

RL and GAN for Sentence Generation and Chat-bot

Hung-yi Lee

Outline

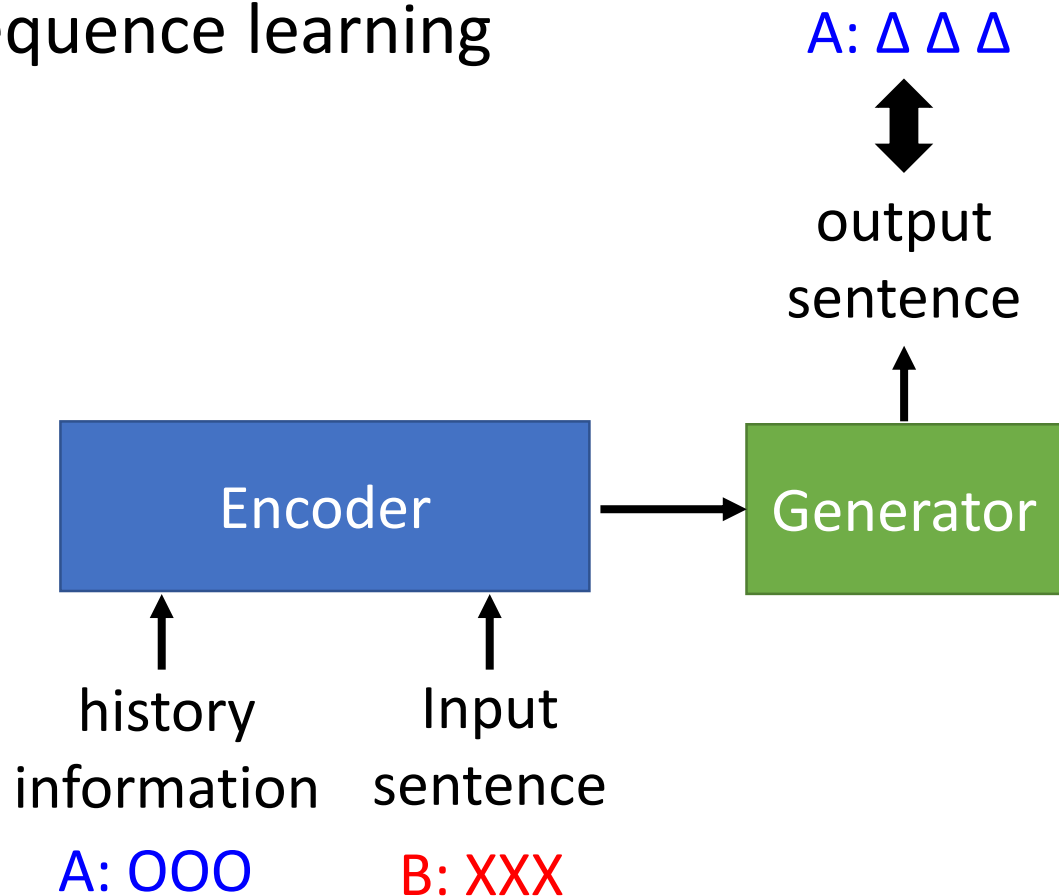
- Policy Gradient
- SeqGAN
 - Two techniques: MCMC, partial
 - Experiments: SeqGAN and dialogue
- Original GAN
 - MadliGAN
 - Gumbel

Review: Chat-bot

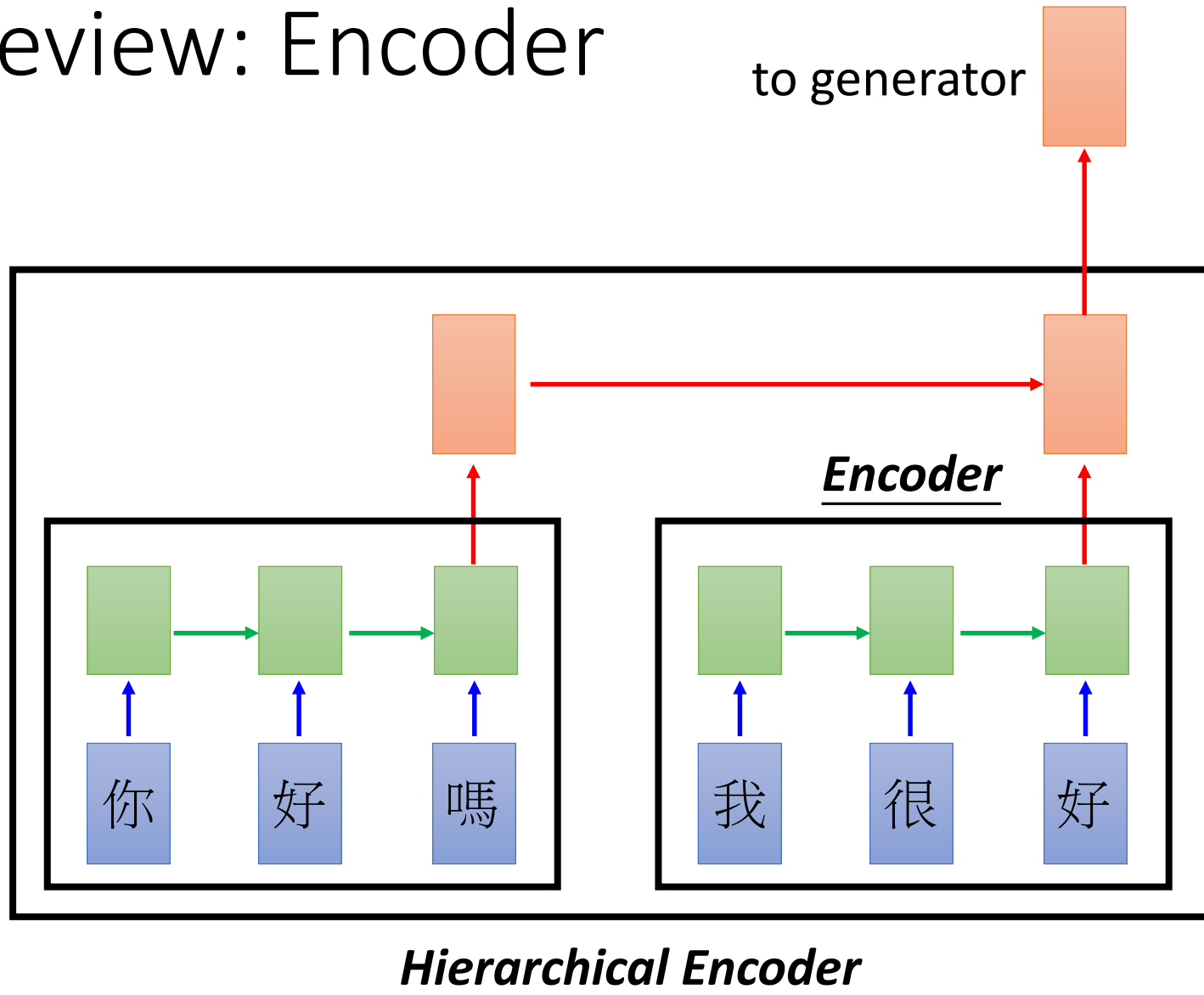
- Sequence-to-sequence learning

Training data:

