# Generative Artificial Intelligence

Decoding Strategies and Evaluations for Natural Language Generation

# Outline

- Recap: Language Generation

- Decoding Strategies

  - Greedy Decoding

  - Beam Search

  - Top-k / Top-p Sampling

- Evaluations

# Natural Language Generation (NLG)

- Natural language generation (NLG) is a **process** that **outputs** text.

- NLG includes a wide variety of NLP tasks.

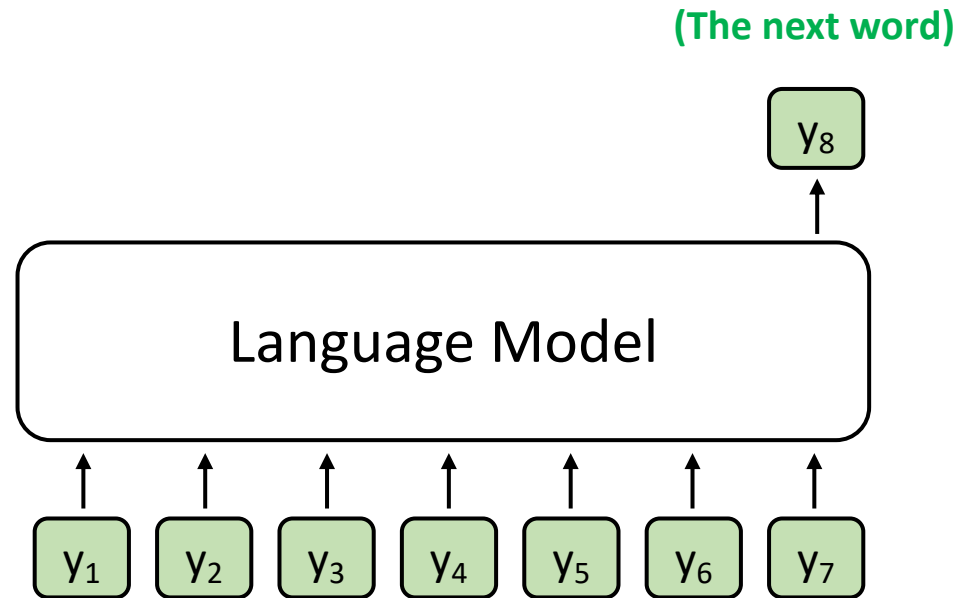Machine Translation   Abstractive Summarization   Dialogue Generation (e.g., ChatGPT)   Story Generation   Image Captioning   ...
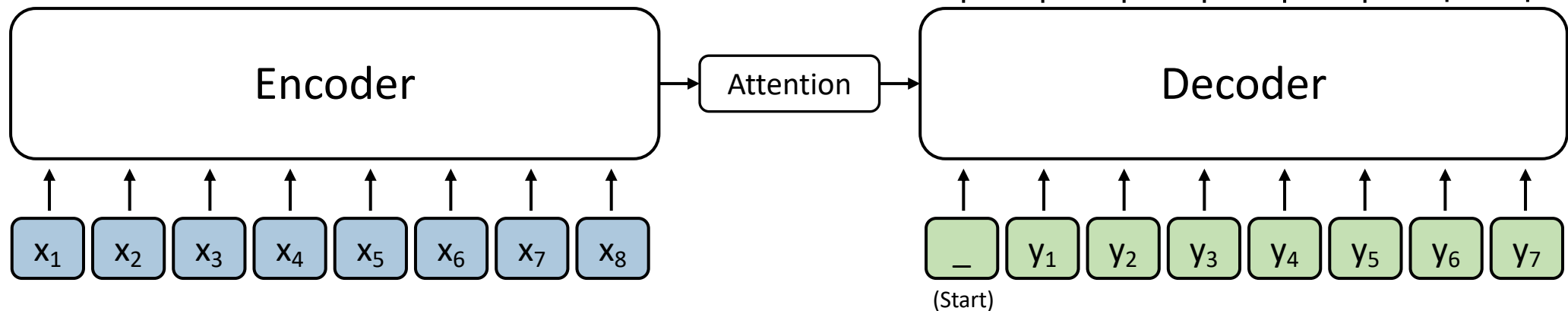
# Recap: Language Model

**(The next word)**

$y_8$

Language Model

$y_1$ $y_2$ $y_3$ $y_4$ $y_5$ $y_6$ $y_7$

$$P(y_t|y_1, y_2, \ldots, y_{t-1})$$

- A model that assigns probabilities to upcoming words is called **a language model**.
- The task involving predictions of upcoming words is **language modeling**.

# Recap: Conditional Language Model

- In addition to previous words, a conditional language model is provided with source text $x$.
- Also referred to sequence-to-sequence models.

$$P(y_t|y_1, y_2, \ldots, y_{t-1}, x)$$

(Target output)

