*Applied Deep Learning*

# Natural Language Generation

**May 12th, 2020**   **http://adl.miulab.tw**

國立臺灣大學
National Taiwan University

# Outline

- ◉ NLG Review
  - ○ Language Modeling
  - ○ Conditional Language Modeling
- ◉ Decoding Algorithm
  - ○ Greedy
  - ○ Beam Search
  - ○ Sampling
- ◉ Evaluation
- ◉ Reinforcement Learning for NLG
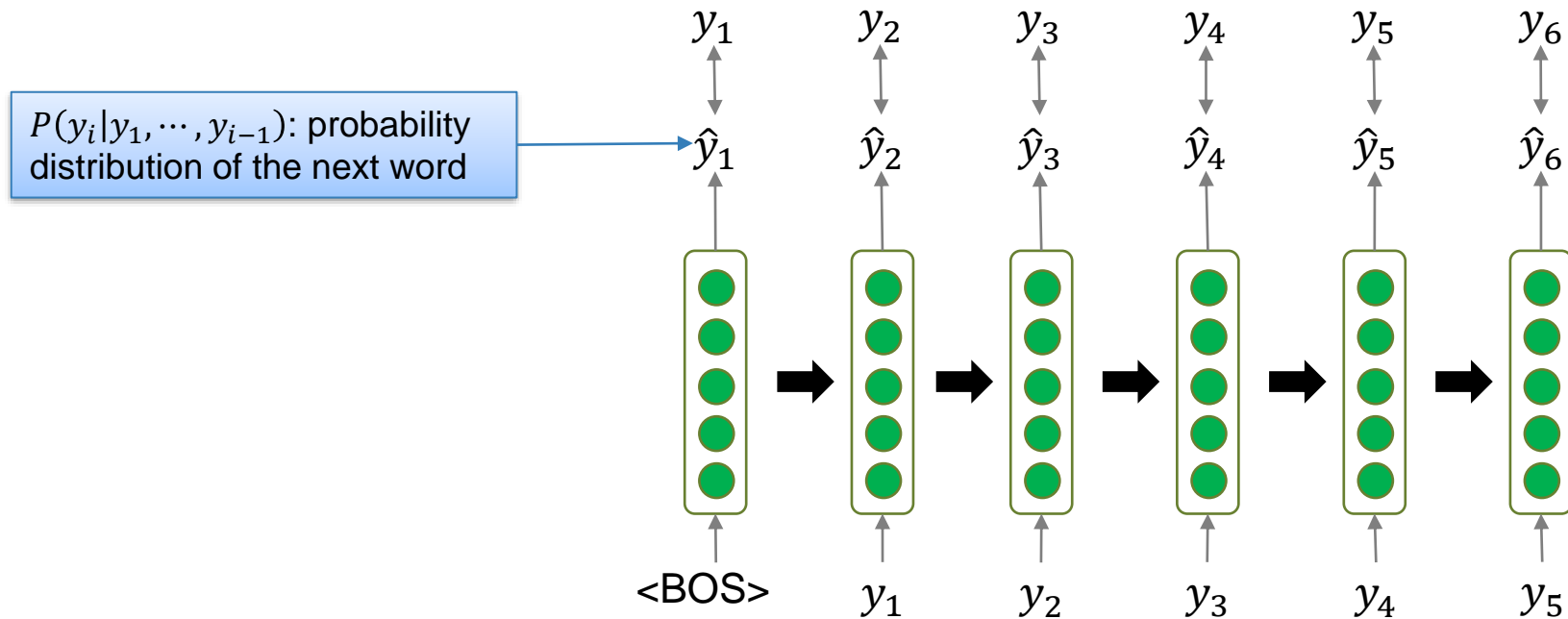
# Natural Language Generation

◉ Many tasks contain NLG
  - ○ Machine Translation
  - ○ Abstractive Summarization
  - ○ Dialogue Generation
  - ○ Image Captioning
  - ○ Creative Writing
    - ■ Storytelling, poetry generation
  - ○ …

# **Language Modeling**

◉ Goal: predicting the next word given the words so far

$$P(y_i|y_1, \cdots, y_{i-1})$$

◉ **Language model** is to estimate the probability distribution
   ○ RNN-LM is to use RNN for modeling the distribution

# RNN-LM

$$P(y_i|y_1,\cdots,y_{i-1}): \text{probability distribution of the next word}$$

$$y_1 \quad y_2 \quad y_3 \quad y_4 \quad y_5 \quad y_6$$

$$\hat{y}_1 \quad \hat{y}_2 \quad \hat{y}_3 \quad \hat{y}_4 \quad \hat{y}_5 \quad \hat{y}_6$$

<BOS>  $y_1 \quad y_2 \quad y_3 \quad y_4 \quad y_5$

Idea: pass the information from the previous hidden layer to leverage all contexts

# Conditional Language Modeling

- Goal: predicting the next word given the words so far, and other input x

$$P(y_i|y_1, \cdots, y_{i-1}, x)$$

- Conditional language modeling tasks
  - Machine translation (x = source sentence, y = target sentence)
  - Summarization (x = document, y = summary)
  - Dialogue (x = dialogue context, y = response)
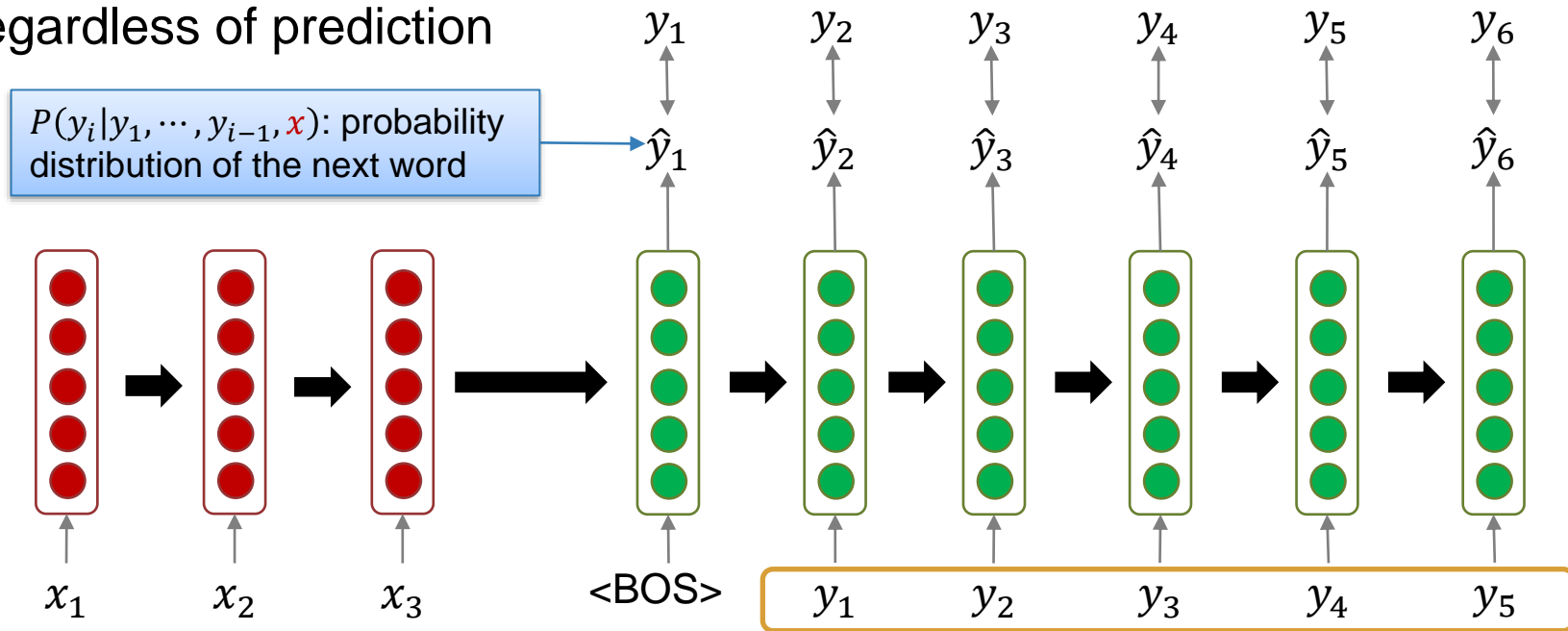  - Image captioning (x = image, y = caption)
  - …