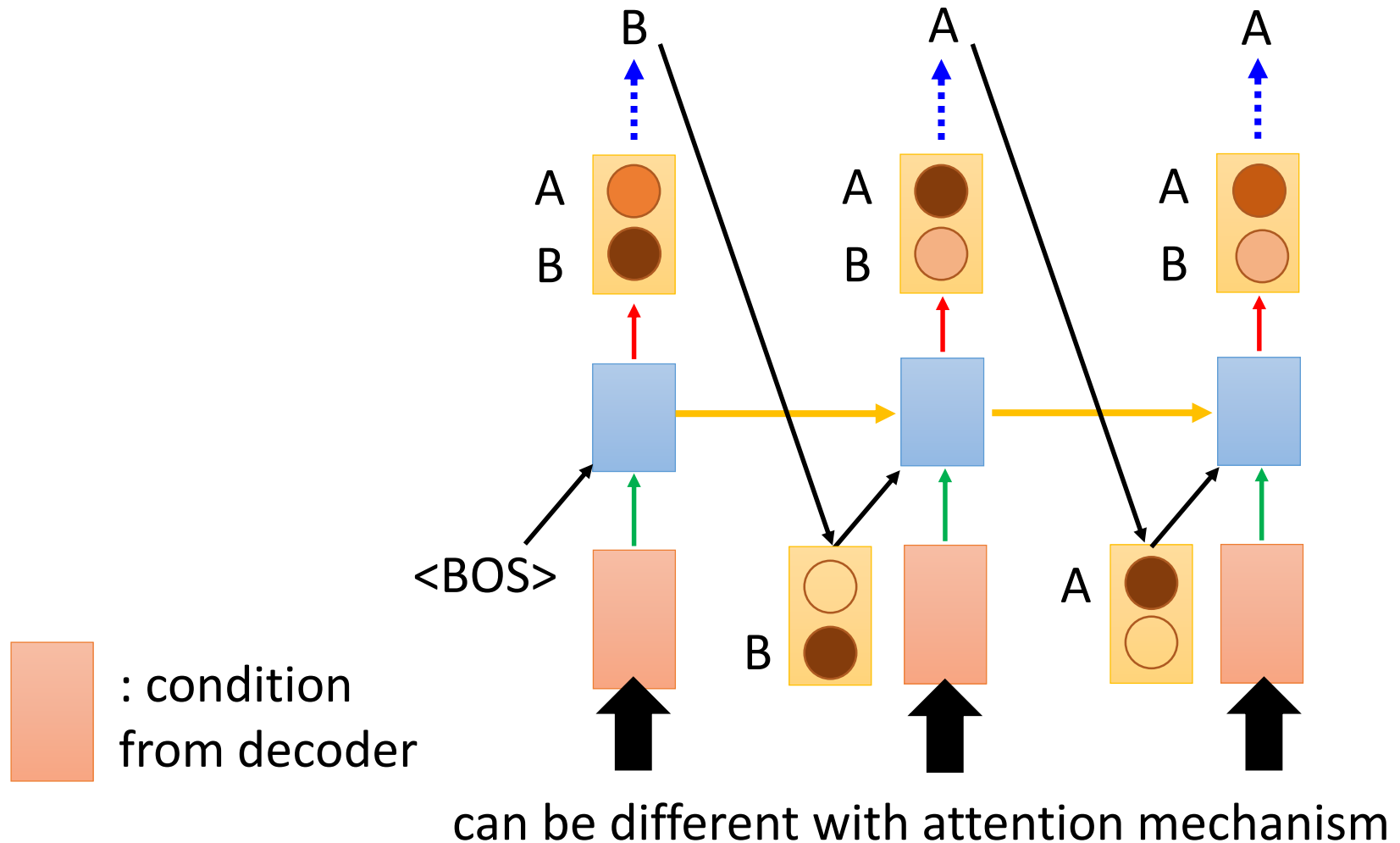
⋮
A: 000
B: XXX
A: Δ Δ Δ
⋮



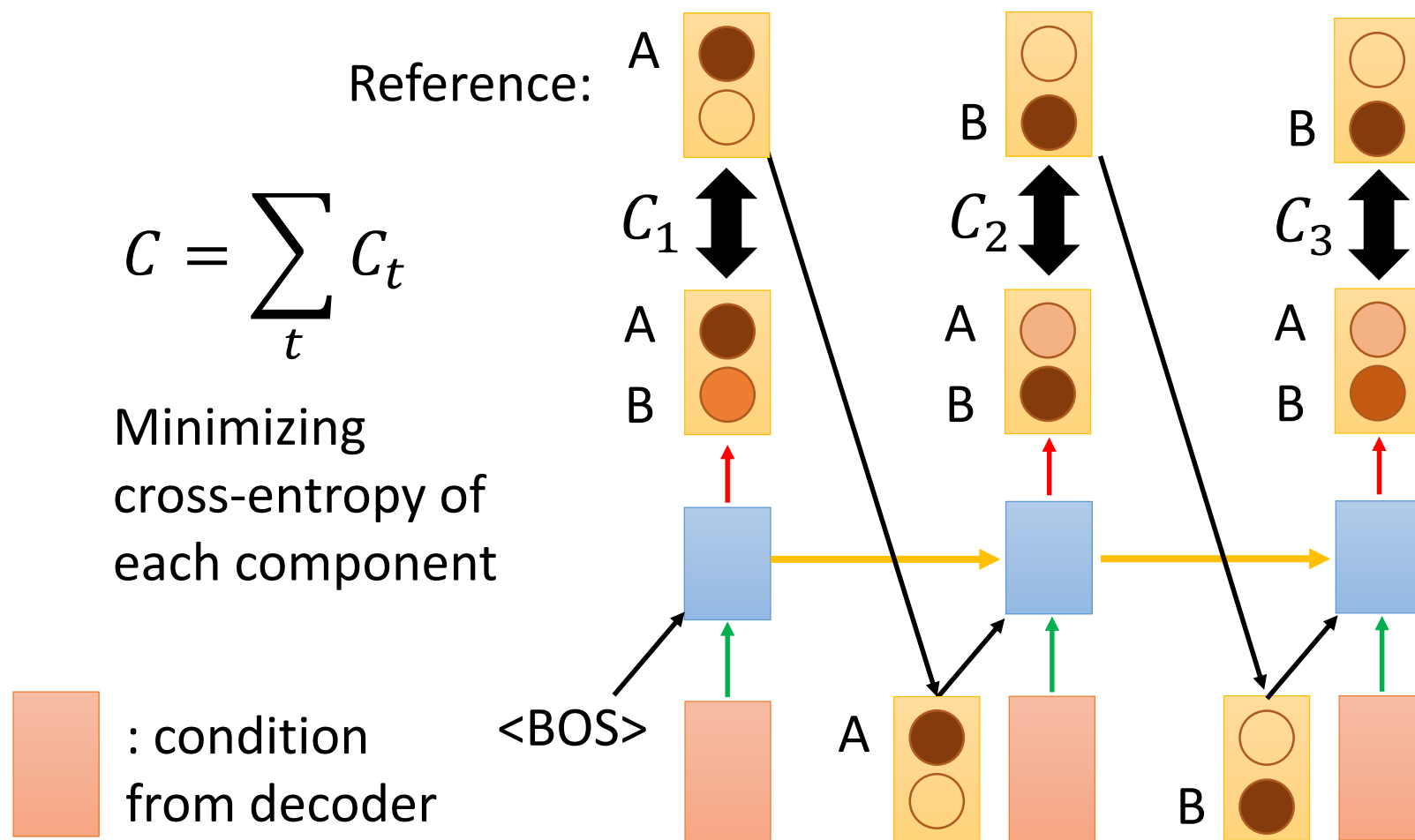
Review: Encoder



Review: Generator



Review: Training Generator



Review: Training Generator

Training data: (h, \hat{x})

h : input sentence and history/context

\hat{x} : correct response (word sequence)

\hat{x}_t : t-th word, $\hat{x}_{1:t}$: first t words of \hat{x}

$$C = \sum_t C_t$$

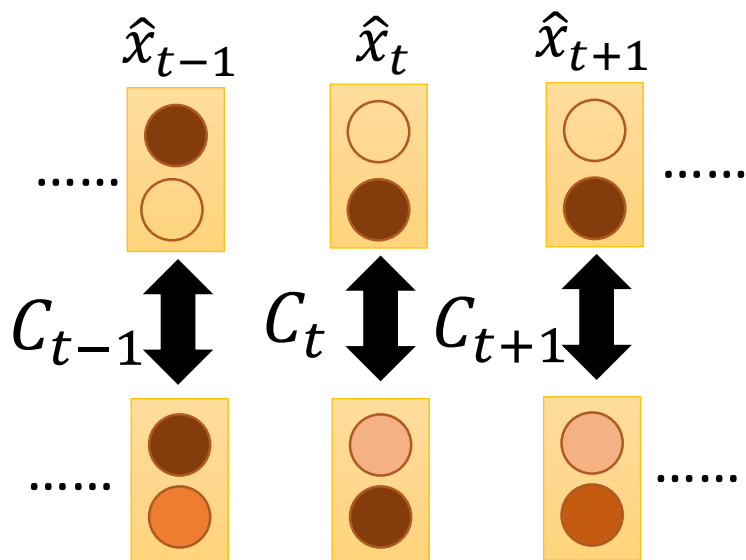
$$C_t = -\log P_\theta(\hat{x}_t | \hat{x}_{1:t-1}, h)$$

