

# COLX-531: Neural Machine Translation

**Muhammad Abdul-Mageed**

[muhammad.mageed@ubc.ca](mailto:muhammad.mageed@ubc.ca)

**Deep Learning & NLP Lab**

The University of British Columbia

# Table of Contents

## 1 MT Evaluation

Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL), Philadelphia, July 2002, pp. 311-318.

## BLEU: a Method for Automatic Evaluation of Machine Translation

**Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu**

IBM T. J. Watson Research Center

Yorktown Heights, NY 10598, USA

{papineni,roukos,toddward,weijing}@us.ibm.com

### Abstract

Human evaluations of machine translation are extensive but expensive. Human evaluations can take months to finish and involve human labor that can not be reused. We propose a method of automatic machine translation evaluation that is quick, inexpensive, and language-independent, that correlates highly with human eval-

the evaluation bottleneck. Developers would benefit from an inexpensive automatic evaluation that is quick, language-independent, and correlates highly with human evaluation. We propose such an evaluation method in this paper.

### 1.2 Viewpoint

How does one measure translation performance?  
*The closer a machine translation is to a professional*

# MT Eval Rationale (Papineni et al., 2002)

## Why MT Eval Needed?

- Evaluating translation is hard. Why?
- Human evaluation is costly
- An automatic method promises to accelerate MT progress
- Measure MT based on **numerical closeness** to a human reference
- The metric can be used to **optimize** MT systems
- **BLEU Main Idea**: Use a weighted avg of variable-length phrase matches against reference translations
- i.e., **compare** n-grams of the candidate with the n-grams of the reference translation, and **count** the number of matches
- Matches are *position-independent*.
- **The more** the matches, **the better** the candidate translation is.

## Precision

- **Precision** counts up the number of candidate translation words (unigrams) which occur in any reference translation and then divides by the total number of words in the candidate translation.

### 1: Precision

$$\frac{\text{count\_candidate\_trans\_words\_in\_ref}}{\text{total\_words\_in\_candidate\_trans}}$$

Candidate: the the the the the the the.

Reference 1: The cat is on the mat.

Reference 2: There is a cat on the mat.

## Modified Precision

- A ref word should be considered *exhausted* after a matching candidate word is identified
- **Modified unigram precision:**
  - ① Count the **maximum** number of times a word occurs **in any single reference translation**.
  - ② **Clip** the **total count** of each candidate word **by its maximum reference count**.  $count_{clip} = \min(count, max\_ref\_count)$
  - ③ **Add clipped counts up**, and **divide** by the **total (unclipped) number of candidate words**

## 2: Modified Precision

$$\frac{count_{clip}}{total\_words\_in\_candidate\_trans}$$

## 3: Modified Precision

$$\text{count}_{\text{clip}} = \min(\text{count}, \text{max\_ref\_count})$$

$$\frac{\text{count}_{\text{clip}}}{\text{total\_words\_in\_candidate\_trans}}$$

Candidate: the the the the the the the.

Reference 1: The cat is on the mat.

Reference 2: There is a cat on the mat.

## Exercise

Calculate modified bi-gram precision for Candidate 1.

Candidate 1: It is a guide to action which ensures that the military always obeys the commands of the party.

Candidate 2: It is to insure the troops forever hearing the activity guidebook that party direct.

Reference 1: It is a guide to action that ensures that the military will forever heed Party commands.

Reference 2: It is the guiding principle which guarantees the military forces always being under the command of the Party.

Reference 3: It is the practical guide for the army always to heed the directions of the party.



# Candidate Bi-grams Count<sub>clip</sub>

Candidate 1: It is a guide to action which ensures that the military always obeys the commands of the party.

- ① It is      ② is a  
③ a guide    ④ guide to  
⑤ to action   ⑥ action which  
⑦ which ensures   ⑧ ensures that  
⑨ that the      ⑩ the military  
⑪ military always   ⑫ always obeys  
⑬ obeys the      ⑭ the commands  
⑮ commands of   ⑯ of the  
⑰ the party

Figure: There are 17 bi-grams in Candidate 1. (Count<sub>clip</sub>)