# Sequence-to-Sequence Modeling

$P(y_i|y_1, \cdots, y_{i-1}, x)$: probability distribution of the next word



Training an encoder-decoder model that generate the next word with condition

# Teacher Forcing

○ During training, feeding the gold target sentence into the decoder regardless of prediction

$P(y_i|y_1, \cdots, y_{i-1}, x)$: probability distribution of the next word

$y_1$    $y_2$    $y_3$    $y_4$    $y_5$    $y_6$

$\hat{y}_1$    $\hat{y}_2$    $\hat{y}_3$    $\hat{y}_4$    $\hat{y}_5$    $\hat{y}_6$

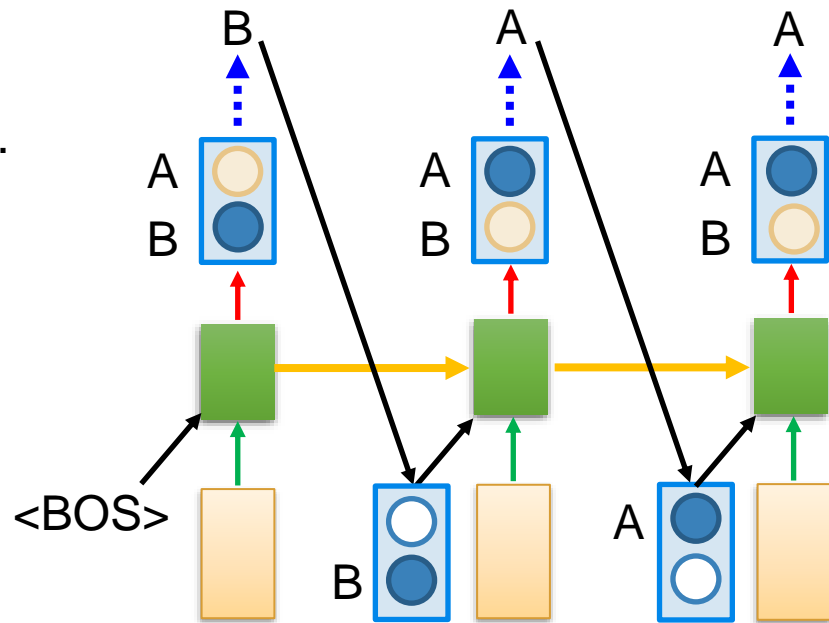$x_1$    $x_2$    $x_3$    <BOS>    $y_1$    $y_2$    $y_3$    $y_4$    $y_5$

Issue: mismatch between training and testing

# Mismatch between Train and Test

○ **_Training_**

$$C = \sum_t C_t$$

minimizing cross-entropy of each word

: condition

Reference:
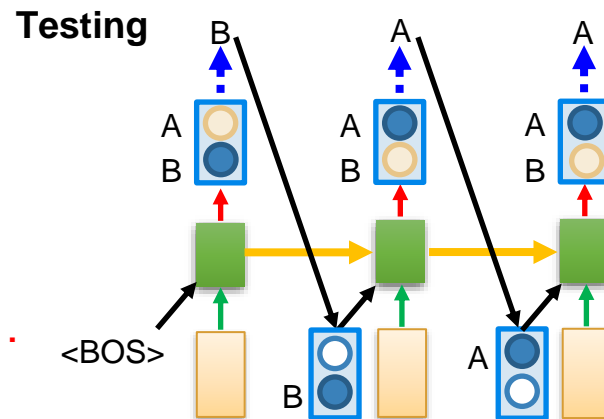
# **Mismatch between Train and Test**