$$C = -\sum_t \log P(\hat{x}_t | \hat{x}_{1:t-1}, h)$$

$$= -\log P(\hat{x}_1 | h) P(\hat{x}_t | \hat{x}_{1:t-1}, h)$$

$$\cdots P(\hat{x}_T | \hat{x}_{1:T-1}, h)$$

$$= -\log P(\hat{x} | h)$$



generator output

Maximizing the likelihood of generating \hat{x} given h

RL for Sentence Generation

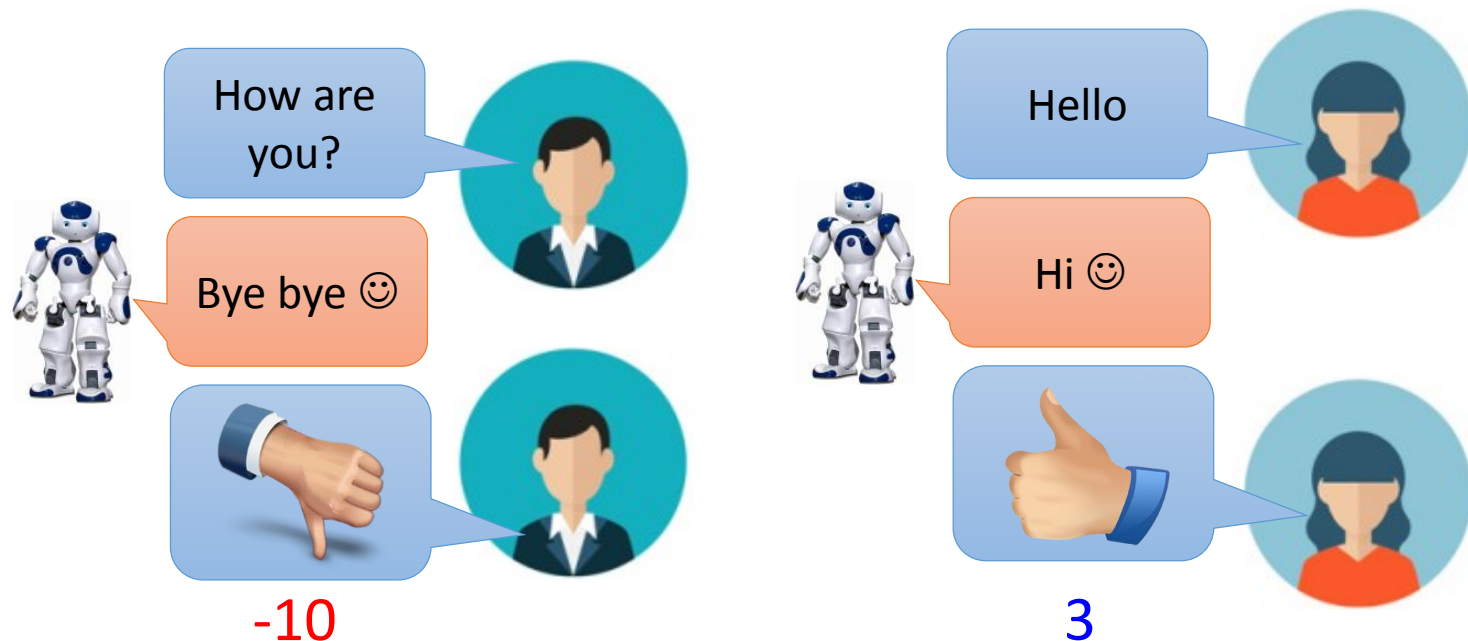
Jiwei Li, Will Monroe, Alan Ritter, Michel Galley, Jianfeng Gao, Dan Jurafsky,
"Deep Reinforcement Learning for Dialogue Generation", EMNLP 2016

https://image.freepik.com/free-vector/variety-of-human-avatars_23-2147506285.jpg

http://www.freepik.com/free-vector/variety-of-human-avatars_766615.htm

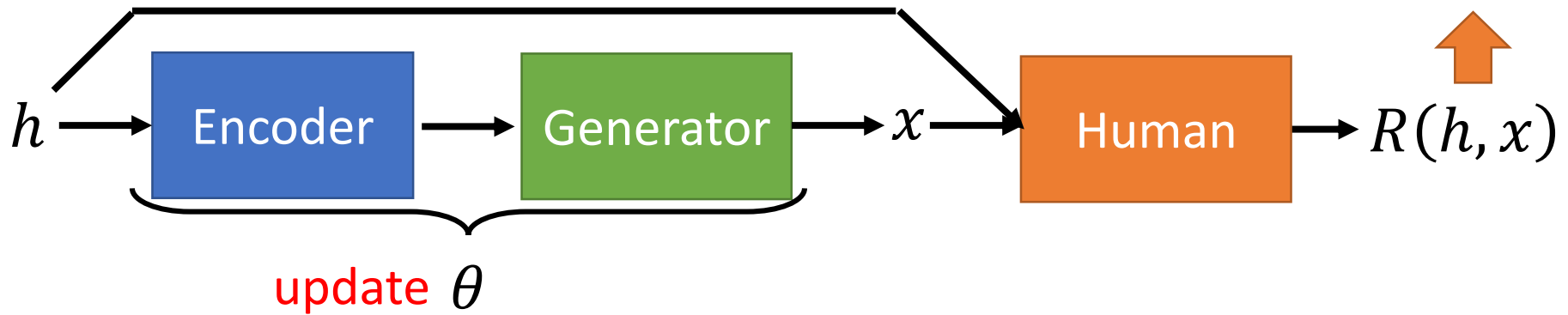
Introduction

- Machine obtains feedback from user



- Chat-bot learns to maximize the *expected reward*

Maximizing Expected Reward



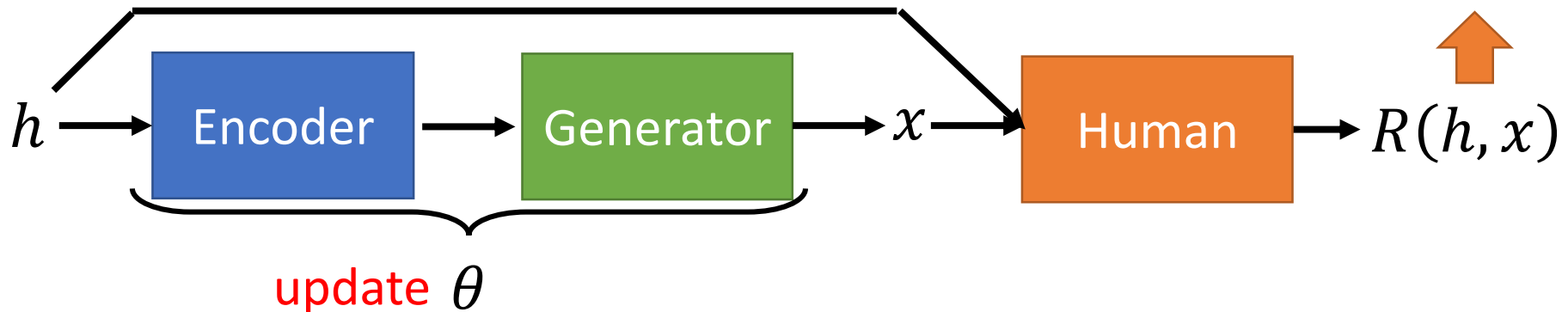
$$\theta^* = \arg \max_{\theta} \bar{R}_{\theta} \quad \leftarrow \text{Maximizing expected reward}$$

$$\bar{R}_{\theta} = \sum_h \underbrace{P(h)} \sum_x R(h, x) \underbrace{P_{\theta}(x|h)}$$

Randomness in generator

Probability that the input/history is h

Maximizing Expected Reward



$$\theta^* = \arg \max_{\theta} \bar{R}_{\theta} \quad \leftarrow \text{Maximizing expected reward}$$

$$\begin{aligned} \bar{R}_{\theta} &= \sum_h P(h) \sum_x R(h, x) P_{\theta}(x|h) = E_{h \sim P(h)} \left[E_{x \sim P_{\theta}(x|h)} [R(h, x)] \right] \\ &= E_{h \sim P(h), x \sim P_{\theta}(x|h)} [R(h, x)] \approx \frac{1}{N} \sum_{i=1}^N R(h^i, x^i) \end{aligned}$$

Sample: $(h^1, x^1), (h^2, x^2), \dots, (h^N, x^N)$

Where
is θ ?