# Tasks of Conditional Language Model

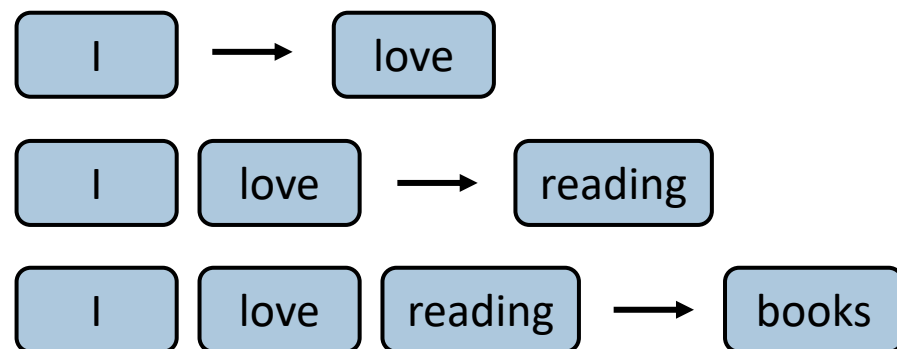- In addition to previous words (target), a conditional language model is provided with source text $x$.

| | Source | Target |
|---|---|---|
| Machine Translation | Language A | Language B |
| Summarization | Long Text | Concise Text |
| Dialogue Generation | User Input | Desired User Input |
| ... | | |

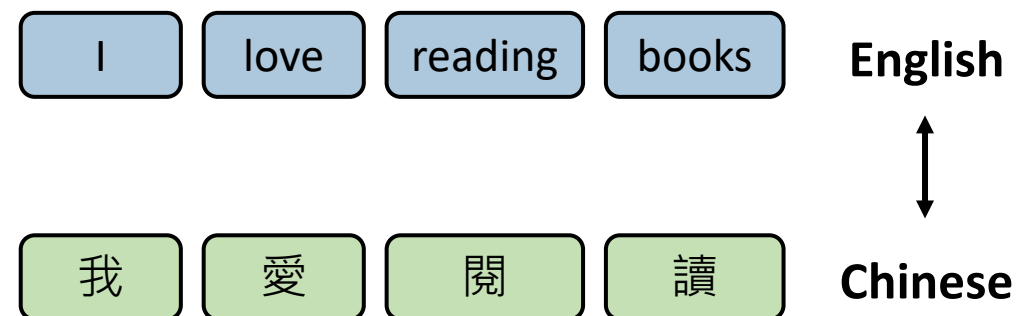# How to train a (Conditional) Language Model?

- First, you need a training corpus.