## ◉ ***Generation***
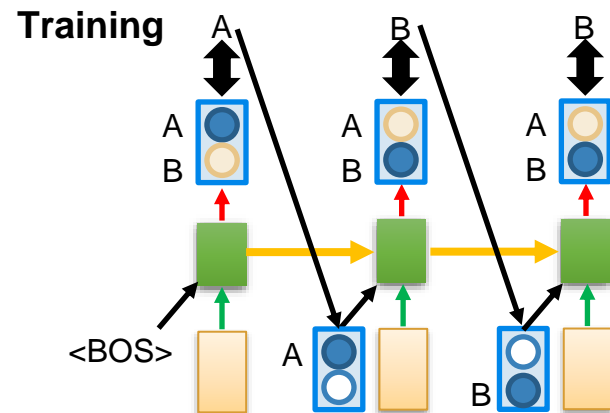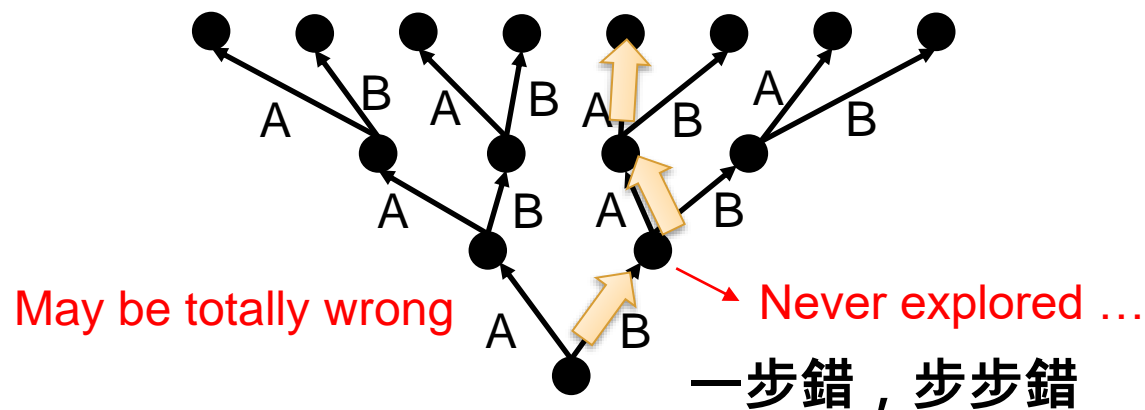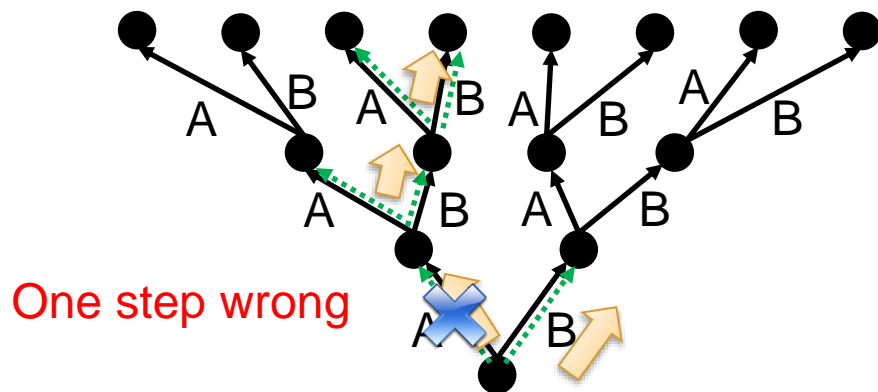
- ○ Testing: Output of model is the input of the next step.
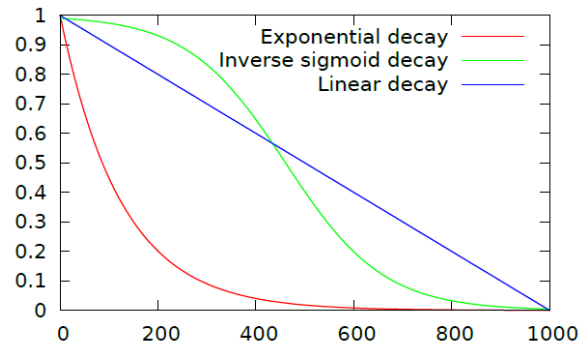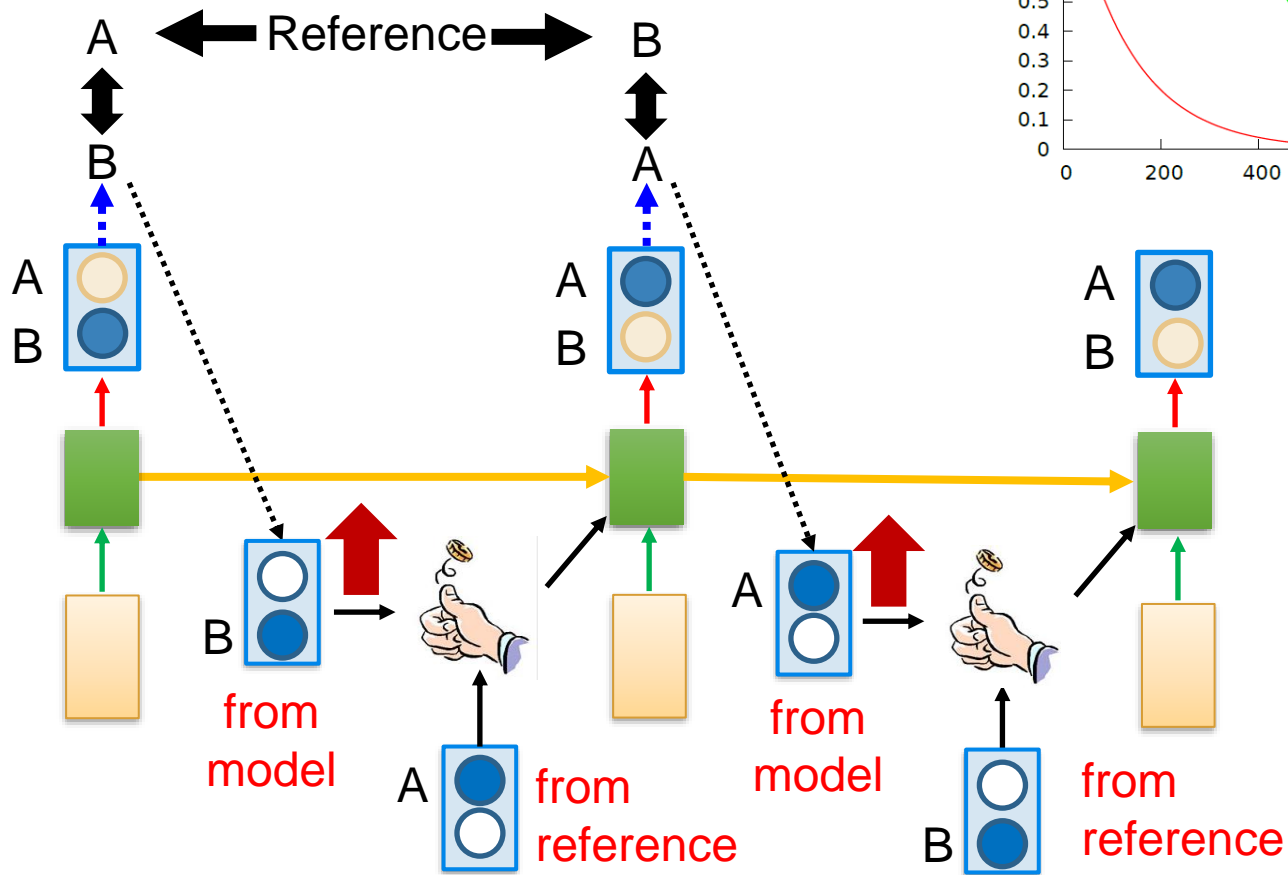  - ■ Reference is unknown
- ○ Training: the inputs are reference.

Exposure Bias

# Exposure Bias



One step wrong

May be totally wrong

一步錯，步步錯

Training

Testing

Never explored …

<BOS>

# Scheduled Sampling

# Scheduled Sampling

◉ Image captioning on MSCOCO

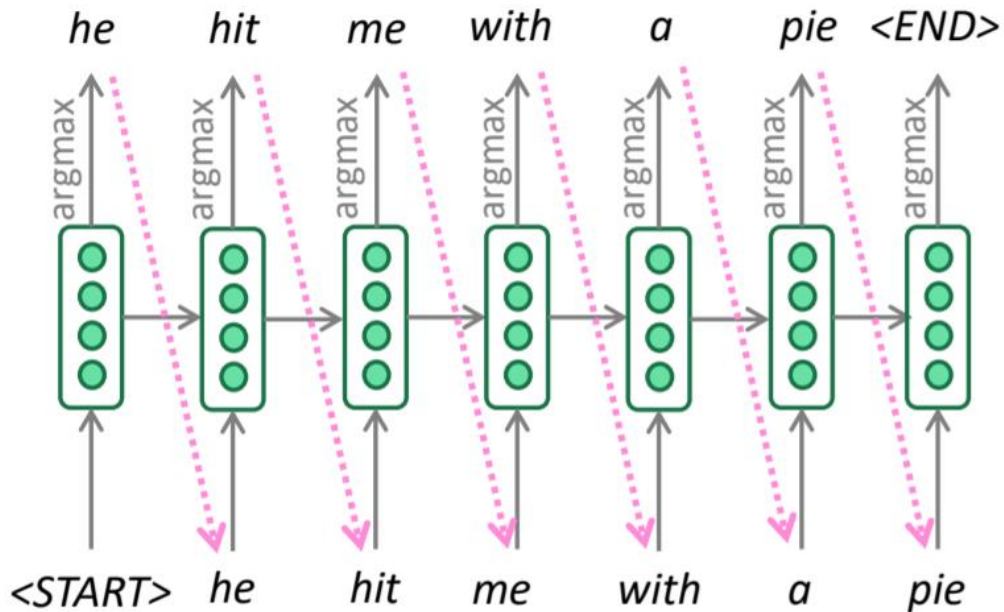|  | BLEU-4 | METEOR | CIDER |
|---|---|---|---|
| Always from reference | 28.8 | 24.2 | 89.5 |
| Always from model | 11.2 | 15.7 | 49.7 |
| Scheduled Sampling | 30.6 | 24.3 | 92.1 |