Policy Gradient

$$\frac{d \log(f(x))}{dx} = \frac{1}{f(x)} \frac{df(x)}{dx}$$

$$\bar{R}_\theta = \sum_h P(h) \sum_x R(h, x) P_\theta(x|h) \approx \frac{1}{N} \sum_{i=1}^N R(h^i, x^i)$$

$$\nabla \bar{R}_\theta = \sum_h P(h) \sum_x R(h, x) \nabla P_\theta(x|h) \approx \frac{1}{N} \sum_{i=1}^N R(h^i, x^i) \nabla \log P_\theta(x|h)$$

$$= \sum_h P(h) \sum_x R(h, x) P_\theta(x|h) \boxed{\frac{\nabla P_\theta(x|h)}{P_\theta(x|h)}}$$


Sampling

$$= \sum_h P(h) \sum_x R(h, x) P_\theta(x|h) \boxed{\nabla \log P_\theta(x|h)}$$

$$= E_{h \sim P(h), x \sim P_\theta(x|h)} [R(h, x) \nabla \log P_\theta(x|h)]$$

Policy Gradient

- Gradient Ascent

$$\theta^{new} \leftarrow \theta^{old} + \eta \nabla \bar{R}_{\theta^{old}}$$

$$\nabla \bar{R}_{\theta} \approx \frac{1}{N} \sum_{i=1}^N R(h^i, x^i) \nabla \log P_{\theta}(x^i | h^i)$$

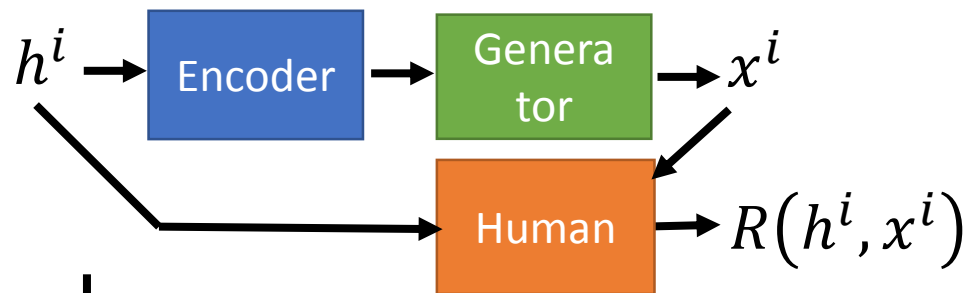
$R(h^i, x^i)$ is positive

➡ After updating θ , $P_{\theta}(x^i | h^i)$ will increase

$R(h^i, x^i)$ is negative

➡ After updating θ , $P_{\theta}(x^i | h^i)$ will decrease

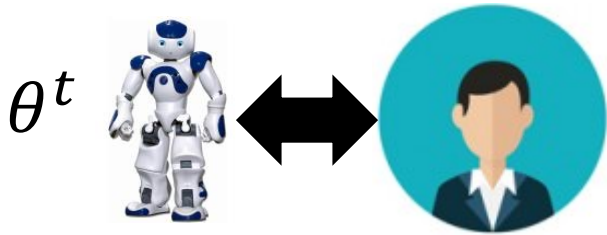
Implementation



	Maximum Likelihood	Reinforcement Learning
Objective Function	$\frac{1}{N} \sum_{i=1}^N \log P_{\theta}(\hat{x}^i h^i)$	$\frac{1}{N} \sum_{i=1}^N R(h^i, x^i) \log P_{\theta}(x^i h^i)$
Gradient	$\frac{1}{N} \sum_{i=1}^N \nabla \log P_{\theta}(\hat{x}^i h^i)$	$\frac{1}{N} \sum_{i=1}^N R(h^i, x^i) \nabla \log P_{\theta}(x^i h^i)$
Training Data	$\{(h^1, \hat{x}^1), \dots, (h^N, \hat{x}^N)\}$ $R(h^i, \hat{x}^i) = 1$	$\{(h^1, x^1), \dots, (h^N, x^N)\}$ <p>Sampling as training data weighted by $R(h^i, x^i)$</p>

Implementation

θ^0 can be well pre-trained from
 $\{(h^1, \hat{x}^1), \dots, (h^N, \hat{x}^N)\}$



$$\begin{array}{ll} (h^1, x^1) & R(h^1, x^1) \\ (h^2, x^2) & R(h^2, x^2) \\ \vdots & \vdots \\ (h^N, x^N) & R(h^N, x^N) \end{array}$$

New Objective:

$$\frac{1}{N} \sum_{i=1}^N R(h^i, x^i) \log P_{\theta}(x^i | h^i)$$

$$\theta^{t+1} \leftarrow \theta^t + \eta \nabla \bar{R}_{\theta^t}$$

$$\frac{1}{N} \sum_{i=1}^N R(h^i, x^i) \nabla \log P_{\theta^t}(x^i | h^i)$$

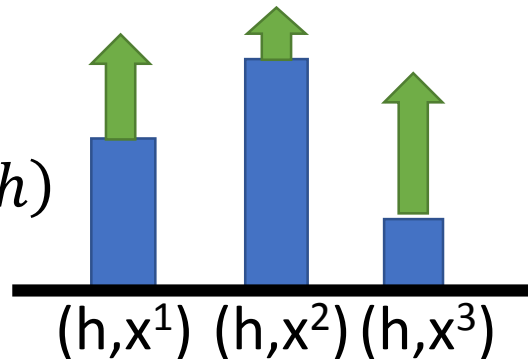
Add a Baseline

If $R(h^i, x^i)$ is always positive

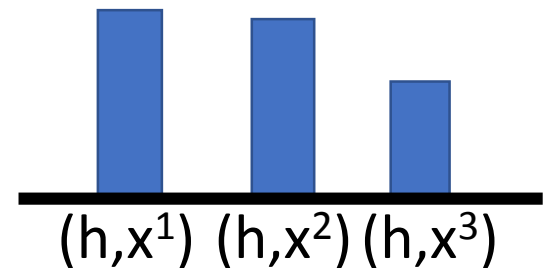
$$\frac{1}{N} \sum_{i=1}^N R(h^i, x^i) \log \nabla P_{\theta}(x^i | h^i)$$

Ideal case

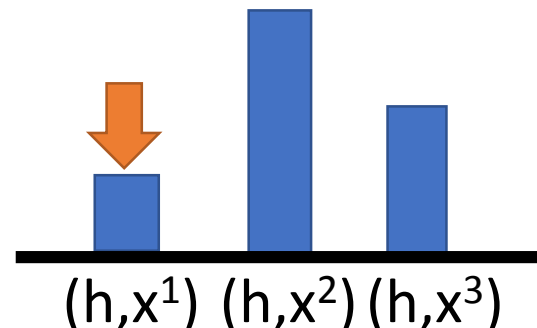
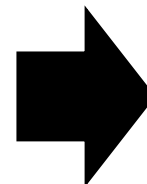
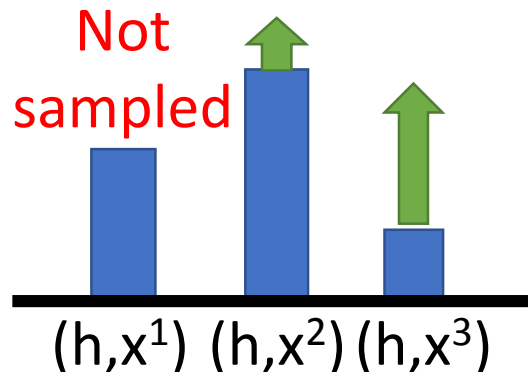
$P_{\theta}(x|h)$



Because it is probability ...



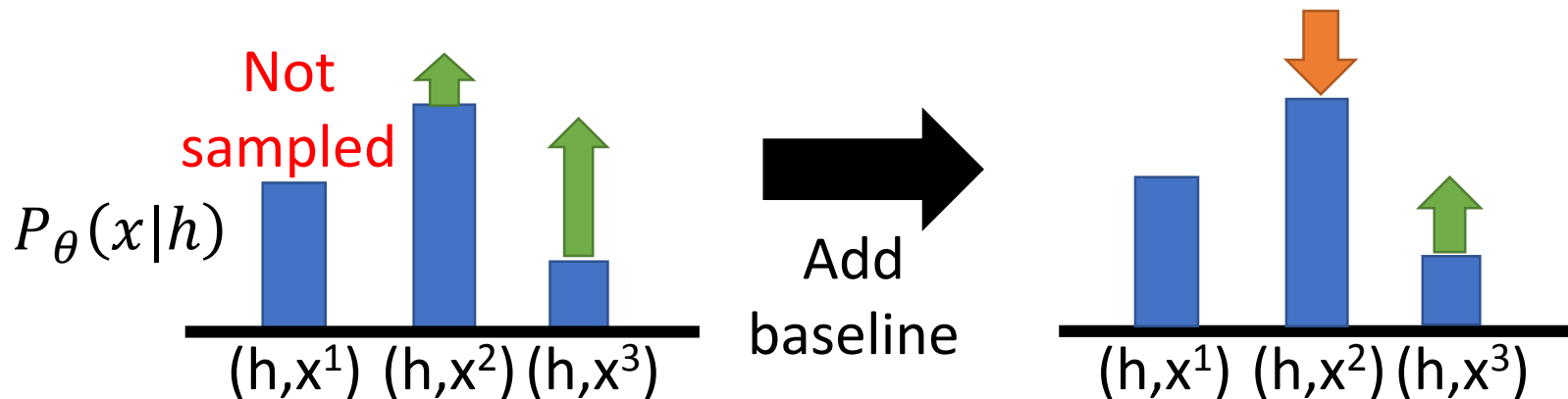
Due to Sampling



Add a Baseline

If $R(h^i, x^i)$ is always positive

$$\frac{1}{N} \sum_{i=1}^N R(h^i, x^i) \log \nabla P_{\theta}(x^i | h^i) \rightarrow \frac{1}{N} \sum_{i=1}^N (R(h^i, x^i) - b) \log \nabla P_{\theta}(x^i | h^i)$$



There are several ways to obtain the baseline b .

