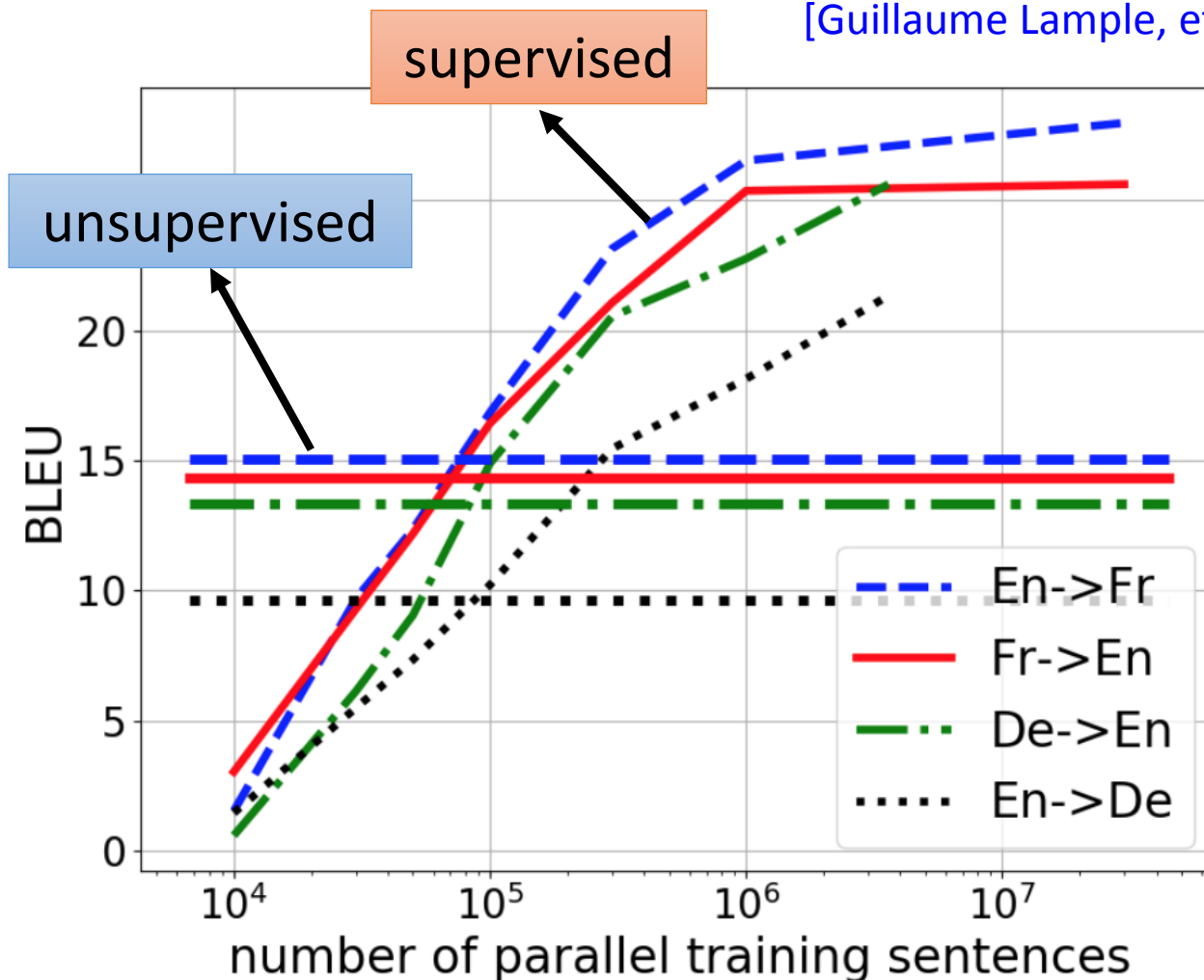


Language 2

Unsupervised Translation

[Alexis Conneau, et al., ICLR, 2018]

[Guillaume Lample, et al., ICLR, 2018]