# Calculating Modified Bi-gram Precision

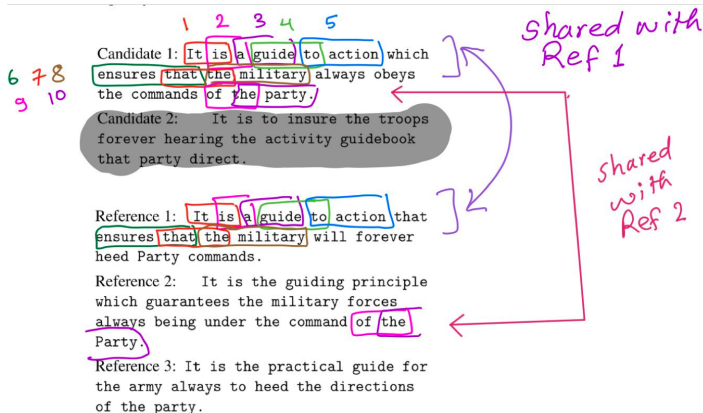


Figure: **Modified precision** =  $10/17$ .  $\text{Count}_{\text{clip}} = 10$ . Total bi-gram count in **Candidate 1** = 17. No need to consider bi-grams in Ref 2 and Ref 3 that already occur in Ref 1.

# Modified N-Gram Precision on **Blocks of Text**

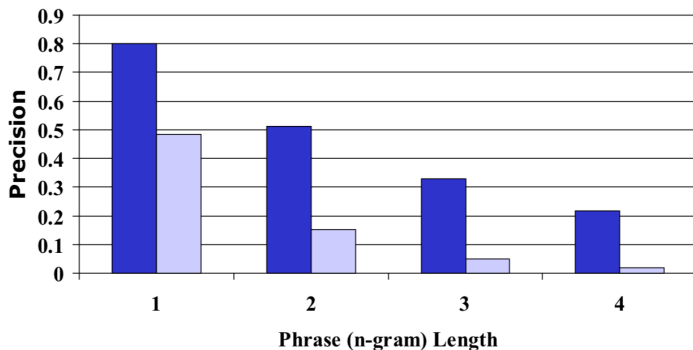
## Blocks of Text

- 1 Compute the n-gram matches sentence by sentence.
- 2 Add the clipped n-gram counts for all the candidate sentences
- 3 Divide by the number of candidate n-grams in the test corpus

4:  $P_n$

$$P_n = \frac{\sum_{C \in \text{Candidates}} \sum_{n\text{-gram} \in C} \text{count}_{\text{clip}}(n\text{-gram})}{\sum_{C' \in \text{Candidates}} \sum_{n\text{-gram}' \in C'} \text{Count}(n\text{-gram}' )}$$

Figure 1: Distinguishing Human from Machine



# Combining the modified n-gram Precisions

## How should we combine the modified precisions for various n-gram sizes?

- **Modified n-gram precision decays roughly exponentially with n:** the modified un-igram precision is much larger than the modified bi-gram precision which in turn is much bigger than the modified tri-gram precision.
- A **weighted average of the logarithm** of modified precisions take this exponential decay into account.
- **BLEU** uses the **average logarithm with uniform weights**, which is equivalent to using the geometric mean of the modified n-gram precisions.

# Sentence Length

## Ensuring Suitable Length

- N-gram precision penalizes spurious words in candidate that do not appear in any ref
- Modified precision is penalized if a word occurs more frequently in a candidate than its max ref count
- However, modified n-gram precision *alone* fails to enforce the proper translation length
- See paper for illustrating examples

## Brevity

- Goal: **Make the brevity penalty 1.0** when the candidate's length is the same as any reference translation's length
- Example: If there are 3 refs with lengths 12, 15, and 17 words and the candidate is 12 words, **we want the brevity penalty to be 1** and call the closest ref sent length the "**best match length.**"
- If we compute brevity penalty sentence by sentence and averaged the penalties, then length deviations on short sentences would be punished harshly.
- Instead, **compute the brevity penalty over the entire corpus** to allow some freedom at the sentence level.
- **First**, compute the test corpus' effective reference length,  $r$ , by summing "**best match length.**" for each candidate sent in corpus.
- **Then**, choose penalty to be a decaying exponential in  $r/c$ , where  $c$  is total length of the candidate corpus.

# BLEU Details

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

Then,

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

The ranking behavior is more immediately apparent in the log domain,

$$\log \text{BLEU} = \min\left(1 - \frac{r}{c}, 0\right) + \sum_{n=1}^N w_n \log p_n.$$

In our baseline, we use  $N = 4$  and uniform weights  $w_n = 1/N$ .