Alpha GO style training !

- Let two agents talk to each other



How old are you?



See you.



How old are you?



I am 16.



See you.



See you.



I thought you were 12.



What make you think so?

Using a pre-defined evaluation function to compute $R(h,x)$

Example Reward

- The final reward $R(h,x)$ is the weighted sum of three terms $r_1(h,x)$, $r_2(h,x)$ and $r_3(h,x)$

$$R(h, x) = \lambda_1 \underline{r_1(h, x)} + \lambda_2 \underline{r_2(h, x)} + \lambda_3 \underline{r_3(h, x)}$$

Ease of
answering



不要成為
句點王

Information
Flow



說點
新鮮的

Semantic
Coherence



不要前言
不對後語

Example Results

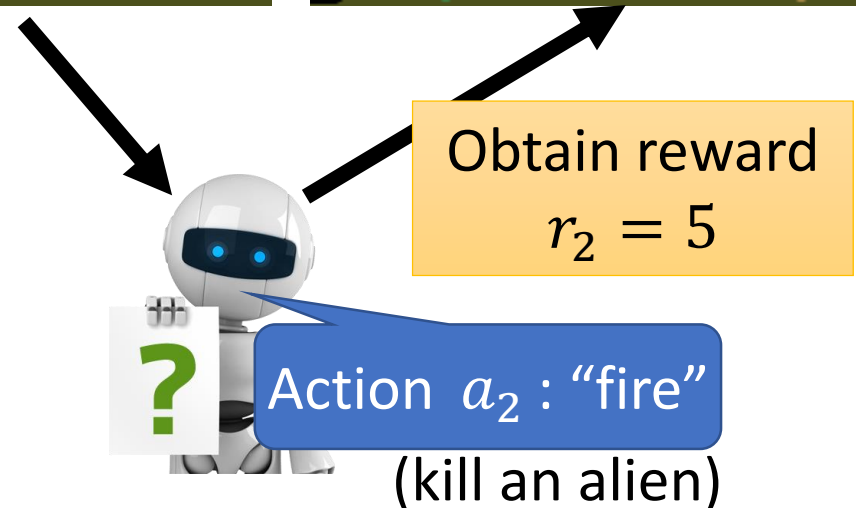
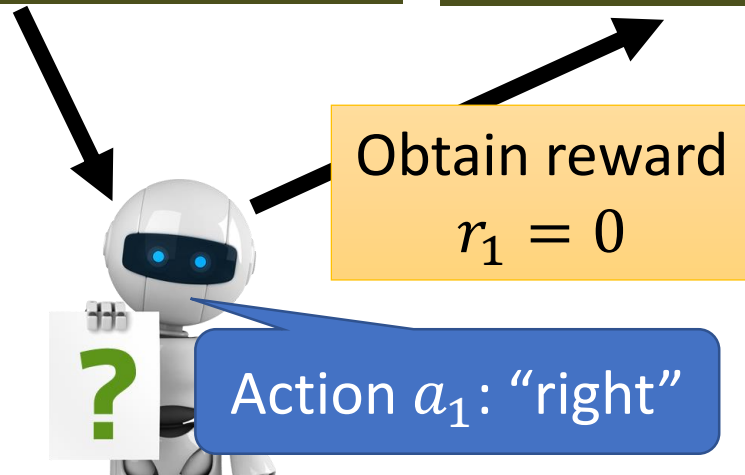
Baseline mutual information model (Li et al. 2015)	Proposed reinforcement learning model
...	...
...	...

Reinforcement learning?

Start with
observation s_1

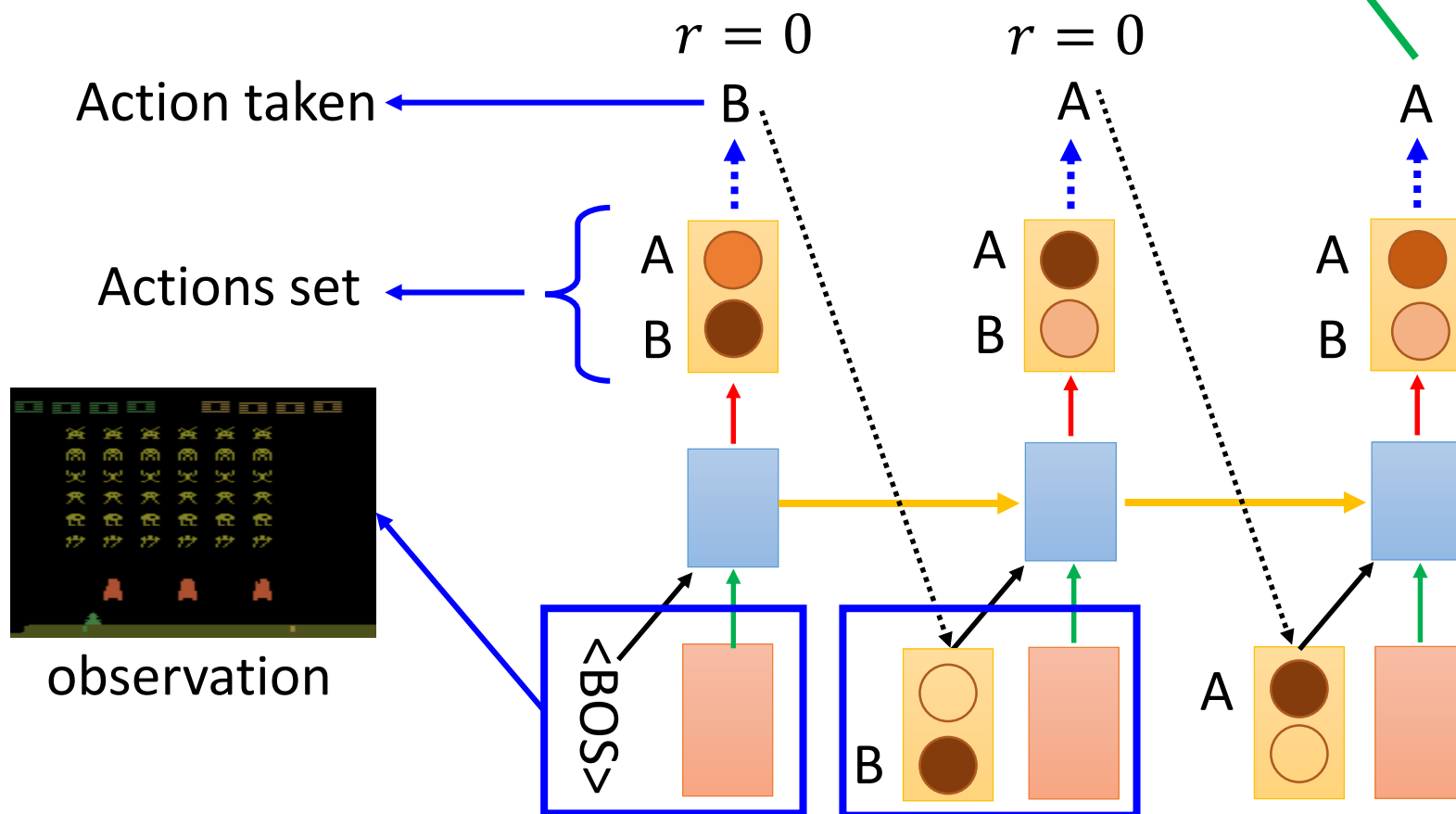
Observation s_2

Observation s_3



Reinforcement learning?