# 14 Decoding Algorithm

Strategy of Word Generation

# Decoding Algorithm

- With a trained (conditional) LM, a _decoding algorithm_ decides how to generate texts from the LM.
- Decoding Algorithms
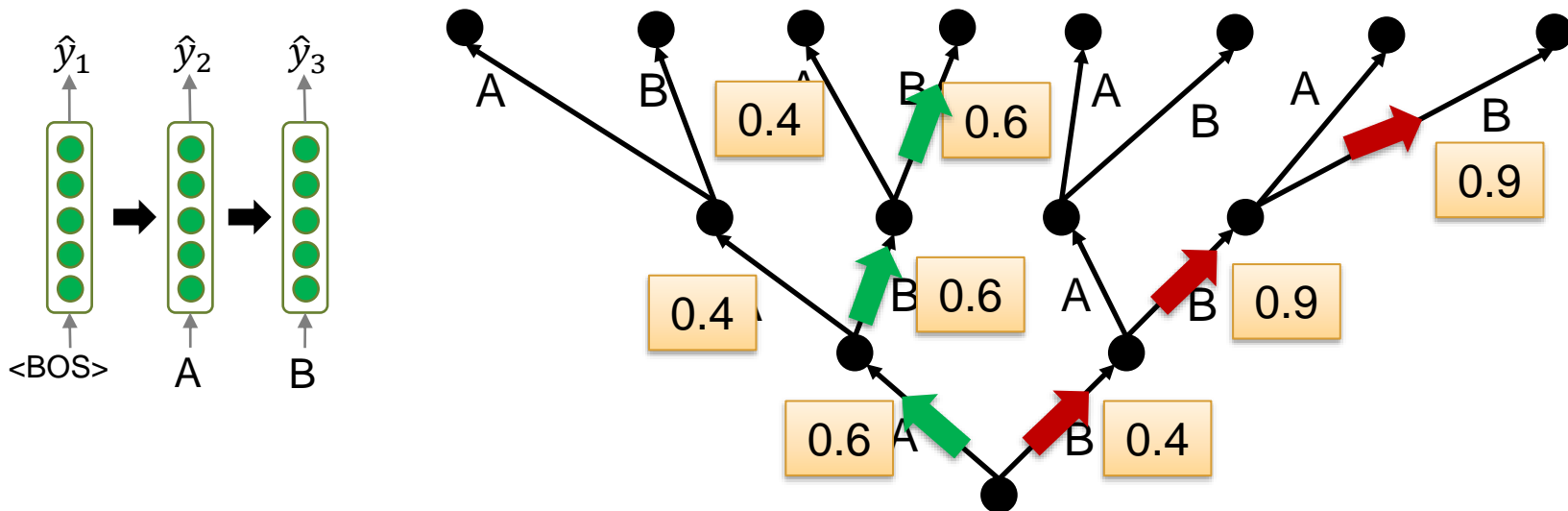  - Greedy
  - Beam Search
  - Sampling

# Greedy

⦿ Strategy: choosing the most probable word (argmax)



Output can be poor due to lack of backtracking

# Suboptimal Issue

◉ Unexplored path may have higher probability.
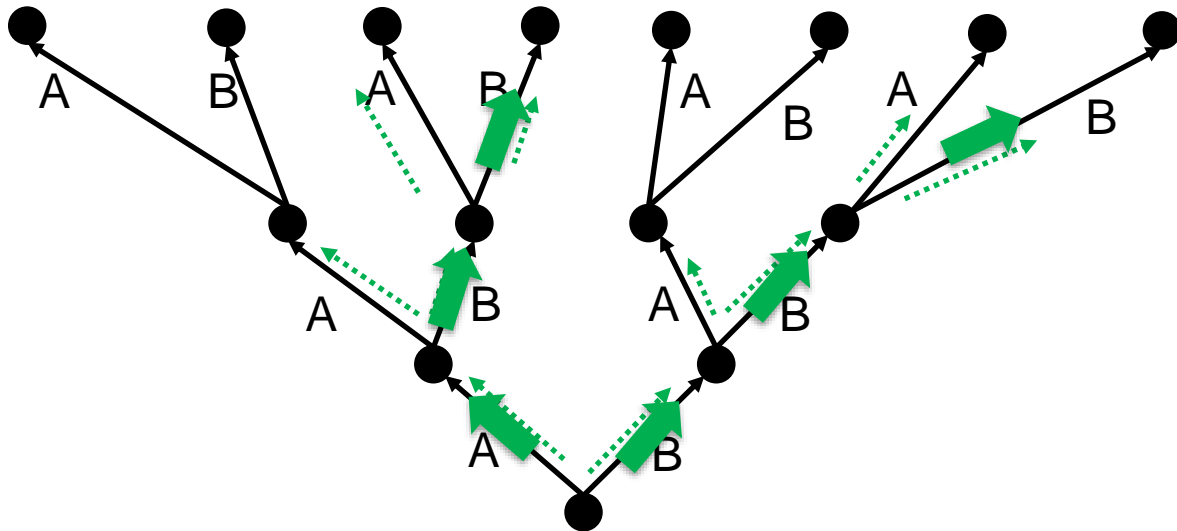


The red path has higher score.

Issue: Impossible to check all paths

# Beam Search

⊙ Strategy: keeping track of the *k* most probable sequences and finding a better one

Keep several best paths at each step (beam size = 2)

# Beam Search



T = 1

current hypotheses   proposed extensions

T = 2

current hypotheses   proposed extensions

T = 3

current hypotheses   proposed extensions

empty string

A standard beam search algorithm with an alphabet of $\{\epsilon, a, b\}$ and a beam size of three.

The size of beam is 3 in this example.

# Effect of Beam Size

- Small $k$
  - Ungrammatical, unnatural, incorrect, etc.

- Large $k$
  - Reduce some above issues
  - Computationally expensive
  - Introduce other issues
    - Chit-chat dialogues with large beam often generate generic sentences

# Effect of Beam Size in Chit-Chat Dialogues

I mostly eat a fresh and raw diet, so I save on groceries

| Beam Size | Model Response |
|---|---|
| 1 | *I love to eat healthy and eat healthy* |
| 2 | *That is a good thing to have* |
| 3 | *I am a nurse so I do not eat raw food* |
| 4 | *I am a nurse so I am a nurse* |
| 5 | *Do you have any hobbies?* |
| 6 | *What do you do for a living?* |
| 7 | *What do you do for a living?* |
| 8 | *What do you do for a living?* |

**Small Beam Size:** More on-topic but nonsensical; bad English

**Large Beam Size:** safe, "correct" response, but generic and less relevant

Finding a proper beam size is not trivial

# Sampling-Based Decoding

◉ Strategy: choosing the next word with randomness (from a distribution)

◉ Sampling

  ○ Randomly sample the word via the <span style="color:red">probability distribution</span> instead of argmax

◉ Top-N Sampling

  ○ Sample the word via distribution but <span style="color:red">restricted to the top-N</span> probable words