Unsupervised learning
with 10M sentences

=

Supervised learning with
100K sentence pairs

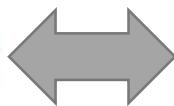
Audio Style



female



positive
sentences



negative sentences

Text Style Transfer



document

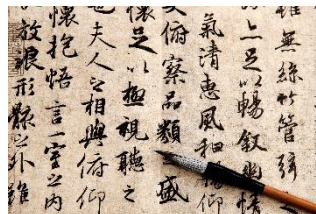


summary

Unsupervised Abstractive Summarization

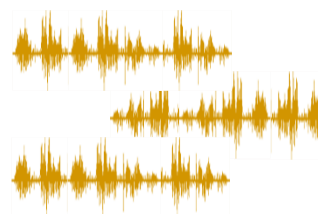


Language 1

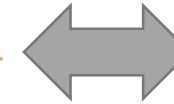


Language 2

Unsupervised Translation



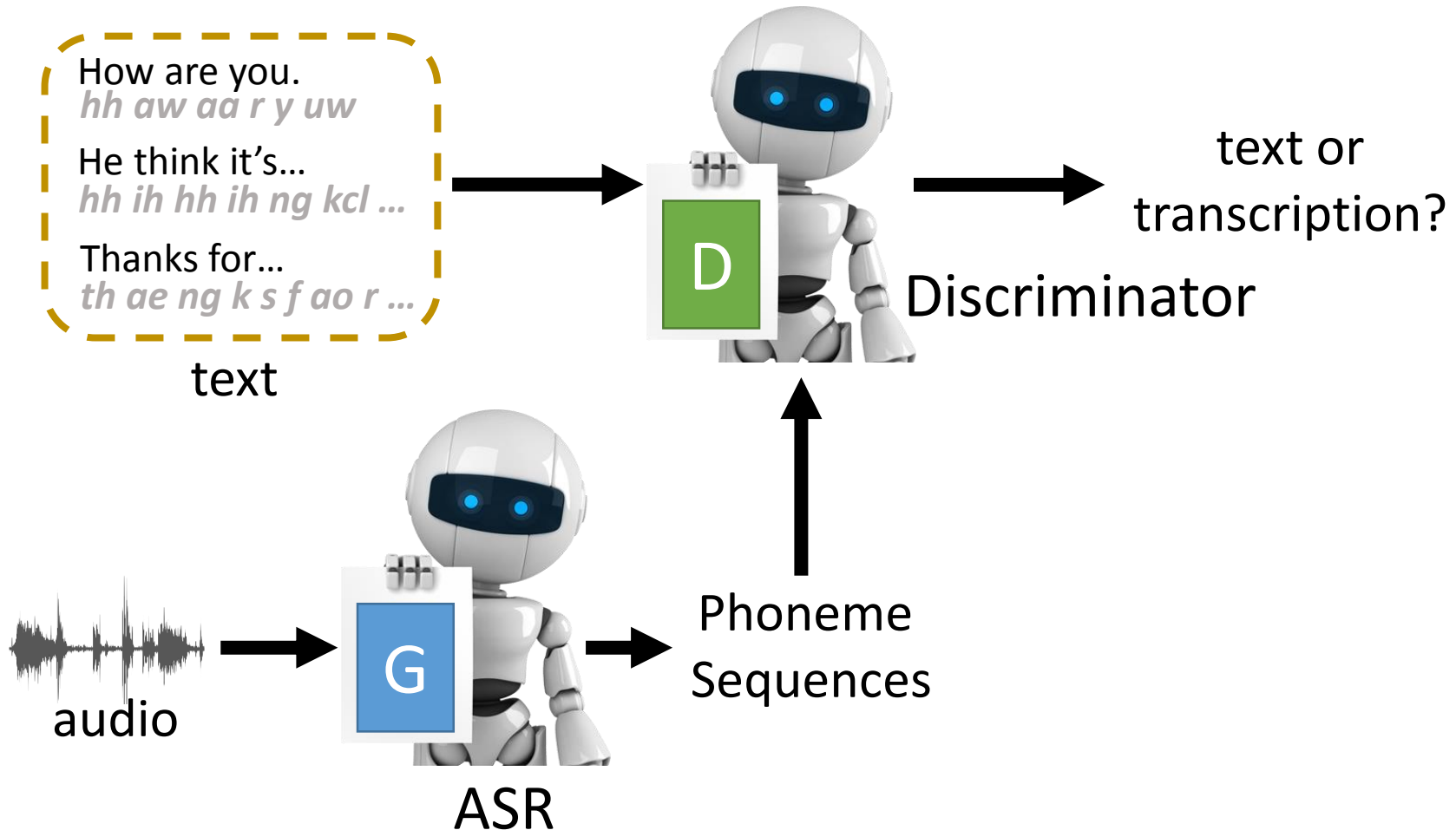
Audio



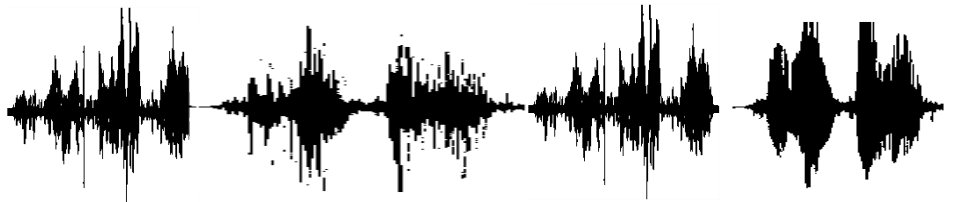
Text