**Example: I love reading books.**

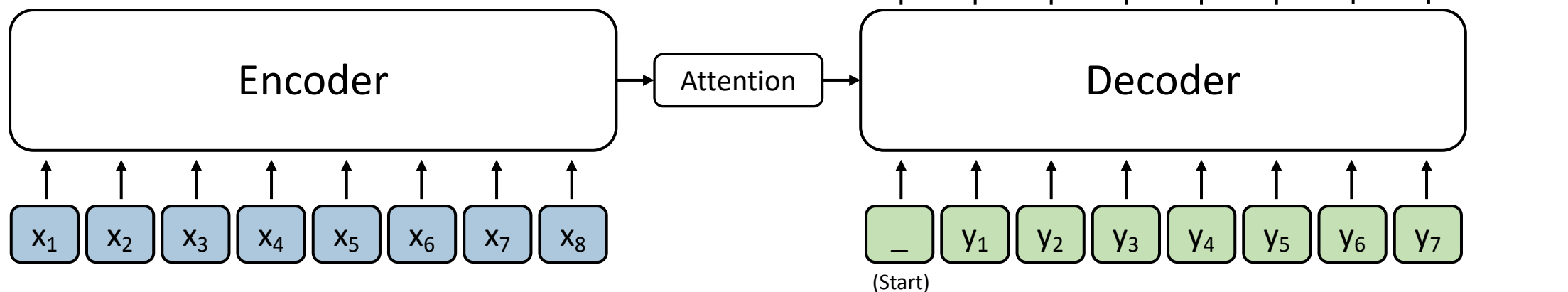**Language modeling (Unsupervised)**



**Machine Translation (Supervised)**
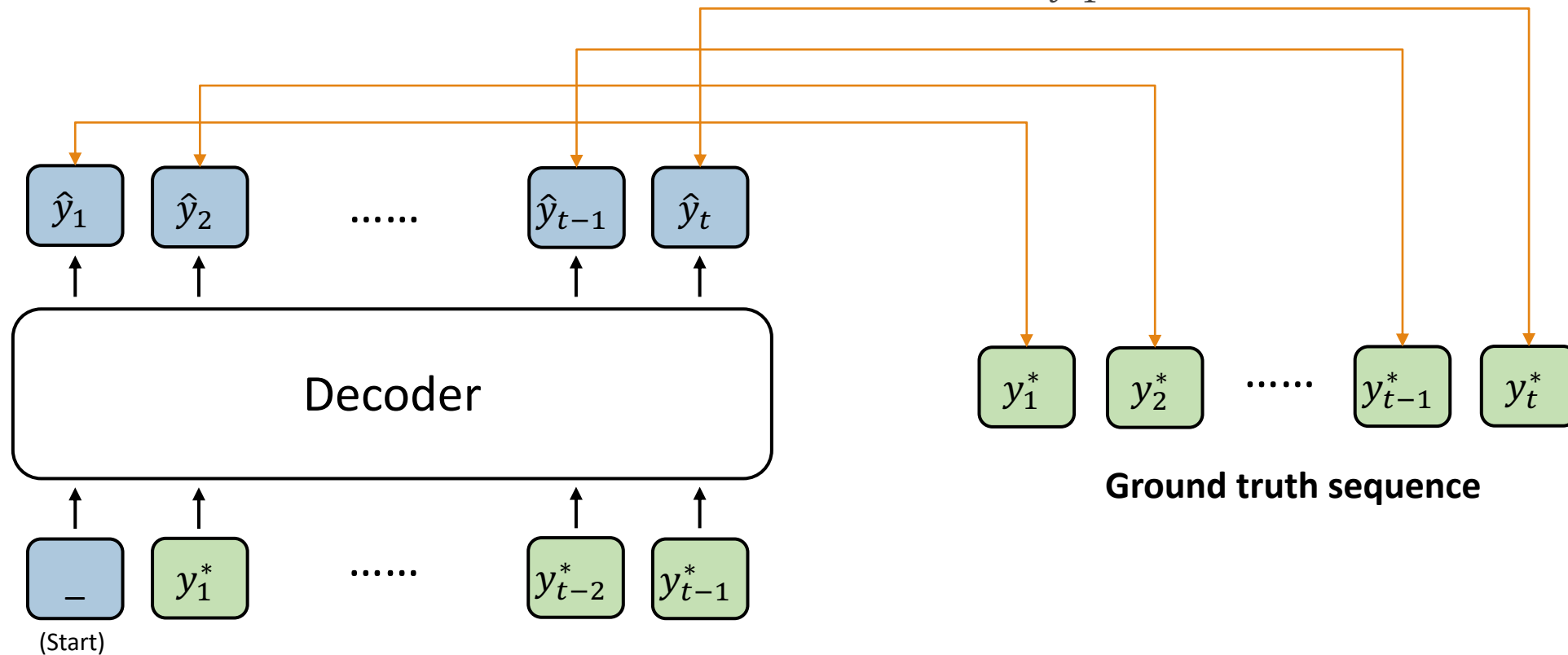
# How to train a (Conditional) Language Model?

- Use the Teacher Forcing technique during training.
- Total loss for a sequence: $\sum_1^T l_t$
  - $T$: Sequence length

**Teacher Forcing**

| $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ | $y_6$ | $y_7$ | $y_8$ |

| $l_1$ | $l_2$ | $l_3$ | $l_4$ | $l_5$ | $l_6$ | $l_7$ | $l_8$ |

Cross-entropy

| $\hat{y}_1$ | $\hat{y}_2$ | $\hat{y}_3$ | $\hat{y}_4$ | $\hat{y}_5$ | $\hat{y}_6$ | $\hat{y}_7$ | $\hat{y}_8$ |

**Encoder**

**Attention**

**Decoder**

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | $x_6$ | $x_7$ | $x_8$ |

| _ | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ | $y_6$ | $y_7$ |

(Start)

# Teacher Forcing – Training stage

**During training:**

$$L_{ml} = -\sum_{t=1}^{n'} \log p(y_t^* | y_1^*, \ldots, y_{t-1}^*, x)$$



$\hat{y}_1$   $\hat{y}_2$   ......   $\hat{y}_{t-1}$   $\hat{y}_t$

Decoder

_   $y_1^*$   ......   $y_{t-2}^*$   $y_{t-1}^*$

(Start)