  ○ N=1 is greedy, N=V is pure sampling

  ○ Increasing N gets more diverse / risky output

  ○ Decreasing N gets more generic / safe output

Balancing between diversity and safety is an important direction

# Probability Distribution

1. Softmax

$$P(w_t) = \frac{e^{s_w}}{\sum_{w' \in V} e^{s_{w'}}}$$

softmax: LM computes a prob dist by applying softmax to a vector of scores

2. Softmax temperature: applying a temperature hyperparameter $\tau$ to the softmax

$$P(w_t) = \frac{e^{s_w/\tau}}{\sum_{w' \in V} e^{s_{w'}/\tau}}$$

- ○ Higher temperature: $P(w_t)$ becomes more uniform → more diversity
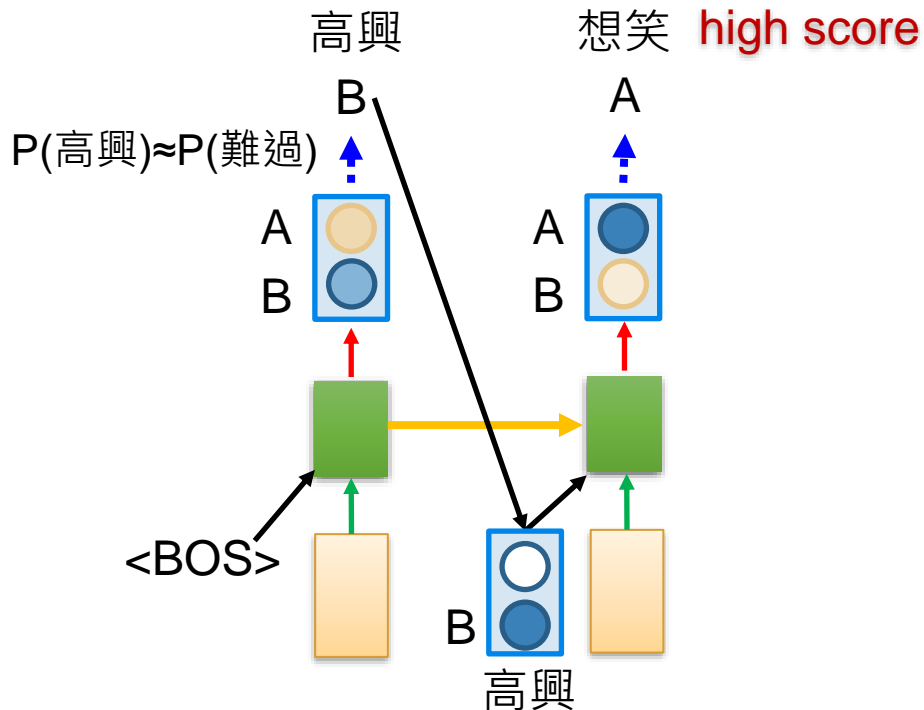- ○ Lower temperature: $P(w_t)$ becomes more spiky → less diversity

softmax temperature is not a decoding algorithm, which is the way of controlling the diversity during testing via any decoding algorithm
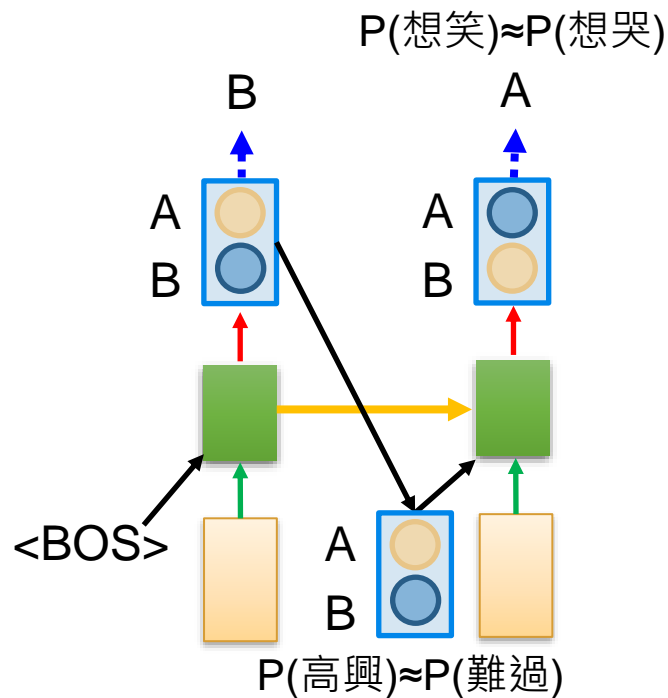
# Distribution Input

U: 你覺得如何?
M: 高興想笑 or 難過想哭

**One-Hot Input**

**Distribution Input**



Distribution input may not be good for NLG

# NLG Evaluation

How Good The Model Performs

# BLEU

◉ N-Gram Precision

$$p_n = \frac{\sum_{ngram \in hyp} count_{clip}(ngram)}{\sum_{ngram \in hyp} count(ngram)}$$

highest count of n-gram in any reference sentence

◉ Brevity Penalty

$$B = \begin{cases} e^{(1-|ref|/|hyp|)}, \text{if } |ref| > |hyp| \\ \quad\quad 1 \quad\quad\quad , \text{otherwise} \end{cases}$$

◉ BLEU
  ○ Often used in machine translation

$$BLEU = B \cdot exp\left[\frac{1}{N}\sum_{n=1}^{N} p_n\right]$$

# ROUGE

◉ ROUGE (Recall-Oriented Understudy for Gisting Evaluation)
  ○ Often used in summarization tasks

ROUGE-N

$$= \frac{\sum_{S\in\{ReferemceSummaries\}} \sum_{gram_n \in S} Count_{match}(gram_n)}{\sum_{S\in\{ReferenceSummaries\}} \sum_{gram_n \in S} Count(gram_n)}$$

# BLEU & ROUGE

- BLEU
  - Based on <u>n-gram overlap</u>
  - Consider precision
  - Reported as a single number
    - Combination of n = 1, 2, 3, 4 n-grams

- ROUGE
  - Based on <u>n-gram overlap</u>
  - Consider recall
  - Reported separately for each n-gram
    - ROUGE-1: unigram overlap
    - ROUGE-2: bigram overlap
    - ROUGE-L: LCS overlap