reward:
 $R(\text{"BAA", reference})$



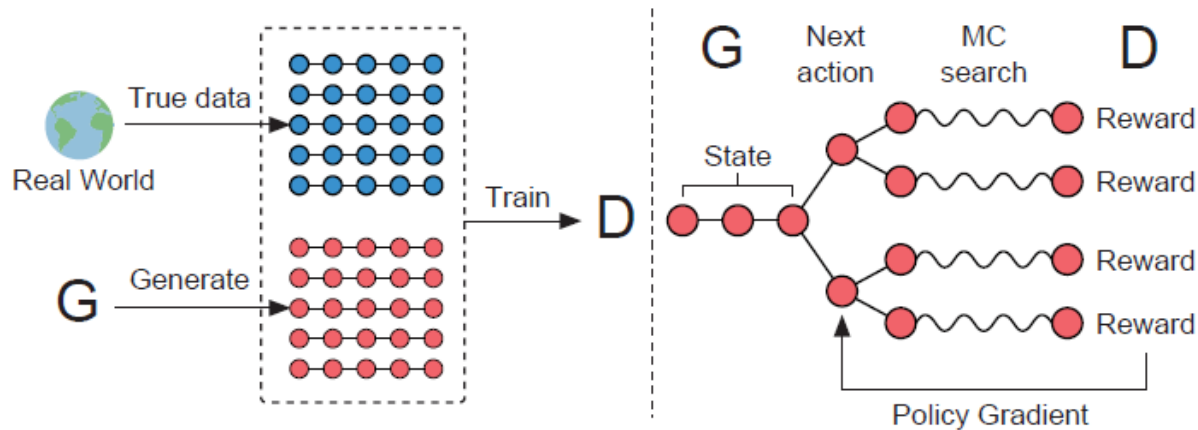
Marc'Aurelio Ranzato, Sumit Chopra, Michael Auli, Wojciech Zaremba, "Sequence Level Training with Recurrent Neural Networks", ICLR, 2016

The action we take influence the observation in the next step

Reinforcement learning?

- One can use any advanced RL techniques here.
- For example, actor-critic
 - Dzmitry Bahdanau, Philemon Brakel, Kelvin Xu, Anirudh Goyal, Ryan Lowe, Joelle Pineau, Aaron Courville, Yoshua Bengio. "An Actor-Critic Algorithm for Sequence Prediction." ICLR, 2017.

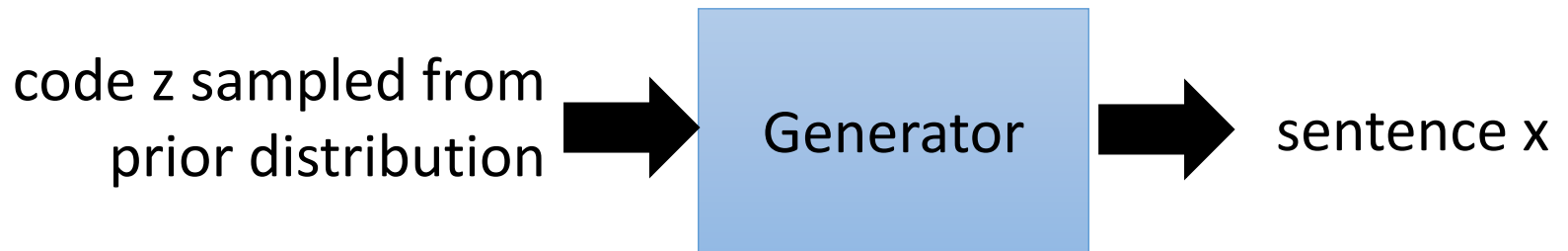
SeqGAN



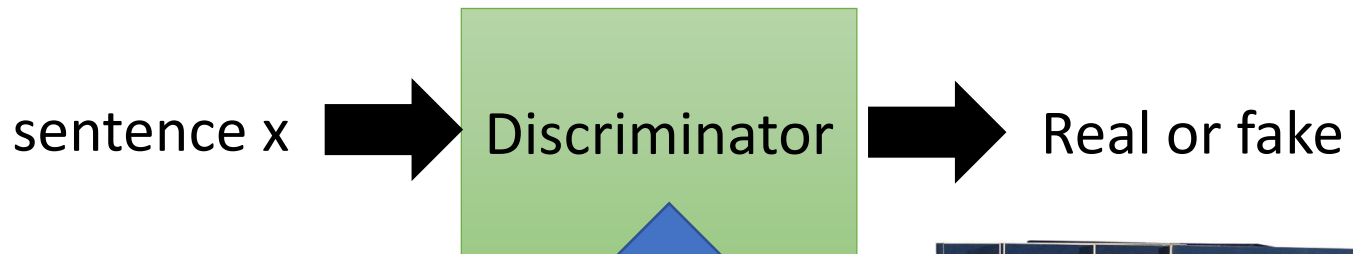
Lantao Yu, Weinan Zhang, Jun Wang, Yong Yu, "SeqGAN: Sequence Generative Adversarial Nets with Policy Gradient", AAAI, 2017

Jiwei Li, Will Monroe, Tianlin Shi, Sébastien Jean, Alan Ritter, Dan Jurafsky, "Adversarial Learning for Neural Dialogue Generation", arXiv preprint, 2017

Basic Idea – Sentence Generation



Sampling from RNN at each time step also provides randomness



Original GAN



Algorithm – Sentence Generation

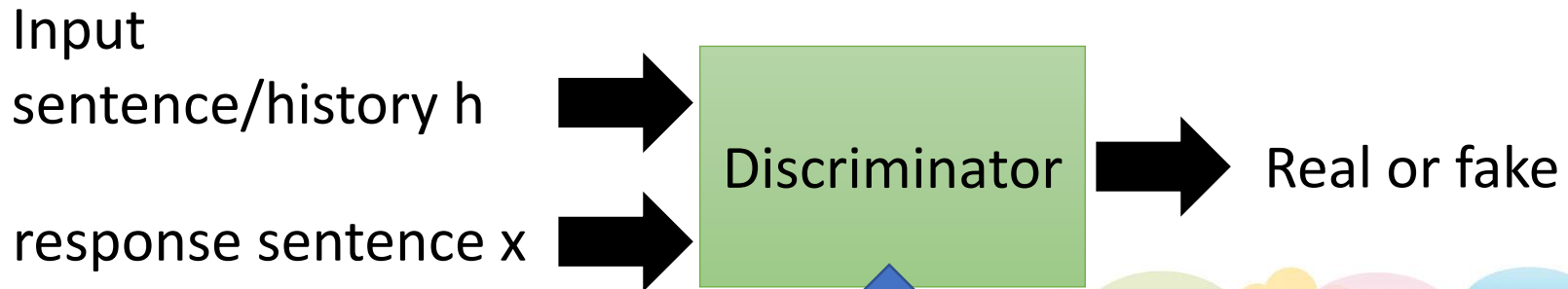
- Initialize generator Gen and discriminator Dis
- In each iteration:

- Sample real sentences x from database
- Generate sentences \tilde{x} by Gen
- Update Dis to increase $Dis(x)$ and decrease $Dis(\tilde{x})$

- Update Gen such that



Basic Idea – Chat-bot



Conditional GAN

human
dialogues



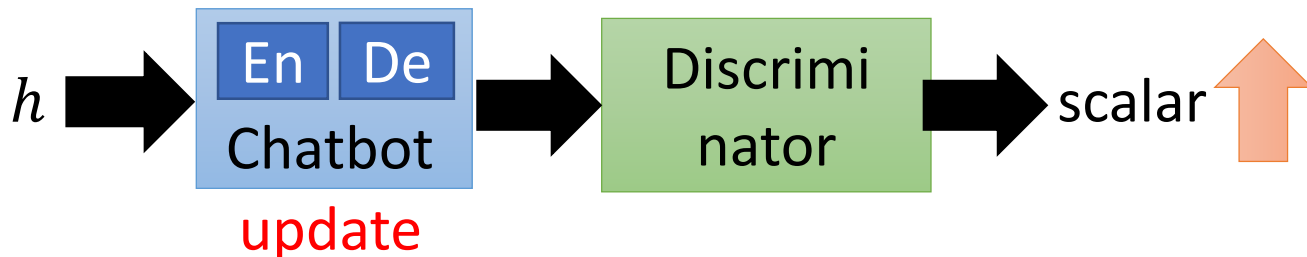
Training data: \vdots

Algorithm – Chat-bot

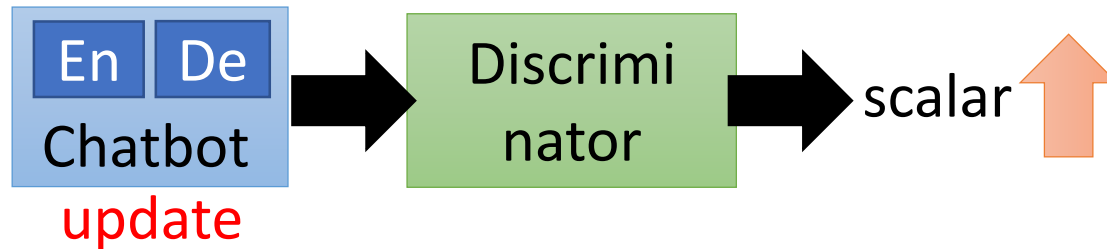
h	A: 000
	B: XXX
x	A: $\Delta \Delta \Delta$
\vdots	