$y_1^*$   $y_2^*$   ......   $y_{t-1}^*$   $y_t^*$

**Ground truth sequence**

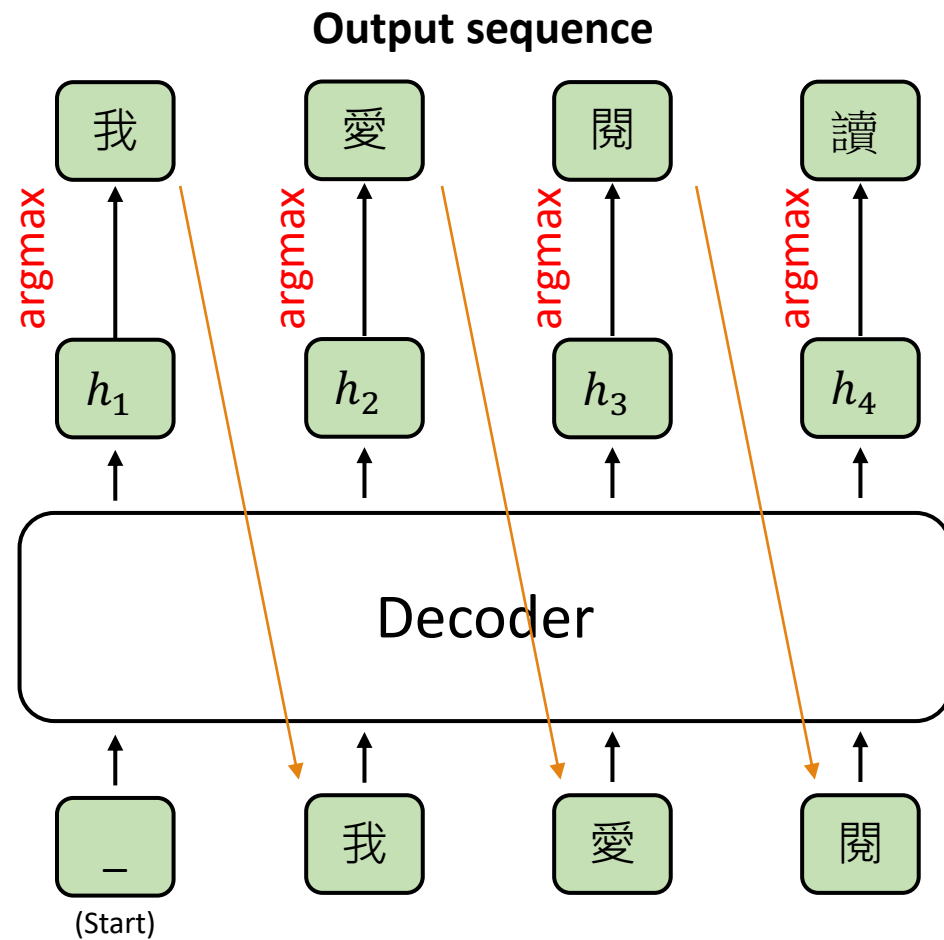# Teacher Forcing – Testing stage
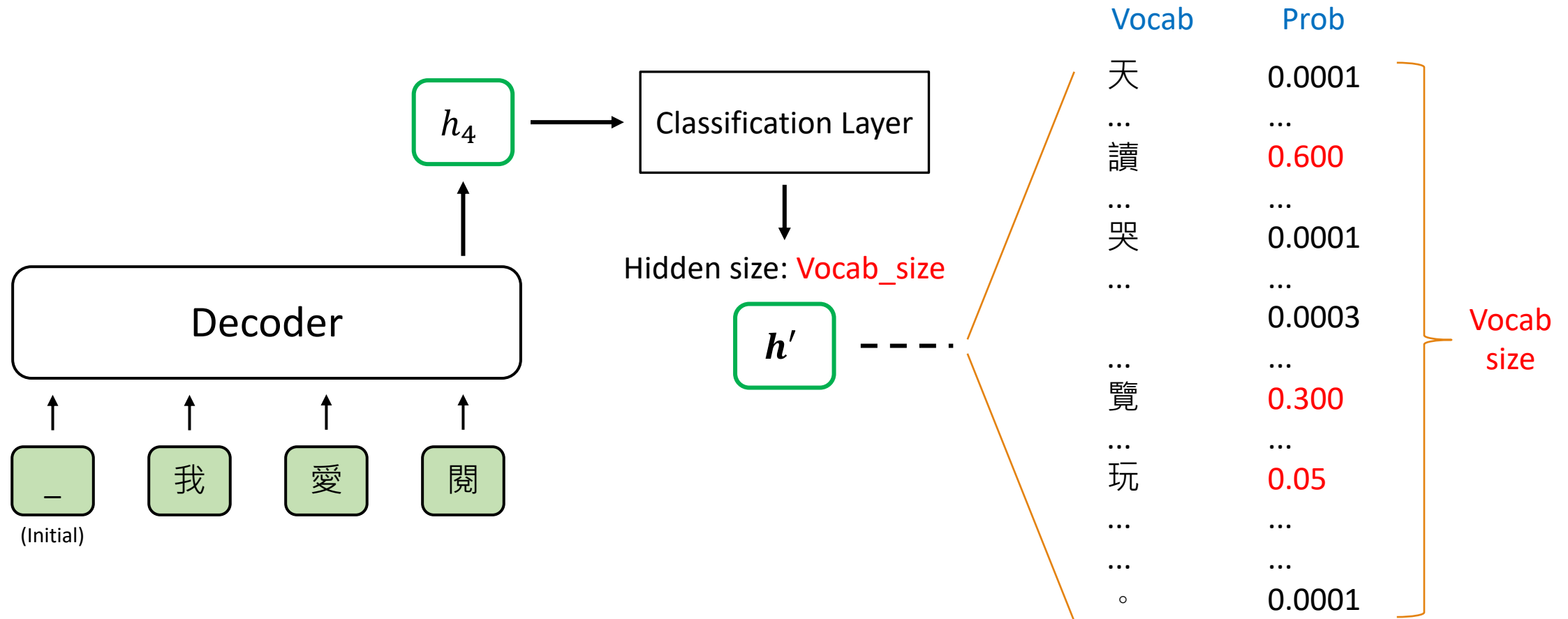
**During testing:**

**Output sequence**



- Advantage: stabilize training and increase performance
- Question: How does the next word be determined?