# Automatic Evaluation Metrics

◉ Word overlap metrics: BLEU, ROUGE, METEOR, etc.
- ○ Not ideal for machine translation
- ○ Much worse for summarization
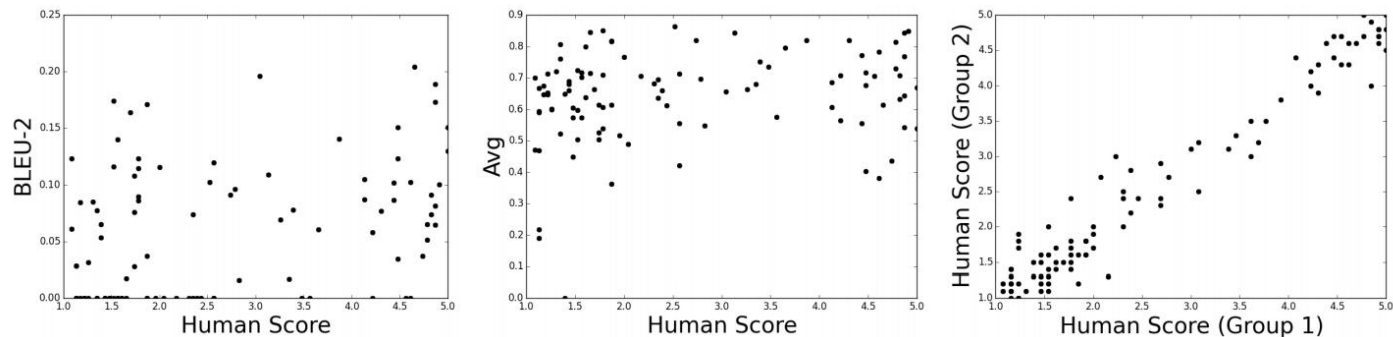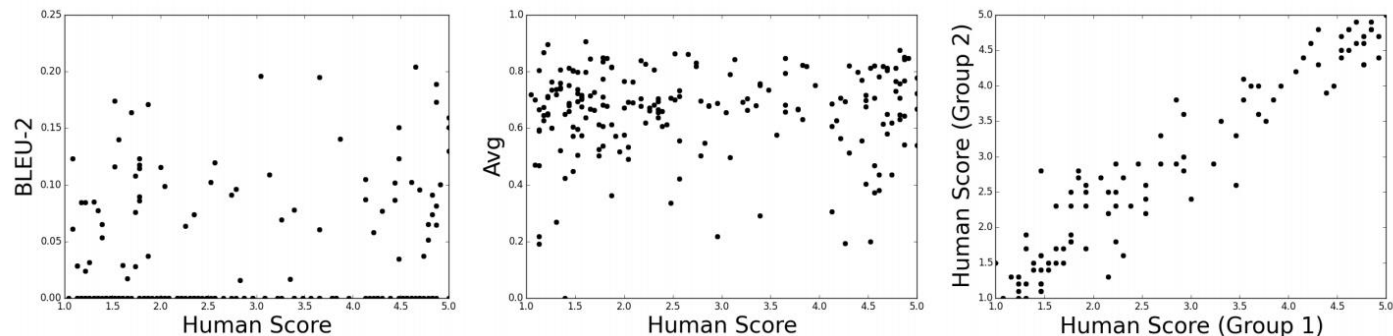- ○ Even worse for dialogue, storytelling

more open-ended

◉ Embedding metrics
- ○ Computing the similarity of word embeddings
- ○ Capturing semantics in a flexible way

# Automatic Metrics v.s. Human Judgement



(a) Twitter

(b) Ubuntu

No agreement between automatic scores and human scores in dialogue quality

# **Focused Metrics for Particular Aspects**

◉ Evaluating a single aspect instead of the overall quality
  - ○ Fluency (compute probability w.r.t. well-trained LM)
  - ○ Correct style (prob w.r.t. LM trained on target corpus)
  - ○ Diversity (rare word usage, uniqueness of n-grams)
  - ○ Relevance to input (semantic similarity measures)
  - ○ Simple things like length and repetition
  - ○ Task-specific metrics e.g. compression rate for summarization

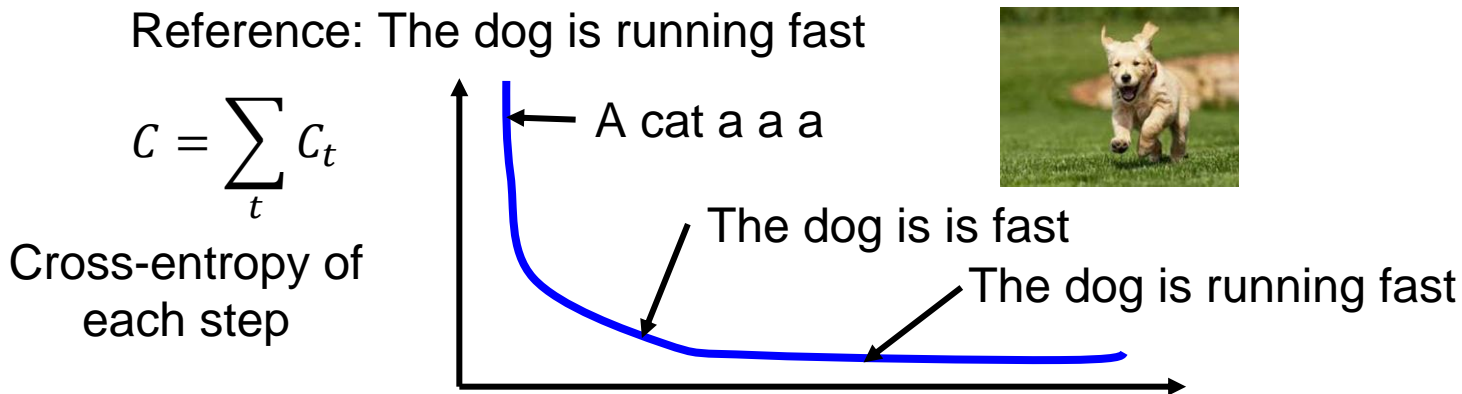Scores help us track some important qualities we care about

# Reinforcement Learning for NLG

Global Optimization
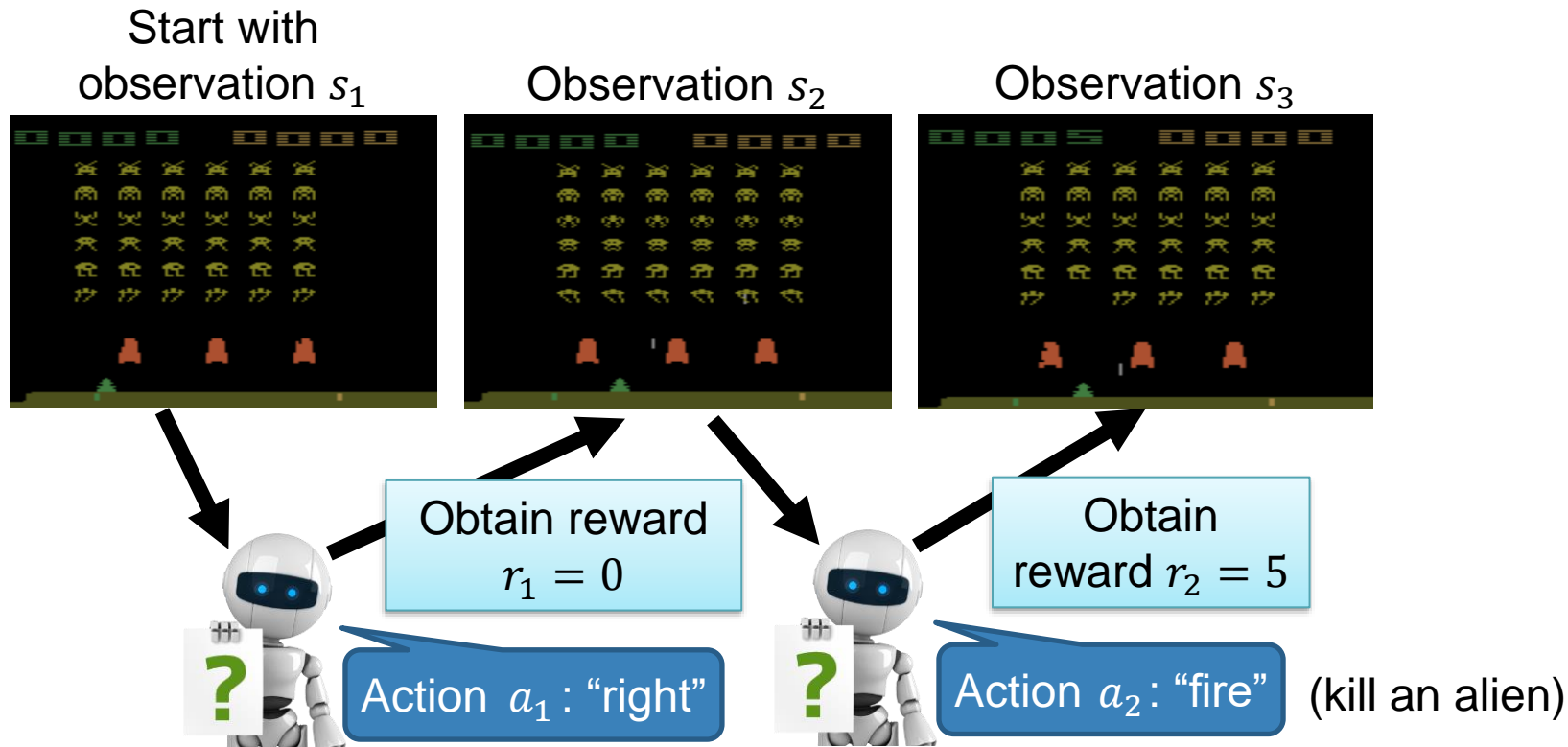
# Global Optimization v.s. Local Optimization