Unsupervised ASR

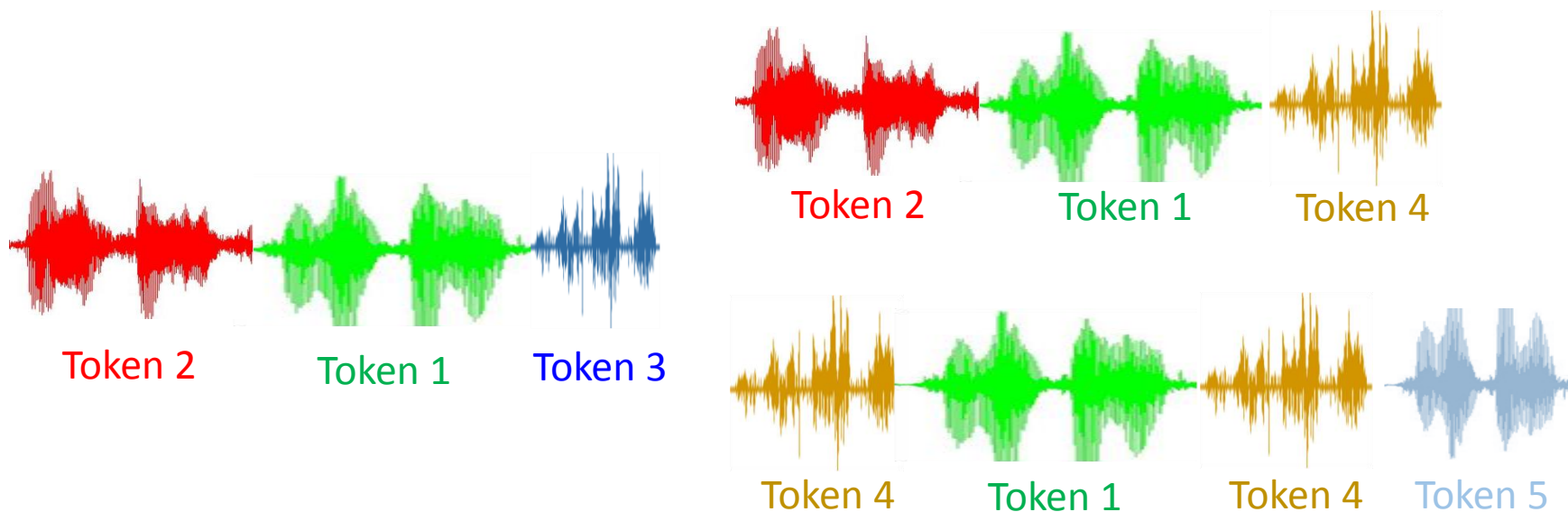
Unsupervised Speech Recognition



Acoustic Token Discovery



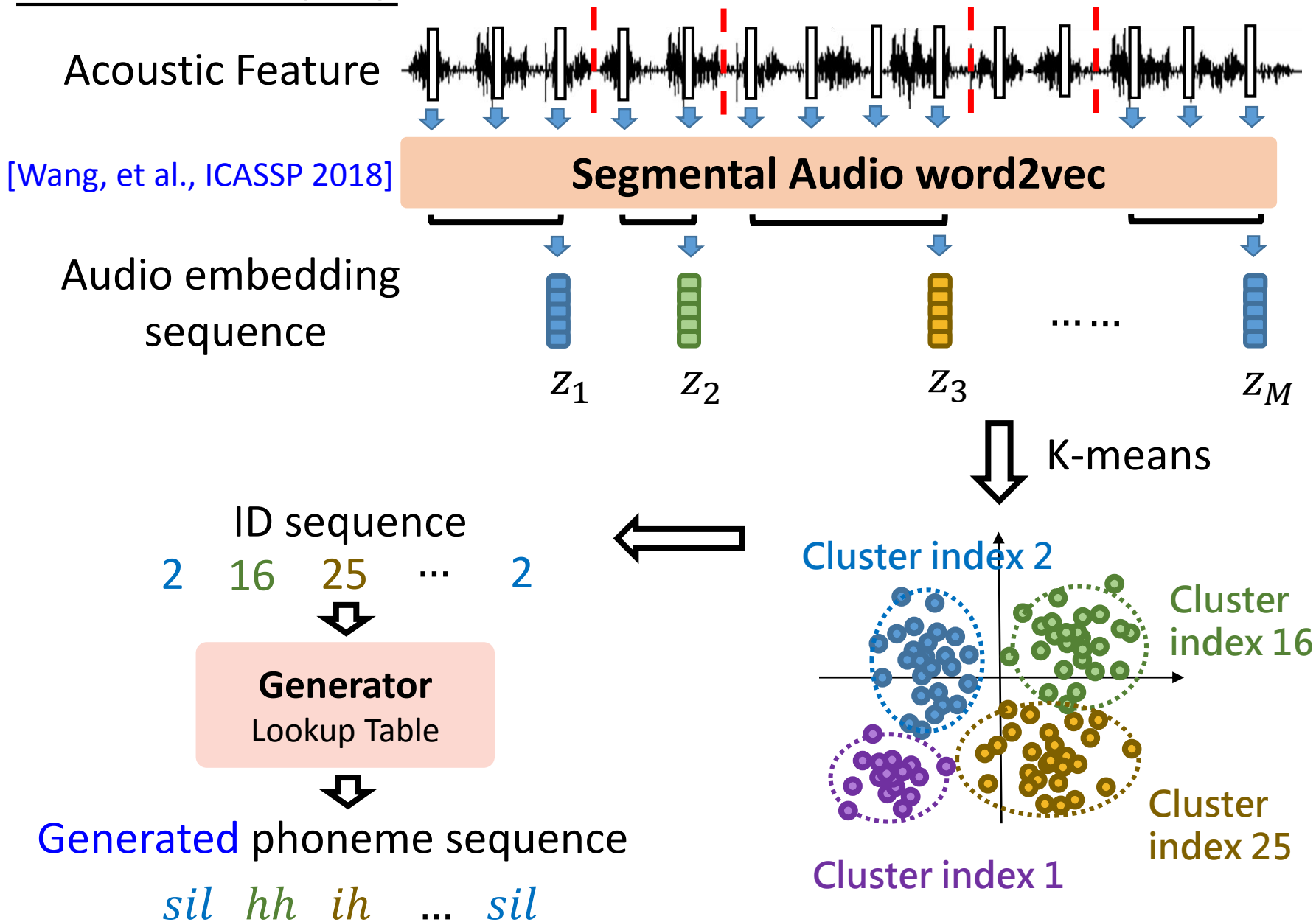
Acoustic Token Discovery



Acoustic tokens can be discovered from audio collection without text annotation.

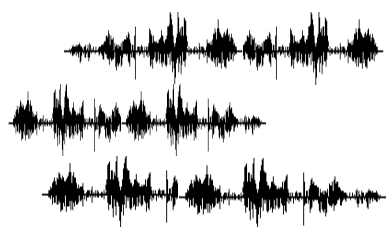
Acoustic tokens: chunks of acoustically similar audio segments with token IDs

Generator (v1)



Experiment

Matched Case (Oracle)

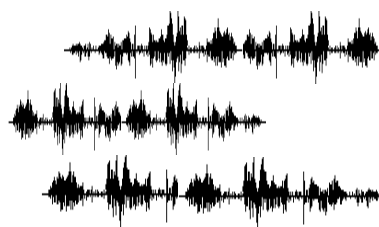


4620
(TIMIT)



4620
(TIMIT)

Nonmatched Case



3620
(TIMIT)



1000
(TIMIT)