# Greedy Decoding

**Example: I love reading books.**

**Output sequence**

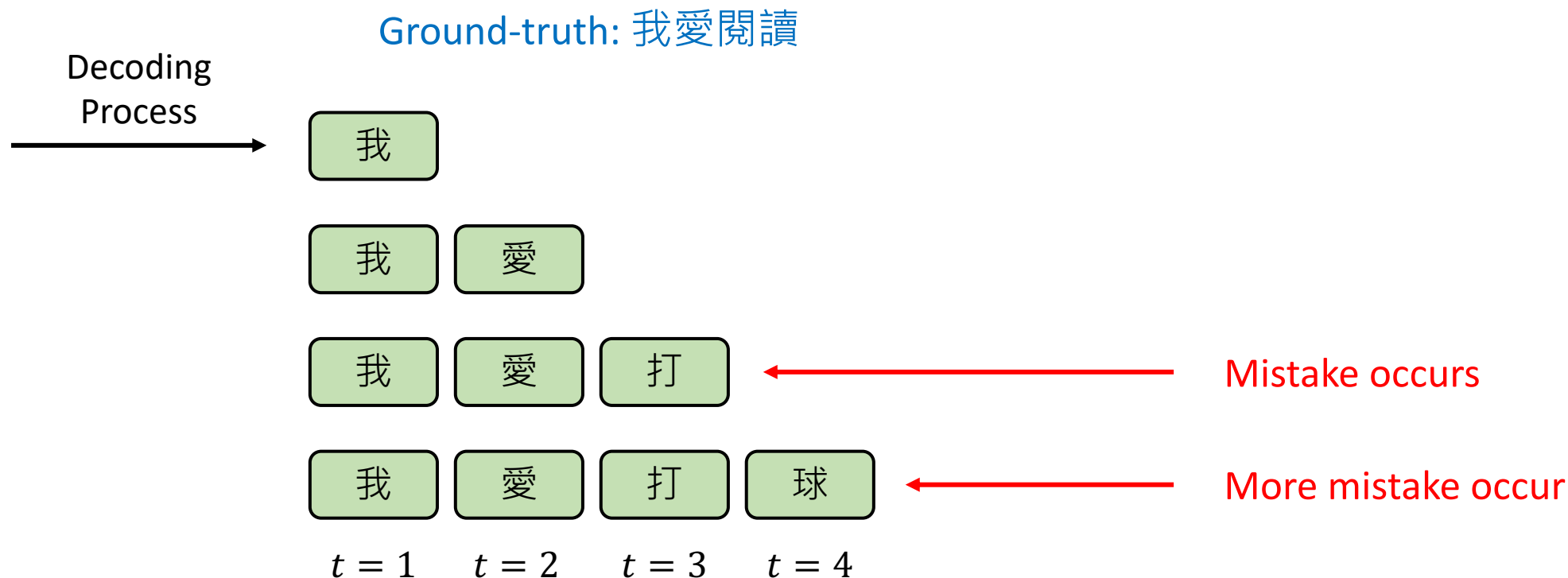# Greedy Decoding – Best Selection Process

# Problem of Greedy Decoding

- Greedy decoding cannot undo!

Ground-truth: 我愛閱讀

Decoding Process

我

我 愛

我 愛 打 ← Mistake occurs

我 愛 打 球 ← More mistake occur

$t = 1$      $t = 2$      $t = 3$      $t = 4$

# Re-thinking Greedy Decoding

- Greedy decoding cannot undo!

- Greedy decoding only provides one best choice at each time step.

- How about providing more than one choices at each time step?
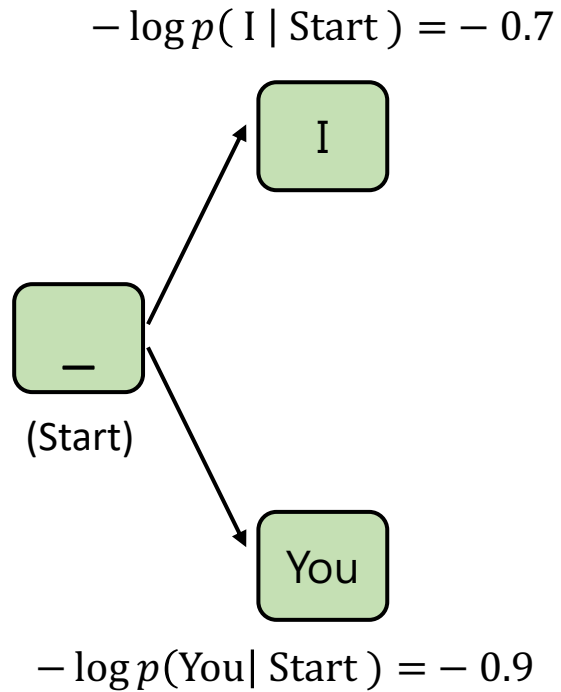
    → **Beam Search**

# Beam Search

- Set the `Beam size` (or `Beam width`) = 2

  - This means that the number of candidates will be preserved at each decoding time.

  - Beam size is a hyperparameter for beam search decoding.

- At each decoding time step, a score is calculated via the following equation:

$$L_{ml} = \sum_{t=1}^{n'} \log p(y_t^* | y_1^*, \dots, y_{t-1}^*, x)$$

# Beam Search ($t = 1$)

$-\log p(\text{I} \mid \text{Start}) = -0.7$
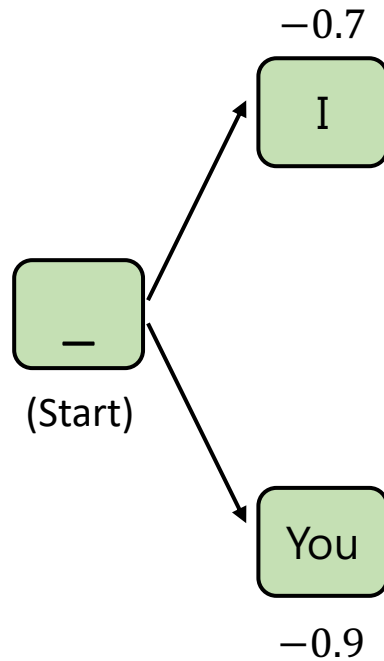
I

- At this decoding step, two choices are preserved.

_

(Start)

You

$-\log p(\text{You} \mid \text{Start}) = -0.9$

# Beam Search ($t = 1$)

$-0.7$

I

_
(Start)

You

$-0.9$

- At this decoding step, two choices are preserved.

