

- Neural Machine Translation
- ✓

Video: Course 4 Introduction
2 min
- ✓

Reading: Connect with your mentors and fellow learners on Slack!
10 min
- ✓

Video: Seq2seq
4 min
- ✓

Video: Alignment
4 min
- ✓

Reading: Background on seq2seq
10 min
- ✓

Reading: (Optional): The Real Meaning of Ich Bin ein Berliner
10 min
- ✓

Video: Attention
6 min
- ✓

Reading: Attention
10 min
- ✓

Video: Setup for Machine Translation
3 min
- 📄

Lab: Ungraded Lab: Stack Semantics
30 min
- ✓

Video: Training an NMT with Attention
6 min
- ✓

Reading: Training an NMT with Attention
10 min
- ✓

Reading: (Optional) What is Teacher Forcing?
10 min
- ✓

Video: Evaluation for Machine Translation
8 min
- 📄

Reading: Evaluation for Machine Translation
10 min
- 📄

Lab: Ungraded Lab: BLEU Score
30 min
- ▶

Video: Sampling and Decoding
9 min
- 📄

Reading: Sampling and Decoding
10 min
- 📄

Reading: Content Resource
10 min
- Assignment
- Heroes of NLP: Oren Etzioni

Evaluation for Machine Translation

The closer the BLEU score is to one, the better your model is. The closer to zero, the worse it is.

To get the BLEU score, the candidates and the references are usually based on an average of uni, bi, tri or even four-gram precision. To demonstrate, I'll use uni-grams as an example. Look at the following table:

Candidate	I	I	am	I	I
Reference 1	Younes	said	I	am	hungry
Reference 2	He	said	I	am	hungry

To calculate the BLEU score you can do the following.

"I" appears at most once in both, so clip to one:

$$m_w = 1$$

(Sum over unique n-gram counts in the candidate)

(total # of words in candidate)

You would sum over the unique n-gram counts in the candidate and divide by the total number of words in the candidate. The same concept could apply to unigrams, bigrams, etc. One issue with the BLEU score is that it does not take into account semantics, so it does not take into account the order of the n-grams in the sentence.

Another similar method for evaluation is the ROUGE score which calculates precision and recall for machine texts by counting the n-gram overlap between the machine texts and a reference text. Here is an example that calculates recall:

Recall in ROUGE

Model	The	cat	had	striped	orange	fur
Reference	The	cat	had	orange	fur	

(Sum of overlapping unigrams in model and reference)

(total # of words in reference)

5

5

Recall = 1

Rouge also allows you to compute precision as follows:

Precision in ROUGE

Model	The	cat	had	striped	orange	fur
Reference	The	cat	had	orange	fur	

(Sum of overlapping unigrams in model and reference)

(total # of words in model)

5

6

Precision = 0.83

Mark as completed