◉ Minimizing the error defined on component level (local) is not equivalent to improving the generated objects (global)

Reference: The dog is running fast

$$C = \sum_t C_t$$

Cross-entropy of each step

A cat a a a

The dog is is fast
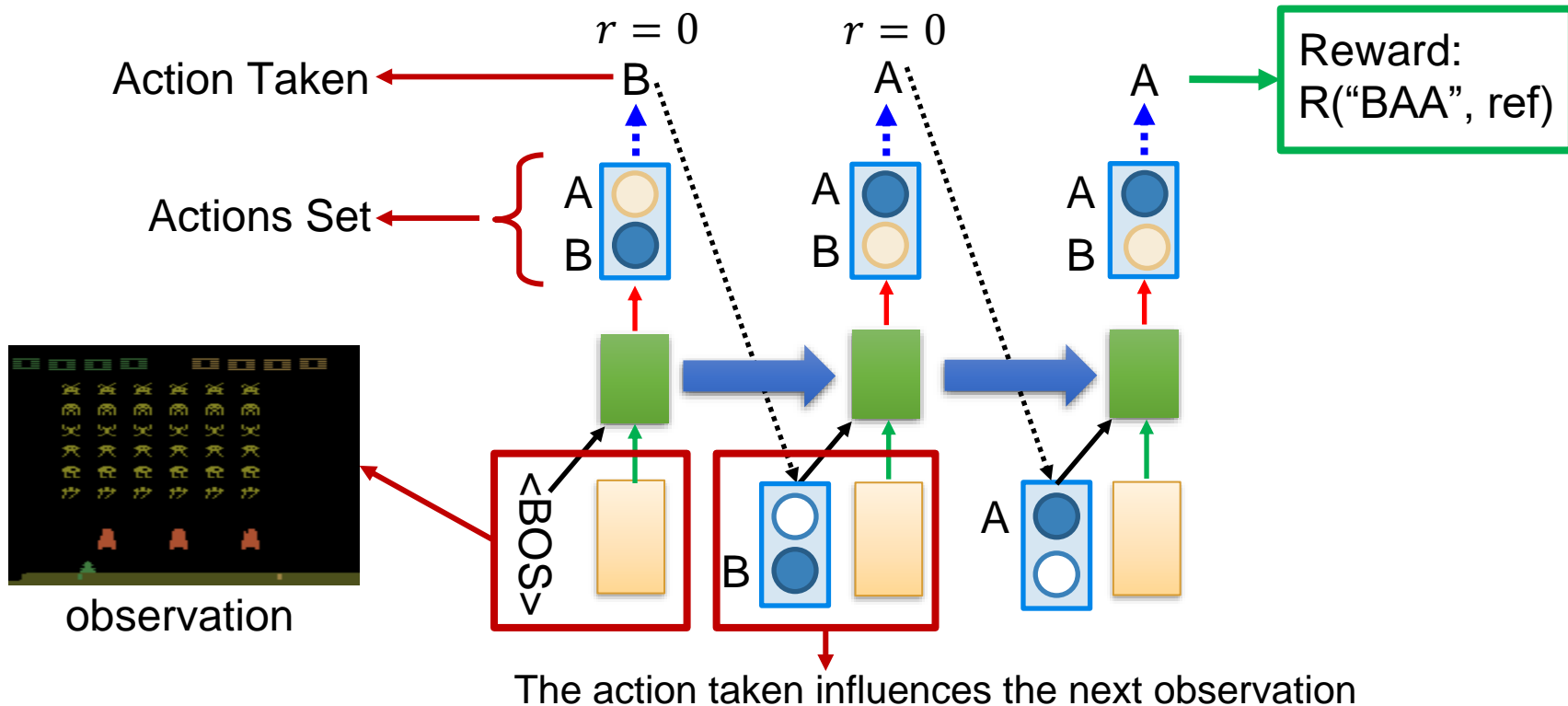
The dog is running fast