# Beam Search ($t = 2$)

$$-0.7 - \log p(\text{ like } | \text{ Start, I}) = \boxed{-1.7}$$

like

$-0.7$

I

want

$$-0.7 - \log p(\text{ want } | \text{ Start, I}) = -2.9$$

$$-0.9 - \log p(\text{ want } | \text{ Start, You}) = \boxed{-1.6}$$

_

(Start)

want

You

are

$-0.9$

$$-0.9 - \log p(\text{ are } | \text{ Start,You}) = -1.8$$

Note the negative loglikelihood! Lower is better!

- At this decoding step, two choices are preserved, and the other two are discarded.

# Beam Search ($t = 2$)

$-1.7$

like

$-0.7$

I

want

$-2.9$

_

(Start)

$-1.6$
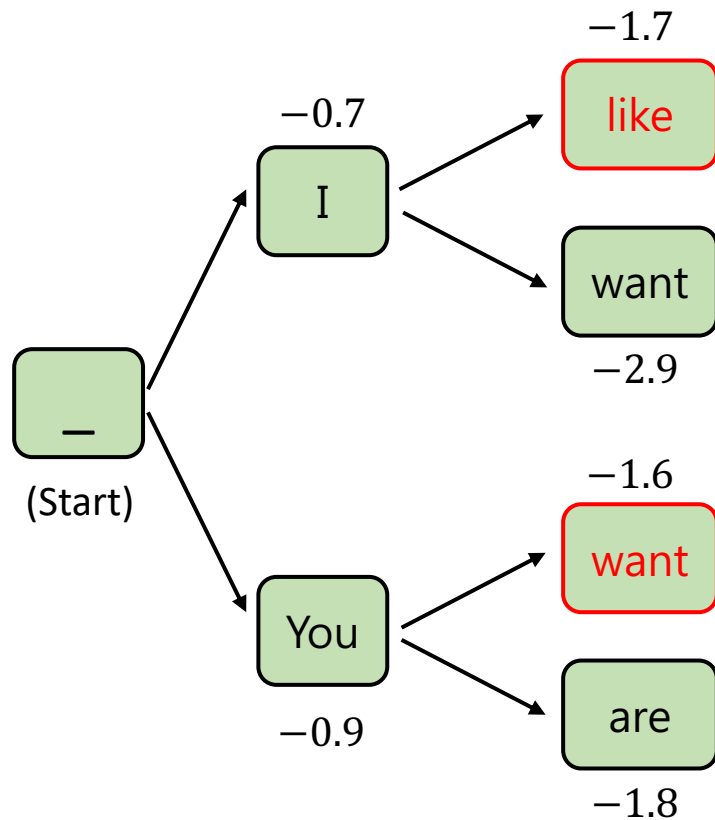
want

You

are

$-0.9$

$-1.8$

# Beam Search ($t = 3$)

$-0.7$

$-1.7$

like

$-1.7 - \log p(\text{ to } | \text{ Start, I, like}) = -2.8$

$-1.7 - \log p(\text{ books } | \text{ Start, I, like}) = -2.5$

to

books

I

want

$-2.9$

_

(Start)

$-1.6$

want

$-1.6 - \log p(\text{ to } | \text{ Start,You, want}) = -2.9$

$-1.7 - \log p(\text{ tea } | \text{ Start, You, want}) = -3.8$

to

tea

You

are

$-0.9$

$-1.8$

# Beam Search ($t = 3$)

$-0.7$

$-1.7$

$-1.7 - \log p(\text{ to } | \text{ Start, I, like}) = -2.8$

$-1.7 - \log p(\text{ books } | \text{ Start, I, like}) = -2.5$

$-2.9$

(Start)

$-1.6$

$-1.6 - \log p(\text{ to } | \text{ Start, You, want}) = -2.9$

$-1.7 - \log p(\text{ tea } | \text{ Start, You, want}) = -3.8$

$-0.9$

$-1.8$

# Beam Search ($t = 3$)

$-2.8$

to

$-1.7$

like

$-0.7$

I

books

$-2.5$

want

$-2.9$

_

(Start)

$-2.9$

to

$-1.6$

want

You

tea

$-3.8$

$-0.9$

are

$-1.8$

# Beam Search ($t = 4$)

# Beam Search ($t = 4$)

−4.1 read

−2.8 to

−1.7 like

−0.7 I

−2.9 want

−3.2 watch

−2.5 books

−3.3 and

−3.5 party

_ (Start)

−0.9 You

−1.6 want

−1.8 are

−2.9 to

−3.8 tea

# Beam Search ($t = 5$)

−4.1
read

−4.8
tv

−2.8
to

−1.7
like

−0.7
I

−3.2
watch

−3.7
horror

−2.5
books

want
−2.9

_
(Start)

−3.3
and

−4.5
write

−3.5
party

−4.3
horror

−2.9
to

−1.6
want

You

−3.8
tea

−0.9

are
−1.8

# Beam Search ($t = 5$)

−4.1 read
−4.8 tv

−2.8 to

−1.7 like

−0.7 I

watch −3.2
horror −3.7

want −2.9

books −2.5

_ (Start)

and −3.3
write −4.5

−1.6 want
−2.9 to

You

party −3.5
horror −4.3

−0.9

are −1.8

tea −3.8

# Beam Search ($t = 6$)

# Beam Search ($t = 6$)

# Stop Criterion

- There are two common stop criterion, either for greedy decoding or beam search decoding:

    - We consider a sequence of generation complete when the <EOS> token is produced by a model. *<EOS>: End of sequence

    - E.g., <Start> I like to watch horror movies <EOS>

- A generated sequence reaches a pre-defined maximal length.