- Initialize generator Gen and discriminator Dis
- In each iteration:
 - Sample real history h and sentence x from database
 - Sample real history h' from database, and generate sentences \tilde{x} by $\text{Gen}(h')$
 - Update Dis to increase $\text{Dis}(h, x)$ and decrease $\text{Dis}(h', \tilde{x})$

- Update Gen such that



Reinforcement Learning



- Consider the output of discriminator as **reward**
 - Update generator to increase discriminator = to get maximum reward

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{i=1}^N \text{reward} \quad D(h^i, x^i) - b \quad \nabla \log P_\theta(x^i | h^i)$$

Discriminator Score

- Different from typical RL
 - The discriminator would update



Reward for Every Generation Step


$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{i=1}^N (D(h^i, x^i) - b) \nabla \log P_\theta(x^i | h^i)$$


h^i = "What is your name?" $D(h^i, x^i) - b$ is negative

x^i = "I don't know" Update θ to decrease $\log P_\theta(x^i | h^i)$

$$\log P_\theta(x^i | h^i) = \log P(x_1^i | h^i) + \log P(x_2^i | h^i, x_1^i) + \log P(x_3^i | h^i, x_{1:2}^i)$$

$P("I" | h^i)$  


 


 


h^i = "What is your name?" $D(h^i, x^i) - b$ is positive

x^i = "I am John" Update θ to increase $\log P_\theta(x^i | h^i)$

$$\log P_\theta(x^i | h^i) = \log P(x_1^i | h^i) + \log P(x_2^i | h^i, x_1^i) + \log P(x_3^i | h^i, x_{1:2}^i)$$

$P("I" | h^i)$ 

Reward for Every Generation Step

h^i = "What is your name?" x^i = "I don't know"

$$\log P_{\theta}(x^i|h^i) = \log P(\underline{x_1^i|h^i}) + \log P(\underline{x_2^i|h^i, x_1^i}) + \log P(\underline{x_3^i|h^i, x_{1:2}^i})$$

$$P(\text{"I"}|h^i) \quad P(\text{"don't"}|h^i, \text{"I"}) \quad P(\text{"know"}|h^i, \text{"I don't"})$$



$$\nabla \bar{R}_{\theta} \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T (Q(h^i, x_{1:t}^i) - b) \nabla \log P_{\theta}(x_t^i|h^i, x_{1:t-1}^i)$$

Method 1. Monte Carlo (MC) Search

Method 2. Discriminator For Partially Decoded Sequences

Monte Carlo Search

- How to estimate $Q(h^i, x_{1:t}^i)$?

$$Q(\underset{h^i}{\text{"What is your name?"}}, \underset{x_1^i}{\text{"I"}})$$

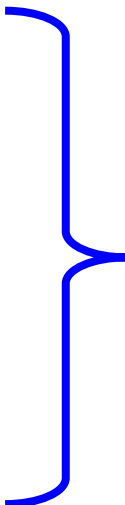
A roll-out generator
for sampling is needed

$$x^A = \text{I am John} \quad D(h^i, x^A) = 1.0$$

$$x^B = \text{I am happy} \quad D(h^i, x^B) = 0.1$$

$$x^C = \text{I don't know} \quad D(h^i, x^C) = 0.1$$

$$x^D = \text{I am superman} \quad D(h^i, x^D) = 0.8$$


$$Q(h^i, \text{"I"}) = 0.5$$

avg

Rewarding Partially Decoded Sequences

- Training a discriminator that is able to assign rewards to both fully and partially decoded sequences
 - Break generated sequences into partial sequences



h="What is your name?", x="I am john"

h="What is your name?", x="I am"

h="What is your name?", x="I"



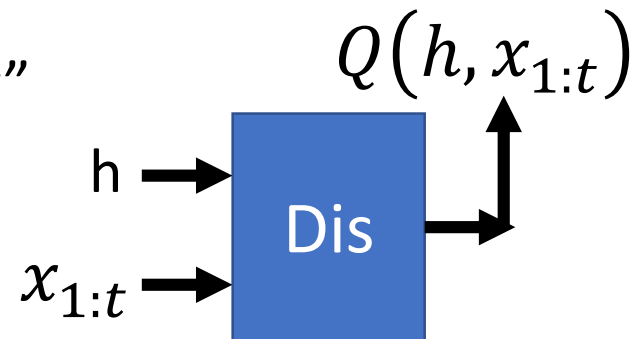
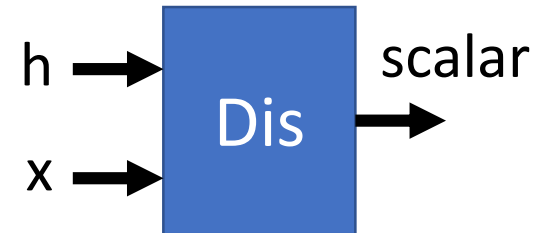
h="What is your name?", x="I don't know"



h="What is your name?", x="I don't"



h="What is your name?", x="I"



Teacher Forcing

- The training of generative model is unstable
 - This reward is used to promote or discourage the generator's own generated sequences.
 - Usually It knows that the generated results are bad, but does not know what results are good.

- Teacher Forcing

Training Data for SeqGAN: $\{(h^1, x^1), \dots, (h^N, x^N)\}$
Obtained by sampling
weighted by $D(h^i, x^i)$

Adding more Data: $\{(h^1, \hat{x}^1), \dots, (h^N, \hat{x}^N)\}$ Real data
Consider $D(h^i, \hat{x}^i) = 1$

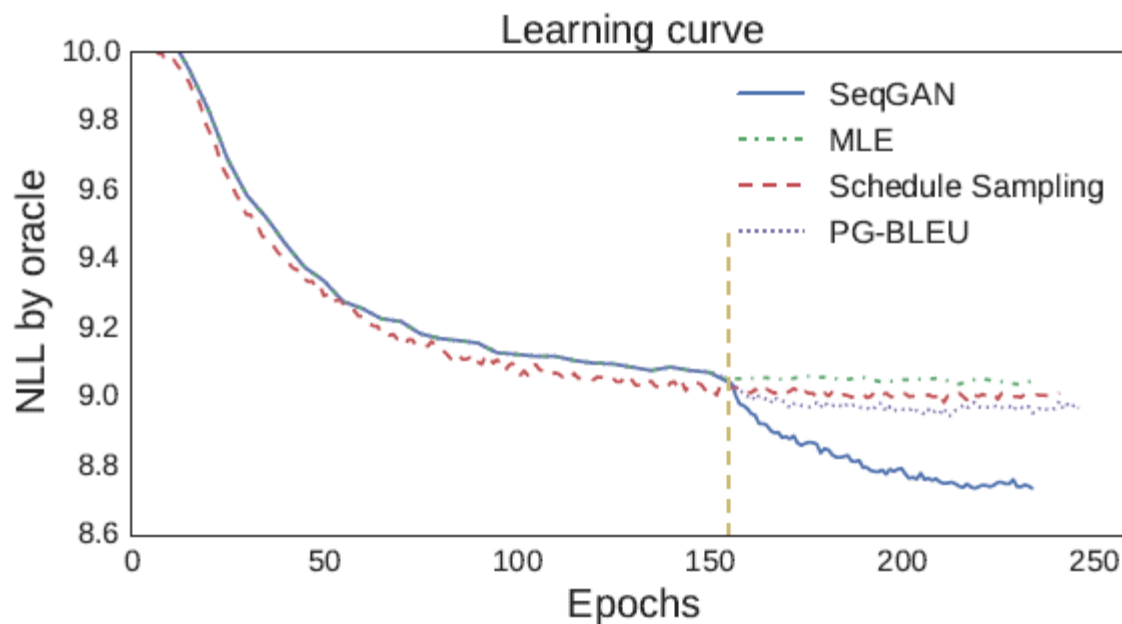
Experiments in paper

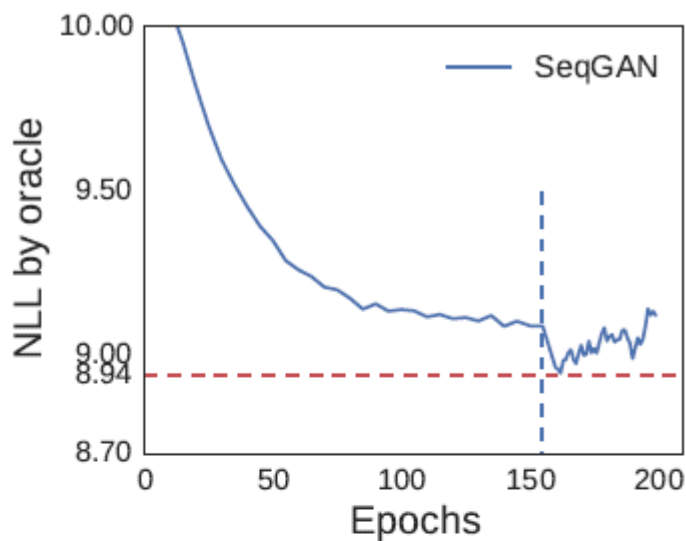
- Sentence generation: Synthetic data
 - Given an LSTM
 - Using the LSTM to generate a lot of sequences as “real data”
 - Generator learns from the “real data” by different approaches
 - Generator generates some sequences
 - Using the LSTM to compute the negative log likelihood (NLL) of the sequences
 - Smaller is better