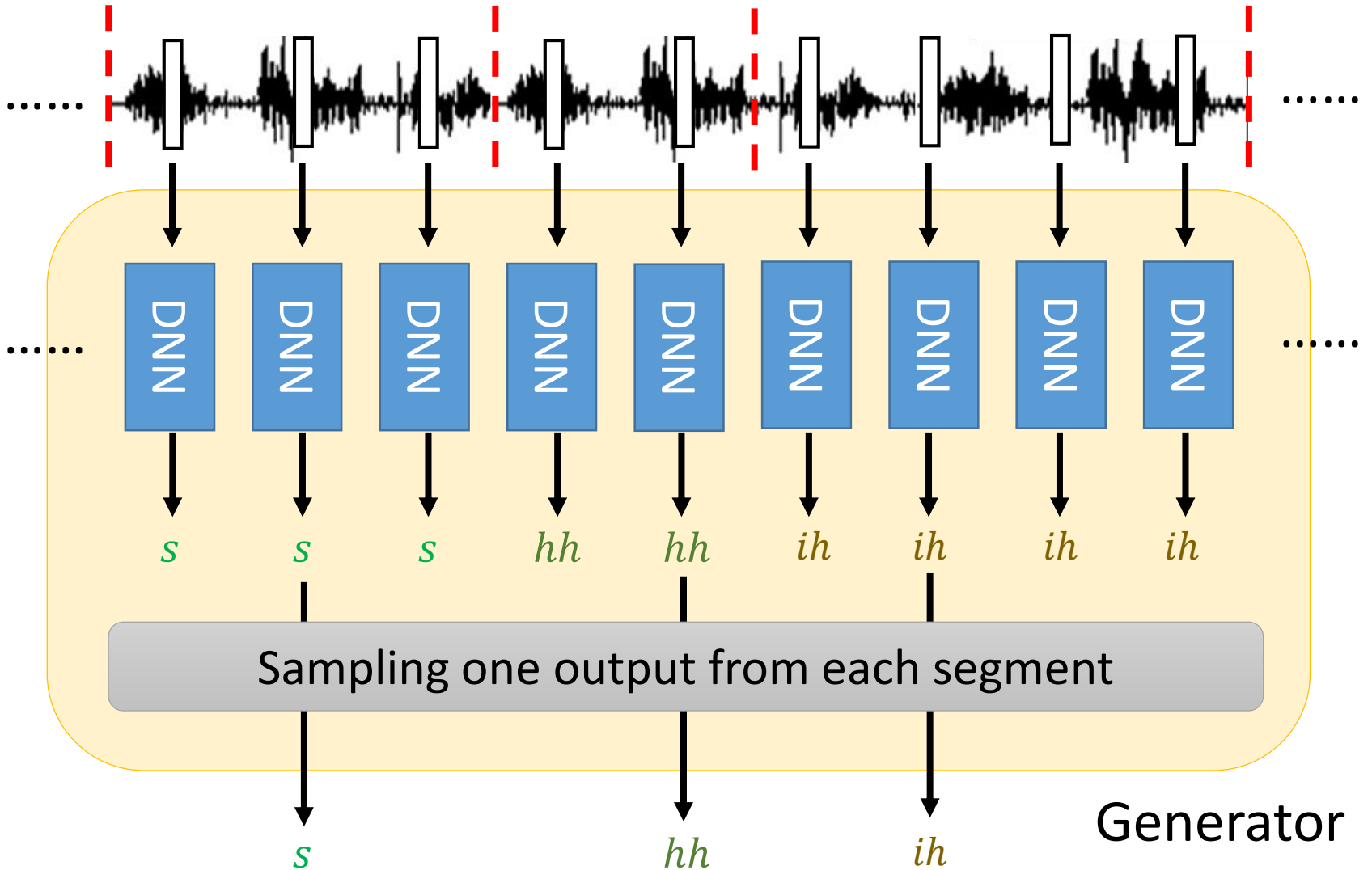
Experimental Results [\[Liu, et al., INTERSPEECH, 2018\]](#)

Approaches	PER	
	Matched	Nonmatched
Supervised		
RNN Transducer	17.7	-
Standard HMMs	21.5	-

Generator (v2)

Phoneme boundaries obtained by
Gate Activation Signals (GAS)

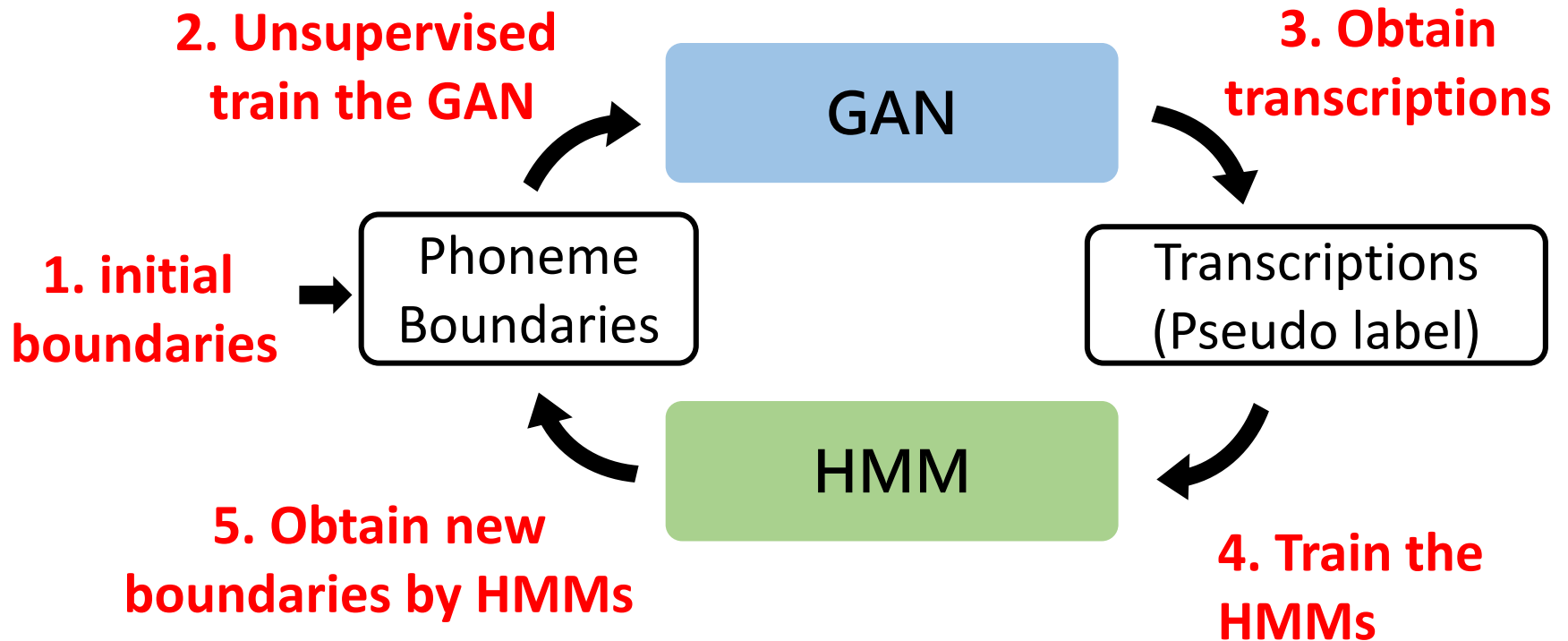
[Wang, et al., INTERSPEECH 2017]