# Problem of Beam Search

- Longer candidates will have lower scores.

- (Let's see again the 6th time step)

# Beam Search ($t = 6$)

# Problem of Beam Search

- Longer candidates will have lower scores.

- Solution: Perform normalization to penalize on length

$$L_{ml} = \frac{1}{T} \sum_{t=1}^{T} \log p(y_t^* | y_1^*, \ldots, y_{t-1}^*, x)$$

# How to evaluate natural language generation?

- Natural language is hard to evaluate due to <u>subjectivity</u> and language <u>diversity</u>.

**For example: Machine Translation**

| 我 | 愛 | 閱 | 讀 |
| --- | --- | --- | --- |

| I | love | reading | books |
| --- | --- | --- | --- |
**(Source language)**

| 我 | 愛 | 讀 | 書 |
| --- | --- | --- | --- |

**(Target language)**

- Human evaluations

- Automatic evaluations (We will focus on this topic.)

# BLEU (Bilingual Evaluation Understudy)

- A word-based metric.

  - It is very sensitive to word tokenization

- Core concept: Compute <span style="color:red">precision</span> for n-grams:

  - Unigrams -> BLEU-1

  - Bigrams -> BLEU-2

  - Trigrams -> BLEU-3

  - 4-grams -> BLEU-4

Papineni, Kishore, et al. "Bleu: a method for automatic evaluation of machine translation."
Proceedings of the 40th annual meeting of the Association for Computational Linguistics. 2002.

# Precision and Recall

$$\text{Precision} = \frac{\text{Relevant and retrieved instances}}{\text{All } \textcolor{red}{\text{retrieved}} \text{ instances}}$$ ← Predicted by a model

$$\text{Recall} = \frac{\text{Relevant and retrieved instances}}{\text{All } \textcolor{red}{\text{relevant}} \text{ instances}}$$ ← Ground-truths

Relevant and retrieved instances: Intersection between predictions and ground-truths

# Calculation of BLEU Score (Example)

Assume we now translate from Chinese to English.

**Calculate BLEU-1 score**

Chinese: 我想要讀那本書

Reference1: I want to read the book.

Reference2: I want to read that book.

Model output: the the the the the the.

# Calculation of BLEU Score (Example)

Assume we now translate from Chinese to English.

**Calculate BLEU-1 score**

Chinese: 我想要讀那本書

Reference1: I want to read the book.

Reference2: I want to read that book.

Model output: the the the the the the.

Precision: $\dfrac{6}{6}$

100%! Can this be true?

# Calculation of BLEU Score (Example)

Assume we now translate from Chinese to English.

**Calculate BLEU-1 score**

Chinese: 我想要讀那本書

Reference1: I want to read <u>the</u> book.

Reference2: I want to read that book.

Model output: <u>the</u> <u>the</u> <u>the</u> <u>the</u> <u>the</u> <u>the</u>.

Precision: $\frac{6}{6}$ ❌

Modified Precision: $\frac{1}{6}$

Papineni, Kishore, et al. "Bleu: a method for automatic evaluation of machine translation."
Proceedings of the 40th annual meeting of the Association for Computational Linguistics. 2002.

# Why should we use modified precision?

- The output sequences can be total mistakes.

  - E.g., the the the the the the

- Original precision is in favor of longer output sequences.

- Therefore, we should use modified precision to prevent bad evaluations.

Papineni, Kishore, et al. "Bleu: a method for automatic evaluation of machine translation."
Proceedings of the 40th annual meeting of the Association for Computational Linguistics. 2002.

# Calculation of BLEU Score (Example)

**Calculate BLEU-2 score**

Reference1: The dog is on the bed.

Reference2: There is a dog on the bed. ← More than one references can be provided for machine translation!

Model output: The dog the dog on the bed.

Papineni, Kishore, et al. "Bleu: a method for automatic evaluation of machine translation."
Proceedings of the 40th annual meeting of the Association for Computational Linguistics. 2002.

# Calculation of BLEU Score (Example)

**Calculate BLEU-2 score**

Reference1: The dog is on the bed.

Reference2: There is a dog on the bed.

Model output: The dog the dog on the bed.

|  | Count |  |
|---|---|---|
| the dog | 2 | (duplicated) |
| dog the | 1 | |
| dog on | 1 | |
| on the | 1 | |
| the bed | 1 | |

Papineni, Kishore, et al. "Bleu: a method for automatic evaluation of machine translation."
Proceedings of the 40th annual meeting of the Association for Computational Linguistics. 2002.

# Calculation of BLEU Score (Example)

Clips to the reference

**Calculate BLEU-2 score**

Reference1: The dog is on the bed.

Reference2: There is a dog on the bed.

Model output: The dog the dog on the bed.

| | Count | Count$_{clip}$ |
|---|---|---|
| the dog | 2 | 1 |
| dog the | 1 | |
| dog on | 1 | |
| on the | 1 | |
| the bed | 1 | |

Papineni, Kishore, et al. "Bleu: a method for automatic evaluation of machine translation."
Proceedings of the 40th annual meeting of the Association for Computational Linguistics. 2002.

# Calculation of BLEU Score (Example)

**Calculate BLEU-2 score**

Reference1: The dog is on the bed.

Reference2: There is a dog on the bed.