Optimize object-level criterion instead of component-level cross-entropy.
Object-level criterion: $R(y, \hat{y})$  $y$: ground truth, $\hat{y}$: generated sentence
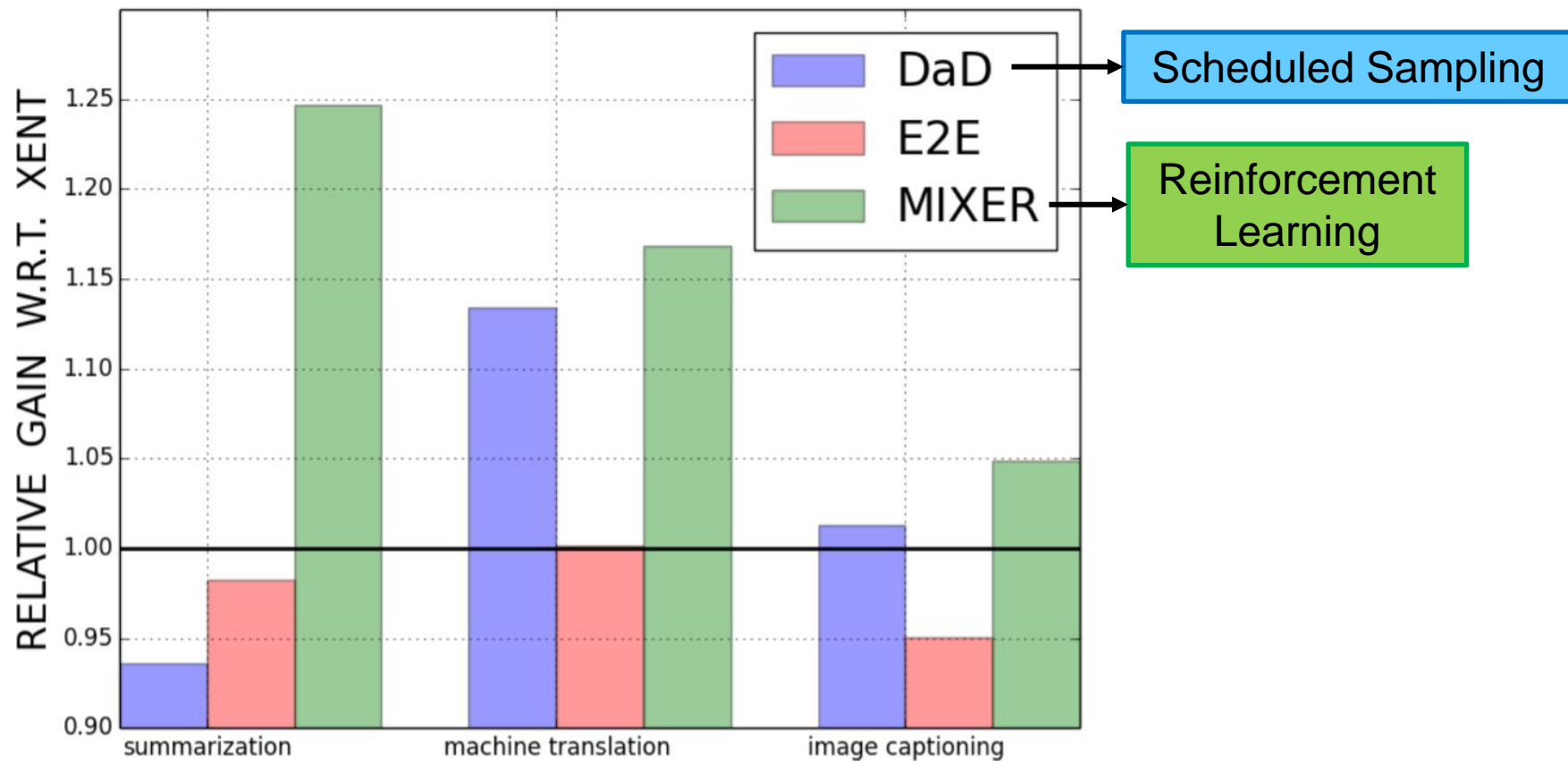
Gradient Descent?
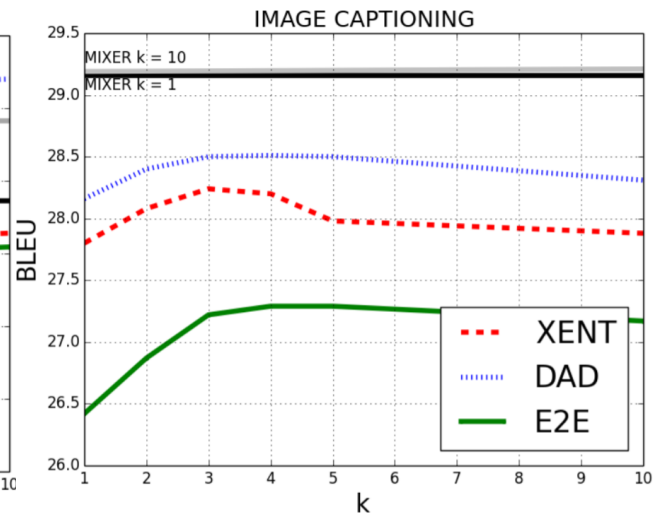
# Reinforcement Learning

Start with
observation $s_1$

Observation $s_2$

Observation $s_3$

Obtain reward
$r_1 = 0$

Obtain
reward $r_2 = 5$

Action $a_1$ : "right"

Action $a_2$ : "fire"   (kill an alien)

# RL for NLG



$r = 0$     $r = 0$

Action Taken ← B     A     A → Reward: R("BAA", ref)

Actions Set ← A B     A B     A B

<BOS>

B

A

observation

The action taken influences the next observation

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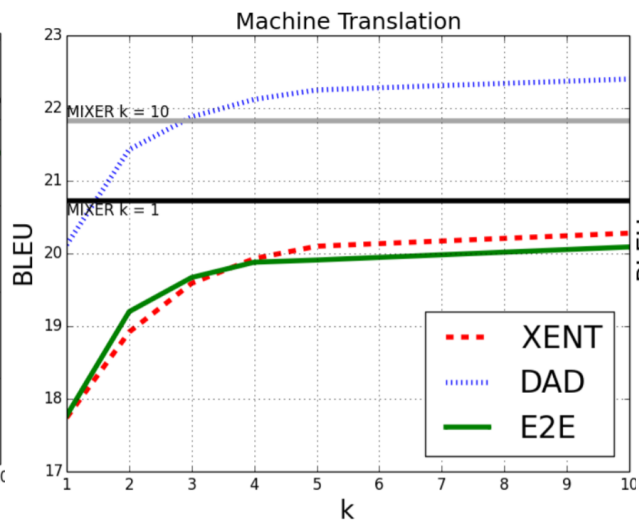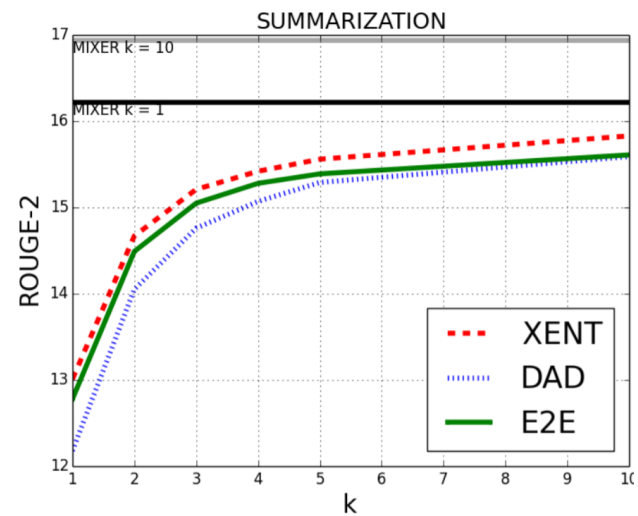
# RL for NLG

# RL for NLG

# RL-Based Summarization

◎ RL: directly optimize ROUGE-L

◎ ML+RL: MLE + RL for optimizing ROUGE-L

**Automatic**

| Model | ROUGE-1 | ROUGE-2 | ROUGE-L |
|---|---|---|---|
| ML, no intra-attention | 44.26 | 27.43 | 40.41 |
| ML, with intra-attention | 43.86 | 27.10 | 40.11 |
| RL, no intra-attention | **47.22** | 30.51 | **43.27** |
| ML+RL, no intra-attention | 47.03 | **30.72** | 43.10 |

**Human**

| Model | Readability | Relevance |
|---|---|---|
| ML | 6.76 | 7.14 |
| RL | 4.18 | 6.32 |
| ML+RL | **7.04** | **7.45** |

Using RL instead of ML achieves higher ROUGE scores, but lower human scores.

Hybrid is the best.

# Concluding Remarks

- ◉ NLG / Conditional NLG
- ◉ Decoding Algorithm
  - ○ Greedy
  - ○ Beam Search
  - ○ Sampling
- ◉ Evaluation
  - ○ Overall Quality → Specific Aspects
- ◉ Reinforcement Learning for NLG
  - ○ Directly optimizing the target score