Experimental Results [\[Chen, et al., INTERSPEECH, 2019\]](#)

Approaches			PER	
			Matched	Nonmatched
Supervised				
RNN Transducer			17.7	-
Standard HMMs			21.5	-
Completely unsupervised (no label at all)				
Generator (v1)			76.0	-
Generator (v2)	Iteration 1	GAN	48.6	50.0

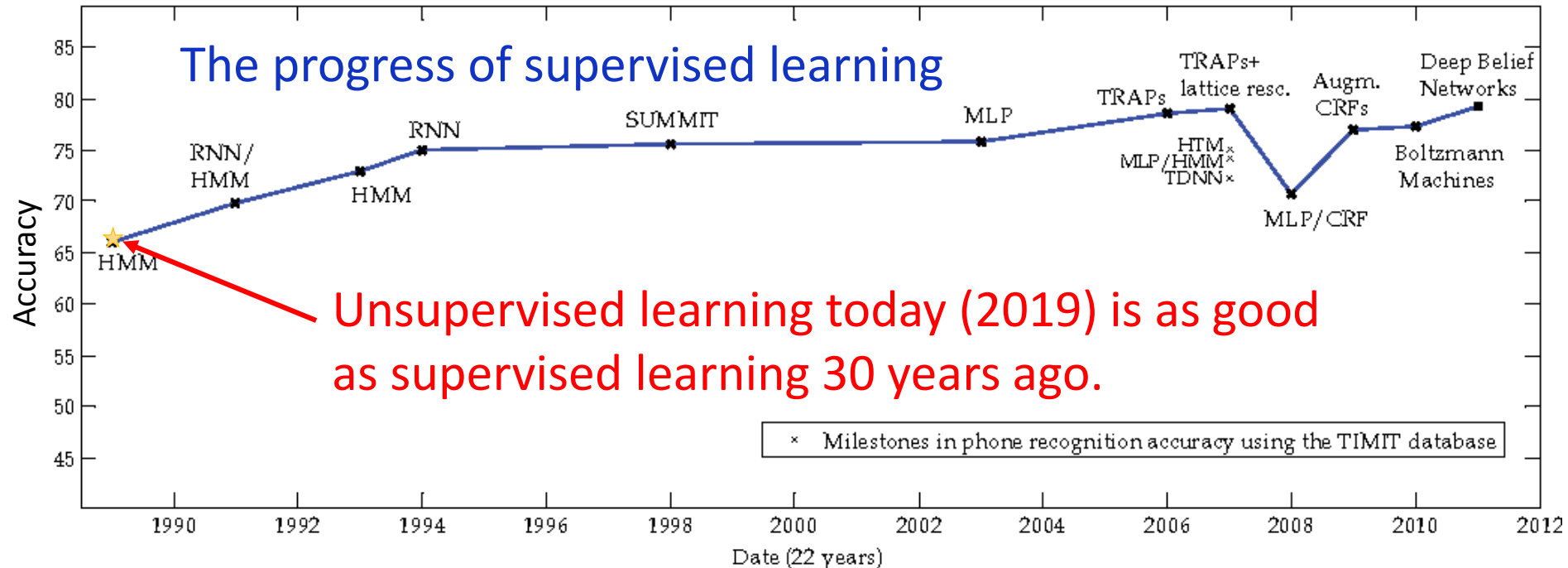
Refining Boundaries



Experimental Results [\[Chen, et al., INTERSPEECH, 2019\]](#)

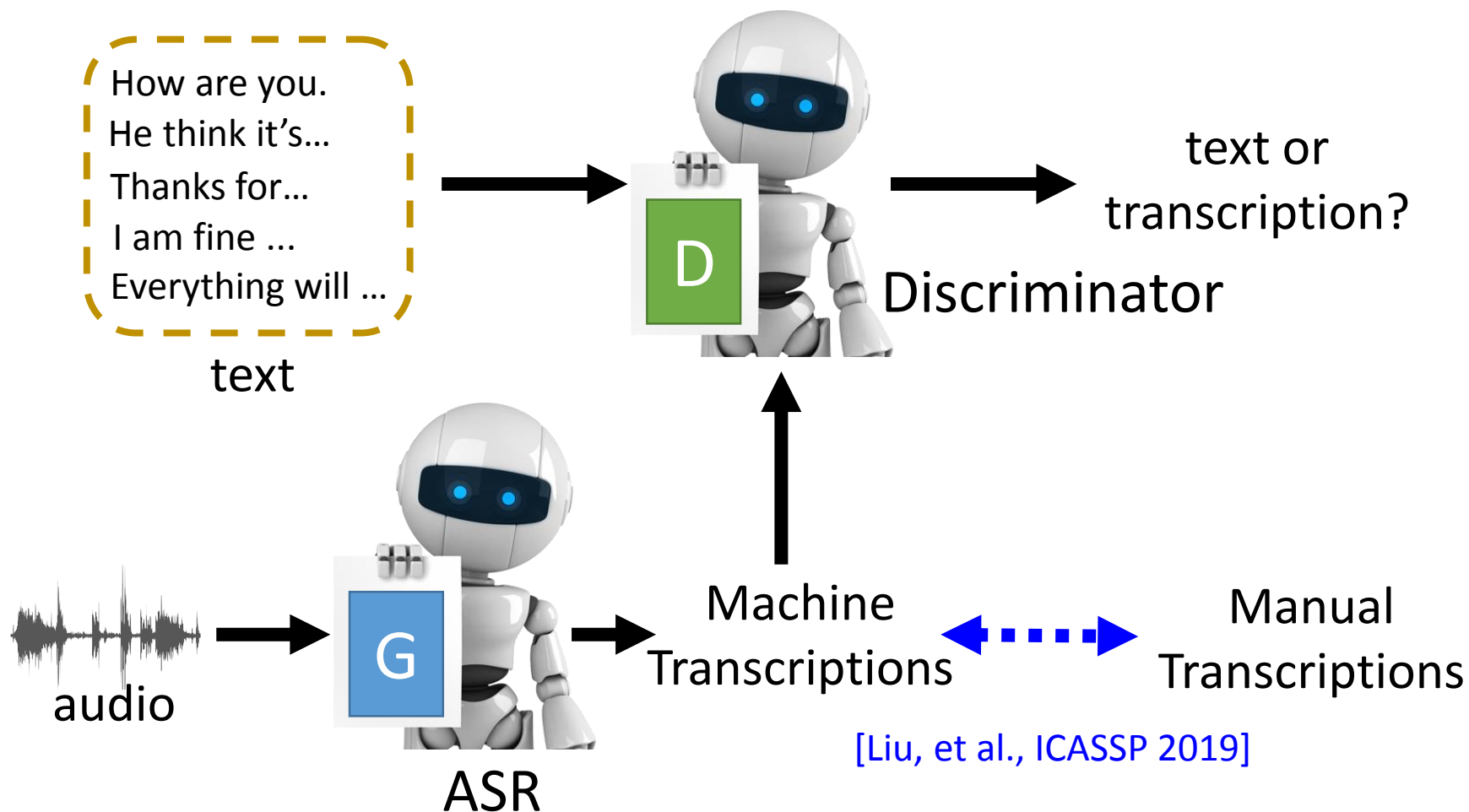
Approaches			PER	
			Matched	Nonmatched
Supervised				
RNN Transducer			17.7	-
Standard HMMs			21.5	-
Completely unsupervised (no label at all)				
Generator (v1)			76.0	-
Generator (v2)	Iteration 1	GAN	48.6	50.0
		HMM	30.7	39.5
	Iteration 2	GAN	41.0	44.3
		HMM	27.0	35.5
	Iteration 3	GAN	38.4	44.2
		HMM	26.1	33.1