Model output: The dog the dog on the bed.

| | Count | Count$_{clip}$ |
|---|---|---|
| the dog | 2 | 1 |
| dog the | 1 | 0 |
| dog on | 1 | |
| on the | 1 | |
| the bed | 1 | |

Papineni, Kishore, et al. "Bleu: a method for automatic evaluation of machine translation." Proceedings of the 40th annual meeting of the Association for Computational Linguistics. 2002.

# Calculation of BLEU Score (Example)

**Calculate BLEU-2 score**

Reference1: The dog is on the bed.

Reference2: There is a <u>dog on</u> the bed.

Model output: The dog the <u>dog on</u> the bed.

|  | Count | Count$_{clip}$ |
|---|---|---|
| the dog | 2 | 1 |
| dog the | 1 | 0 |
| dog on | 1 | <span style="color:red">1</span> |
| on the | 1 | |
| the bed | 1 | |

Papineni, Kishore, et al. "Bleu: a method for automatic evaluation of machine translation."
Proceedings of the 40th annual meeting of the Association for Computational Linguistics. 2002.

# Calculation of BLEU Score (Example)

**Calculate BLEU-2 score**

Reference1: The dog is <u>on the</u> bed.

Reference2: There is a dog <u>on the</u> bed.

Model output: The dog the dog <u>on the</u> bed.

Count <span style="color:red">only one time</span> even mapped to both references.

|  | Count | Count$_{clip}$ |
|---|---|---|
| the dog | 2 | 1 |
| dog the | 1 | 0 |
| dog on | 1 | 1 |
| on the | 1 | 1 |
| the bed | 1 |  |

Papineni, Kishore, et al. "Bleu: a method for automatic evaluation of machine translation."
Proceedings of the 40th annual meeting of the Association for Computational Linguistics. 2002.

# Calculation of BLEU Score (Example)

**Calculate BLEU-2 score**

Reference1: The dog is on the bed.

Reference2: There is a dog on the bed.

Model output: The dog the dog on the bed.

Count only one time even mapped to both references.

|  | Count | Count$_{clip}$ |
|---|---|---|
| the dog | 2 | 1 |
| dog the | 1 | 0 |
| dog on | 1 | 1 |
| on the | 1 | 1 |
| the bed | 1 | 1 |

Modified Precision: $\frac{4}{6}$

Papineni, Kishore, et al. "Bleu: a method for automatic evaluation of machine translation."
Proceedings of the 40th annual meeting of the Association for Computational Linguistics. 2002.

# Formula of BLEU Score

Summation for unigram, bigram, tri-gram, and 4-gram

$$p_n = \frac{\displaystyle\sum_{C \in \{Candidates\}} \sum_{n\text{-}gram \in C} Count_{clip}(n\text{-}gram)}{\displaystyle\sum_{C' \in \{Candidates\}} \sum_{n\text{-}gram' \in C'} Count(n\text{-}gram')}$$

Summation for all candidates (model outputs)
of each translation

Papineni, Kishore, et al. "Bleu: a method for automatic evaluation of machine translation."
Proceedings of the 40th annual meeting of the Association for Computational Linguistics. 2002.

# What we've learned BLEU so far

- The BLEU score is calculated from the summation of 1-gram to 4-gram.

  - You can also measure n-gram individually.

- We use modified precision to prevent bad evaluations.

- What will happen if a model tends to generate really short sentences?

  ➡️ **More penalty for calculating BLEU score!**

Papineni, Kishore, et al. "Bleu: a method for automatic evaluation of machine translation."
Proceedings of the 40th annual meeting of the Association for Computational Linguistics. 2002.

# Brevity Penalty (BP)

- BP is used to penalize short candidates.

$c$: The length of a candidate sequence
$r$: The length of a reference sequence that is closest to $c$ (shorter one)

$$\text{BP} = \begin{cases} 1 & \text{if } c > r \\ e^{(1-r/c)} & \text{if } c \leq r \end{cases}$$

Then,

$$\text{BLEU} = \text{BP} \cdot \exp\left( \sum_{n=1}^{N} w_n \log p_n \right)$$

$N$=4 to include 1-gram to 4-gram

Weight for each $n$-gram (was set 1/4 in the original paper)

Papineni, Kishore, et al. "Bleu: a method for automatic evaluation of machine translation."
Proceedings of the 40th annual meeting of the Association for Computational Linguistics. 2002.

# (Recap) Perplexity

Perplexity (PPL) is a quantitative criterion used to evaluate the capacities of language modeling models.

- Given the sequence of words $W = w_1 w_2 \ldots w_N$ and an N-gram model. The PPL of the model was computed by:

$$Perplexity(W) = P(w_1 w_2 \ldots w_N)^{-\frac{1}{N}} = \sqrt[N]{\prod_{k=1}^{n} \frac{1}{P(w_k | w_{k-N+1:k})}}$$

The lower the value of perplexity, the better the language modeling capability of the model.

# Comparison for Human and Automatic Evaluations

- <span style="color:red">Human evaluations</span>

  - Pros: More accurate for subjectivity, flexibility for any desired comparison

  - Cons: Less objective, time-consuming, expensive

- Automatic evaluations

  - Pros: Objective enough to serve as common evaluation metrics, fast

  - Cons: Cannot meet language diversity

    - Take machine translation for instance, there are always other valid ways to translate the source sentence.