

## Machine Translation



Dan Klein  
UC Berkeley

Many slides from John DeNero and  
Philip Koehn

### Translation Task

- Text is both the input and the output.
- Input and output have roughly the same information content.
- Output is more predictable than a language modeling task.
- Lots of naturally occurring examples.

### Translation Examples

### English-German News Test 2013 (a standard dev set)

Republican leaders justified their policy by the need to combat electoral fraud.

Die Führungskräfte der Republikaner  
The Executives of the republican  
rechtfertigen ihre Politik mit der  
justify your politics With of the  
Notwendigkeit, den Wahlbetrug zu  
need, the election fraud to  
bekämpfen.  
fight.

### Variety in Translations?

#### Human-generated reference translation

A small planet, whose is as big as could destroy a middle sized city, passed by the earth with a distance of 463 thousand kilometers. This was not found in advance. The astronomers got to know this incident 4 days later. This small planet is 50m in diameter. The astronomers are hard to find it for it comes from the direction of sun.

A volume enough to destroy a medium city small planet is 100,000 km<sup>3</sup> with within 463,000 kilometres of close however were not in advance discovered, astronomer just knew this matter after four days. This small planet diameter is about 50 metre, from the direction at sun, therefore astronomer very hard to discovers it.

From <https://nlp.stanford.edu/IRN2013/>

### Evaluation

### BLEU Score

BLEU score: geometric mean of 1-, 2-, 3-, and 4-gram precision vs. a reference, multiplied by brevity penalty (harshly penalizes translations shorter than the reference).

$$\text{Matched}_i = \sum_{t_i} \min \left\{ C_h(t_i), \max_j C_j(t_i) \right\}$$

Let "of the" appears twice in hypothesis h but only at most once in a reference, then only the first is "correct"

$$P_i = \frac{\text{Matched}_i}{H_i}$$

"clipped" precision of n-gram tokens

$$B = \exp \left\{ \min \left( 0, \frac{n-L}{n} \right) \right\}$$

brevity penalty only matters if the hypothesis corpus is shorter than the sum of (shortest) references

$$\text{BLUE} = B \left( \prod_{i=1}^4 P_i \right)^{\frac{1}{4}}$$

BLUE is a mean of clipped precisions, scaled down by the brevity penalty.

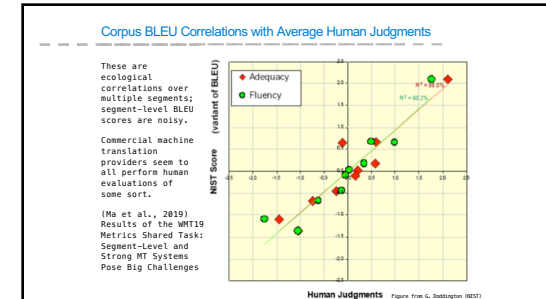
### Evaluation with BLEU

In this sense, the measures will partially undermine the American democratic system.

In this sense, these measures partially undermine the democratic system of the United States.

BLEU = 26.52, 75.0/40.0/21.4/7.7 (BP=1.000, ratio=1.143, hyp\_len=16, ref\_len=14)

(Papinen et al., 2002) BLEU: a method for automatic evaluation of machine translation.



### Human Evaluations

**Direct assessment: adequacy & fluency**

- Monolingual: Ask humans to compare machine translation to a human-generated reference. (Easier to source annotators)
- Bilingual: Ask humans to compare machine translation to the source sentence that was translated. (Compares to human quality)
- Annotators can assess segments (sentences) or whole documents.
- Segments can be assessed with or without document context.

**Ranking assessment:**

- Raters are presented with 2 or more translations.
- A human-generated reference may be provided, along with the source.
- <https://www.isi-ranking-experiments.org/> **Ranking** Document Source: **WMT2019** Document Source: **WMT2019** Findings of the 2021 Conference on Machine Translation

### Translationese and Evaluation

Translated text can: (Baker et al., 1993; Graham et al., 2019)

- be more explicit than the original source
- be less ambiguous
- be simplified (lexically, syntactically, and stylistically)
- display a preference for conventional grammaticality
- avoid repetition
- exaggerate target language features
- display features of the source language

"If we consider only original source text (i.e. not translated from another language, or translationese), then we find evidence showing that human parity has not been achieved." (Toral et al., 2018)

(Baker et al., 2003) Corpus Linguistics and Translation Studies: Implications for Machine Translation  
(Toral et al., 2018) Translationese in Machine Translation Evaluation  
(Toral et al., 2020) Enhancing the Effectiveness of Assessing Claims of Human Parity in Neural Machine Translation

### How are We Doing? Example: WMT 2019 Evaluation

**2019 segment-in-context direct assessment (Barrault et al., 2019):**

- ✓ German to English: many systems are tied with human performance;
- ✗ English to Gujarati: all systems are outperformed by the human translator;
- ✗ English to Chinese: all systems are outperformed by the human translator;
- ✗ English to Czech: all systems are outperformed by the human translator;
- ✗ English to Lithuanian: all systems are outperformed by the human translator;
- ✗ English to Finnish: all systems are outperformed by the human translator;
- ✓ English to German: Facebook-FAIR achieves super-human translation performance; several systems are tied with human performance;
- ✗ English to Russian: Facebook-FAIR is tied with human performance.

(Barrault et al., 2019) Findings of the 2019 Conference on Machine Translation (EMT19)

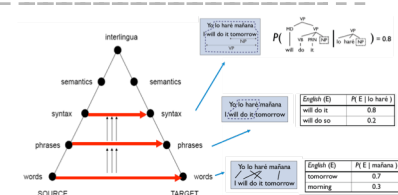
### Statistical Machine Translation (1990 - 2015)



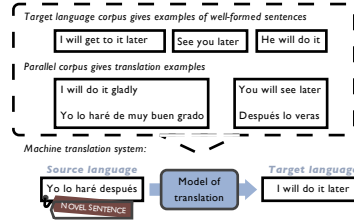
When I look at an article in Russian, I say: "This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode."

Warren Weaver (1949)

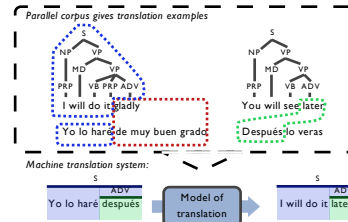
### Levels of Transfer: Vauquois Triangle (1968)



### Data-Driven Machine Translation



### Stitching Together Fragments



### Evolution of the Noisy Channel Model

$$P(e|f) \propto P(f|e) \cdot P(e)$$

$$P(e|f) \propto P(f|e)^{\phi_{lm}} \cdot P(e)^{\phi_{lm}}$$

$$P(e|f) \propto \exp \left\{ \sum_i w_i \cdot f_i(e, f) \right\}$$

Chosen to minimize loss

e.g.,  $\sum \log P(e)$

[illegible]