The progress of supervised learning



The image is modified from: Phone recognition on the TIMIT database Lopes, C. and Perdigão, F., 2011. Speech Technologies, Vol 1, pp. 285--302.

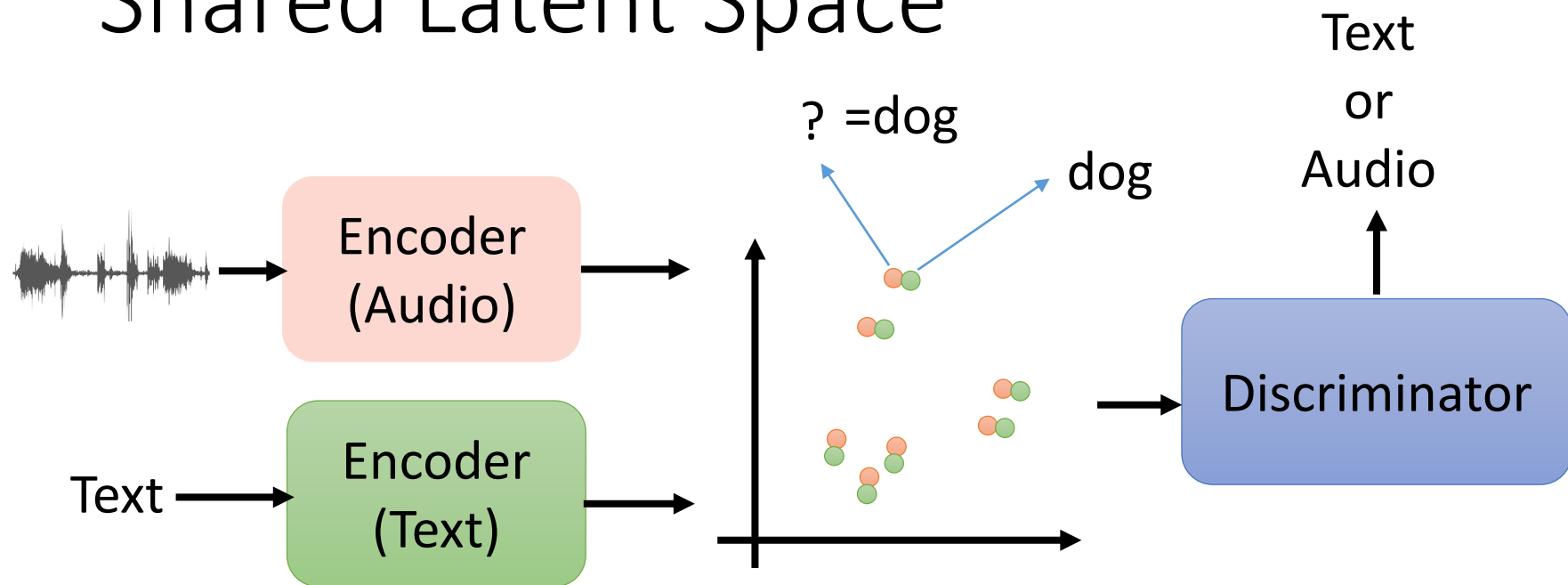
Semi-supervised Speech Recognition



Using 100 hours pairs annotated audio from Librispeech,
and text without audio

21.7% WER → 18.7% WER

Shared Latent Space



- Initial attempt [Chen, et al., SLT, 2018]
- 76.3% WER on Librispeech [Chung, et al., NIPS 2018]
- Unsupervised speech translation is possible [Chung, et al., ICASSP 2019]
- WSJ with 2.5 hours paired data: 64.6% WER
[Jennifer Drexler, et al., SLT 2018]
- LJ speech with 20 mins paired data: 11.7% PER [Ren, et al., ICML 2019]