Experiments in paper

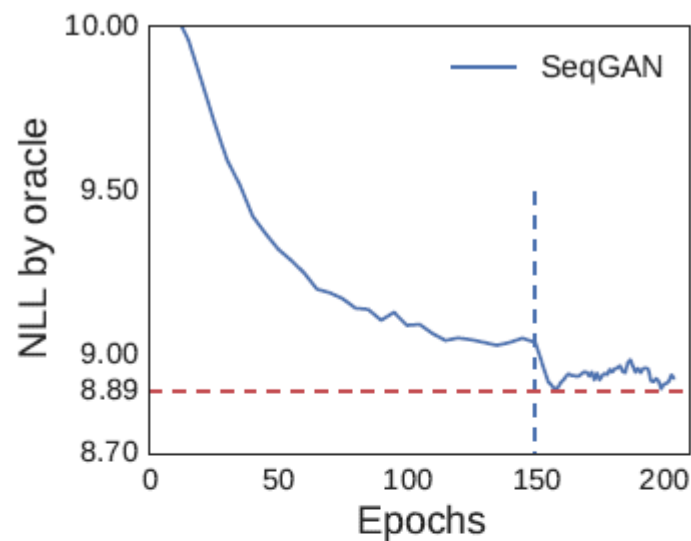
- Synthetic data

Algorithm	Random	MLE	SS	PG-BLEU	SeqGAN
NLL	10.310	9.038	8.985	8.946	8.736
<i>p</i> -value	$< 10^{-6}$	$< 10^{-6}$	$< 10^{-6}$	$< 10^{-6}$	

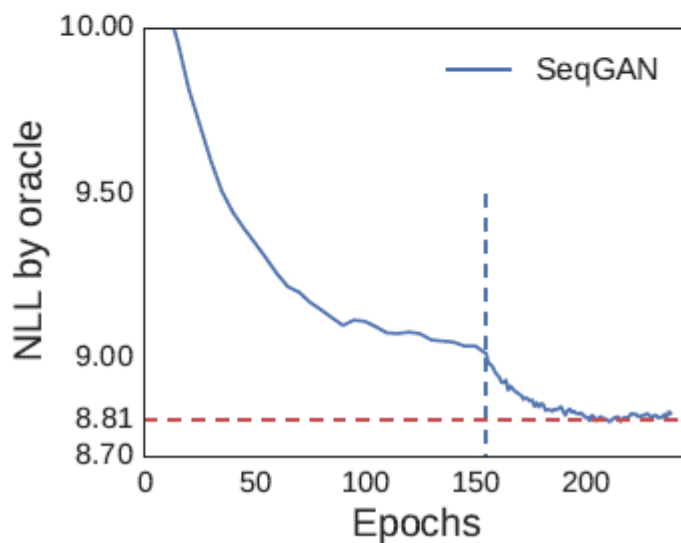




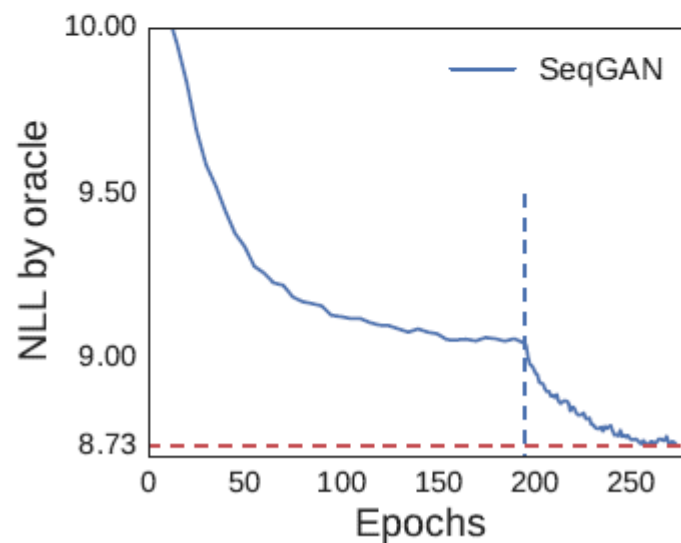
(a) $g\text{-steps}=100$, $d\text{-steps}=1$, $k=10$



(b) $g\text{-steps}=30$, $d\text{-steps}=1$, $k=30$



(c) $g\text{-steps}=1$, $d\text{-steps}=1$, $k=10$



(d) $g\text{-steps}=1$, $d\text{-steps}=5$, $k=3$

Experiments in paper

- Real data

Table 2: Chinese poem generation performance comparison.

Algorithm	Human score	p -value	BLEU-2	p -value
MLE	0.4165	0.0034	0.6670	$< 10^{-6}$
SeqGAN	0.5356		0.7389	
Real data	0.6011		0.746	

Table 3: Obama political speech generation performance.

Algorithm	BLEU-3	p -value	BLEU-4	p -value
MLE	0.519	$< 10^{-6}$	0.416	0.00014
SeqGAN	0.556		0.427	

Table 4: Music generation performance comparison.

Algorithm	BLEU-4	p -value	MSE	p -value
MLE	0.9210	$< 10^{-6}$	22.38	0.00034
SeqGAN	0.9406		20.62	

Results - Chat-bot

Setting	adver-win	adver-lose
single-turn	0.62	0.18
multi-turn	0.72	0.10

Input	tell me ... how long have you had this falling sickness ?
Vanilla-MLE	i 'm not a doctor .
Vanilla-Sample	well everything you did was totally untrue .
REINFORCE	i don 't know how long it 's been .
REGS Monte Carlo	A few months, I guess .
Input	so i had the doctors test sammy 's response to conditioning .
Vanilla-MLE	sammy wrote the test sammy wrote the test .
Vanilla-Sample	objects pick objects up objects objects objects objects
REINFORCE	i 'm not sure that 's relevant .
REGS Monte Carlo	so he took the pills .
Input	what are you going to do ?
Vanilla-MLE	i 'm going to the movies .
Vanilla-Sample	get him outta here first !
REINFORCE	i 'm going to get you .
REGS Monte Carlo	i 'm going to get you a new car .
Input	they fear your power your intellect .
Vanilla-MLE	you 're the only one who knows what 's going on .
Vanilla-Sample	when they are conquered and you surrender they will control all of us .
REINFORCE	i 'm afraid i 'm not ready yet .
REGS Monte Carlo	i 'm not afraid of your power .

To Learn More ...

Algorithm – MaliGAN

Maximum-likelihood Augmented Discrete GAN

- Initialize generator Gen and discriminator Dis
- In each iteration:

- Sample real sentences x from database
- Generate sentences \tilde{x} by Gen
- Update Dis to maximize

$$\sum_x \log D(x) + \sum_{\tilde{x}} \log(1 - D(\tilde{x}))$$

- Update Gen by gradient

$$\frac{1}{N} \sum_{i=1}^N \left(\frac{r_D(x^i)}{\sum_{i=1}^N r_D(x^i)} - b \right) \nabla \log P_{\theta}(x^i)$$

$$D(h^i, x^i)$$

$$r_D(x^i) = \frac{D(x^i)}{1 - D(x^i)}$$

To learn more

- Professor forcing
 - Alex Lamb, Anirudh Goyal, Ying Zhang, Saizheng Zhang, Aaron Courville, Yoshua Bengio, “Professor Forcing: A New Algorithm for Training Recurrent Networks”, NIPS, 2016
- Handling discrete output by methods other than policy gradient
 - MaliGAN, Boundary-seeking GAN
 - Yizhe Zhang, Zhe Gan, Lawrence Carin, “Generating Text via Adversarial Training”, Workshop on Adversarial Training, NIPS, 2016
 - Matt J. Kusner, José Miguel Hernández-Lobato, “GANS for Sequences of Discrete Elements with the Gumbel-softmax Distribution”, arXiv preprint, 2016