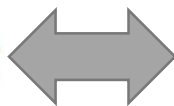
Audio Style



female



positive
sentences



negative
sentences

Text Style Transfer



document



summary

Unsupervised Abstractive Summarization

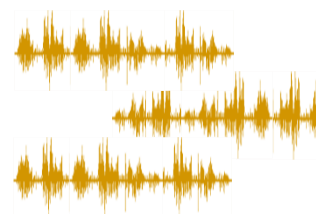


Language 1

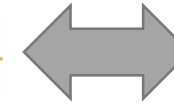


Language 2

Unsupervised Translation



Audio



Text

Unsupervised ASR

Reference

- **Unsupervised Speech Recognition**

- [Liu, et al., ICASSP'18] Alexander H. Liu, Hung-yi Lee, Lin-shan Lee, Adversarial Training of End-to-end Speech Recognition Using a Criticizing Language Model, ICASSP 2018
- [Liu, et al., INTERSPEECH'18] Da-Rong Liu, Kuan-Yu Chen, Hung-Yi Lee, Lin-shan Lee, Completely Unsupervised Phoneme Recognition by Adversarially Learning Mapping Relationships from Audio Embeddings, INTERSPEECH, 2018
- [Chen, et al., INTERSPEECH'19] Kuan-yu Chen, Che-ping Tsai, Da-Rong Liu, Hung-yi Lee and Lin-shan Lee, "Completely Unsupervised Phoneme Recognition By A Generative Adversarial Network Harmonized With Iteratively Refined Hidden Markov Models", INTERSPEECH, 2019
- [Chen, et al., SLT'18] Yi-Chen Chen, Sung-Feng Huang, Chia-Hao Shen, Hung-yi Lee, Lin-shan Lee, "Phonetic-and-Semantic Embedding of Spoken Words with Applications in Spoken Content Retrieval", SLT, 2018
- [Yeh, et al., ICLR'19] Chih-Kuan Yeh, Jianshu Chen, Chengzhu Yu, Dong Yu, Unsupervised Speech Recognition via Segmental Empirical Output Distribution Matching, ICLR, 2019

Reference

- **Unsupervised Speech Recognition**

- Takaaki Hori, Ramon Astudillo, Tomoki Hayashi, Yu Zhang, Shinji Watanabe, Jonathan Le Roux, Cycle-consistency training for end-to-end speech recognition, ICASSP 2019
- Murali Karthick Baskar, Shinji Watanabe, Ramon Astudillo, Takaaki Hori, Lukáš Burget, Jan Černocký, Semi-supervised Sequence-to-sequence ASR using Unpaired Speech and Text, INTERSPEECH 2019
- Andros Tjandra, Sakriani Sakti, Satoshi Nakamura, Listening while Speaking: Speech Chain by Deep Learning, ASRU 2017
- [\[Chung, et al., NIPS 2018\]](#) Yu-An Chung, Wei-Hung Weng, Schrasing Tong, James Glass, Unsupervised Cross-Modal Alignment of Speech and Text Embedding Spaces, NIPS, 2018
- [\[Chung, et al., ICASSP 2019\]](#) Yu-An Chung, Wei-Hung Weng, Schrasing Tong, James Glass, Towards Unsupervised Speech-to-Text Translation, ICASSP 2019
- [\[Ren, et al., ICML 2019\]](#) Yi Ren, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, Tie-Yan Liu, Almost Unsupervised Text to Speech and Automatic Speech Recognition, ICML 2019

Reference

- **Unsupervised Speech Recognition**

- Shigeki Karita , Shinji Watanabe, Tomoharu Iwata, Atsunori Ogawa, Marc Delcroix, Semi-Supervised End-to-End Speech Recognition, INTERSPEECH, 2018
- [\[Jennifer Drexler, et al., SLT 2018\]](#) Jennifer Drexler, James R. Glass, “Combining End-to-End and Adversarial Training for Low-Resource Speech Recognition”, SLT 2018
- Tomoki Hayashi, Shinji Watanabe, Yu Zhang, Tomoki Toda, Takaaki Hori, Ramon Astudillo, Kazuya Takeda, Back-Translation-Style Data Augmentation for End-to-End ASR, SLT, 2018

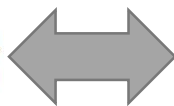
Audio Style



female



positive
sentences



negative
sentences

Text Style Transfer



document

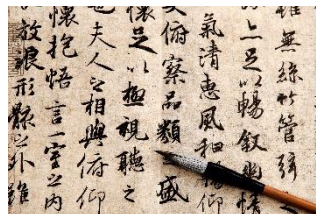
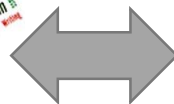


summary

Unsupervised Abstractive Summarization

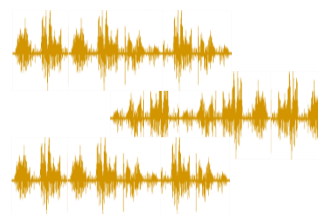


Language 1

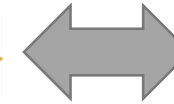


Language 2

Unsupervised Translation



Audio



Text

Unsupervised ASR

Reference

- [Lee, et al., ICASSP'18] Chih-Wei Lee, Yau-Shian Wang, Tsung-Yuan Hsu, Kuan-Yu Chen, Hung-Yi Lee, Lin-shan Lee, Scalable Sentiment for Sequence-to-sequence Chatbot Response with Performance Analysis, ICASSP, 2018
- [Ning Dai, et al., ACL'19] Ning Dai, Jianze Liang, Xipeng Qiu, Xuanjing Huang, Style Transformer: Unpaired Text Style Transfer without Disentangled Latent Representation, ACL, 2019
- [Lample, et al., ICLR'19] Guillaume Lample, Sandeep Subramanian, Eric Smith, [Ludovic Denoyer](#), [Marc'Aurelio Ranzato](#), Y-Lan Boureau, Multiple-Attribute Text Rewriting, ICLR, 2019

Reference

- [Hu, et al., ICML'17] Zhiting Hu, [Zichao Yang](#), [Xiaodan Liang](#), [Ruslan Salakhutdinov](#), [Eric P. Xing](#), **Toward Controlled Generation of Text**, ICML, 2017
- [Fu, et al., AACL'17] Zhenxin Fu, Xiaoye Tan, Nanyun Peng, Dongyan Zhao, and Rui Yan, "Style transfer in text: Exploration and evaluation," AACL, 2017.
- [Shen, et al., NIPS'17] Tianxiao Shen, [Tao Lei](#), [Regina Barzilay](#), [Tommi Jaakkola](#), **Style Transfer from Non-Parallel Text by Cross-Alignment**, NIPS, 2017
- [Li, et al., NAACL'18] Juncen Li, [Robin Jia](#), [He He](#), [Percy Liang](#), [Delete, Retrieve, Generate: a Simple Approach to Sentiment and Style Transfer](#), NAACL, 2018
- [Xu, et al., ACL'18] Jingjing Xu, [Xu Sun](#), [Qi Zeng](#), [Xuancheng Ren](#), [Xiaodong Zhang](#), [Houfeng Wang](#), [Wenjie Li](#), **Unpaired Sentiment-to-Sentiment Translation: A Cycled Reinforcement Learning Approach**, ACL, 2018

Reference

- [Wang, Lee, EMNLP'18] Yau-Shian Wang, Hung-Yi Lee, "Learning to Encode Text as Human-Readable Summaries using Generative Adversarial Networks", EMNLP, 2018
- [Chu, et al., ICML'19] Eric Chu, Peter Liu, "MeanSum: A Neural Model for Unsupervised Multi-Document Abstractive Summarization", ICML, 2019
- [Baziotis, et al., NAACL'19] Christos Baziotis, Ion Androutsopoulos, Ioannis Konstas, Alexandros Potamianos, "SEQ³: Differentiable Sequence-to-Sequence-to-Sequence Autoencoder for Unsupervised Abstractive Sentence Compression", NAACL 2019

Reference

- [Alexis Conneau, et al., ICLR'18] Alexis Conneau, Guillaume Lample, Marc'Aurelio Ranzato, Ludovic Denoyer, Hervé Jégou, Word Translation Without Parallel Data, ICLR 2018
- [Guillaume Lample, et al., ICLR'18] Guillaume Lample, Ludovic Denoyer, Marc'Aurelio Ranzato, Unsupervised Machine Translation Using Monolingual Corpora Only, ICLR, 2018

Reference

- [\[Li, et al., EMNLP, 2017\]](#) Jiwei Li, Will Monroe, Tianlin Shi, Sébastien Jean, Alan Ritter, Dan Jurafsky, Adversarial Learning for Neural Dialogue Generation, EMNLP, 2017
- [\[Matt J. Kusner, et al., arXiv, 2016\]](#) Matt J. Kusner, José Miguel Hernández-Lobato, GANS for Sequences of Discrete Elements with the Gumbel-softmax Distribution, arXiv 2016
- [\[Tong Che, et al, arXiv, 2017\]](#) Tong Che, Yanran Li, Ruixiang Zhang, R Devon Hjelm, Wenjie Li, Yangqiu Song, Yoshua Bengio, Maximum-Likelihood Augmented Discrete Generative Adversarial Networks, arXiv 2017
- [\[Yu, et al., AAAI, 2017\]](#) Lantao Yu, Weinan Zhang, Jun Wang, Yong Yu, SeqGAN: Sequence Generative Adversarial Nets with Policy Gradient, AAAI 2017
- [\[Sai Rajeswar, et al., arXiv, 2017\]](#) Sai Rajeswar, Sandeep Subramanian, Francis Dutil, Christopher Pal, Aaron Courville, Adversarial Generation of Natural Language, arXiv, 2017
- [\[Ofir Press, et al., ICML workshop, 2017\]](#) Ofir Press, Amir Bar, Ben Bogin, Jonathan Berant, Lior Wolf, Language Generation with Recurrent Generative Adversarial Networks without Pre-training, ICML workshop, 2017

Reference

- [\[Zhen Xu, et al., EMNLP, 2017\]](#) Zhen Xu, Bingquan Liu, Baoxun Wang, Chengjie Sun, Xiaolong Wang, Zhuoran Wang, Chao Qi , Neural Response Generation via GAN with an Approximate Embedding Layer, EMNLP, 2017
- [\[Alex Lamb, et al., NIPS, 2016\]](#) Alex Lamb, Anirudh Goyal, Ying Zhang, Saizheng Zhang, Aaron Courville, Yoshua Bengio, Professor Forcing: A New Algorithm for Training Recurrent Networks, NIPS, 2016
- [\[Yizhe Zhang, et al., ICML, 2017\]](#) Yizhe Zhang, Zhe Gan, Kai Fan, Zhi Chen, Ricardo Henao, Dinghan Shen, Lawrence Carin, Adversarial Feature Matching for Text Generation, ICML, 2017
- [\[Jiaxian Guo, et al., AACL, 2018\]](#) Jiaxian Guo, Sidi Lu, Han Cai, Weinan Zhang, Yong Yu, Jun Wang, Long Text Generation via Adversarial Training with Leaked Information, AACL, 2018
- [\[Kevin Lin, et al, NIPS, 2017\]](#) Kevin Lin, Dianqi Li, Xiaodong He, Zhengyou Zhang, Ming-Ting Sun, Adversarial Ranking for Language Generation, NIPS, 2017
- [\[William Fedus, et al., ICLR, 2018\]](#) William Fedus, Ian Goodfellow, Andrew M. Dai, MaskGAN: Better Text Generation via Filling in the _____, ICLR, 2018
- [\[Cyprien de Masson d'Autume, et al., arXiv 2019\]](#) Cyprien de Masson d'Autume, Mihaela Rosca, Jack Rae, Shakir Mohamed, Training language GANs from Scratch, arXiv 2019

Reference

- [\[Tua, Lee, TASLP, 2019\]](#) Yi-Lin Tuan, Hung-Yi Lee, Improving Conditional Sequence Generative Adversarial Networks by Stepwise Evaluation, TASLP, 2019
- [\[Xu, et al., EMNLP, 2018\]](#) Jingjing Xu, Xuancheng Ren, Junyang Lin, Xu Sun, Diversity-Promoting GAN: A Cross-Entropy Based Generative Adversarial Network for Diversified Text Generation, EMNLP, 2018
- [\[Weili Nie, et al. ICLR, 2019\]](#) Weili Nie, Nina Narodytska, Ankit Patel, RelGAN: Relational Generative Adversarial Networks for Text Generation, ICLR 2019