

Machine Translation



Dan Klein
UC Berkeley

Slides from John
DeNero

1

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 (but not much metadata).

2

Translation Examples

3

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	.		

4

Variety in Human-Generated Translations

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 astronomists are hard to find it for it comes from the direction of sun.

A volume enough to destroy a medium city small planet is big, flit earth 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://nmtatlas.fds.unimelb.edu.au/LC2000117>

5

Variety in Machine 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 astronomists are hard to find it for it comes from the direction of sun.

A Google Translate translation

A volume enough to destroy a medium city small planet is big, flit earth 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://nmtatlas.fds.unimelb.edu.au/LC2000117>

6

Evaluation

7

BLEU Score

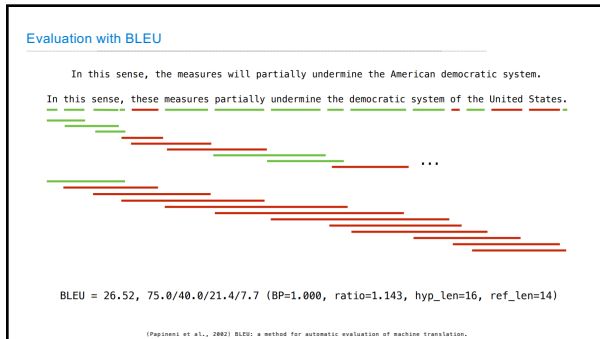
BLEU score: geometric mean of 1-, 2-, 3-, and 4-gram precision vs. a reference, multiplied by brevity penalty (harshly penalizes translations shorter than the reference).

$$\begin{aligned}
 \text{Matched}_i &= \sum_{t_i} \min \left\{ C_h(t_i), \max_j C_j(t_i) \right\} \\
 P_i &= \frac{\text{Matched}_i}{H_i} \\
 B &= \exp \left\{ \min \left(0, \frac{n-L}{n} \right) \right\} \\
 \text{BLUE} &= B \left(\prod_{i=1}^4 P_i \right)^{\frac{1}{4}}
 \end{aligned}$$

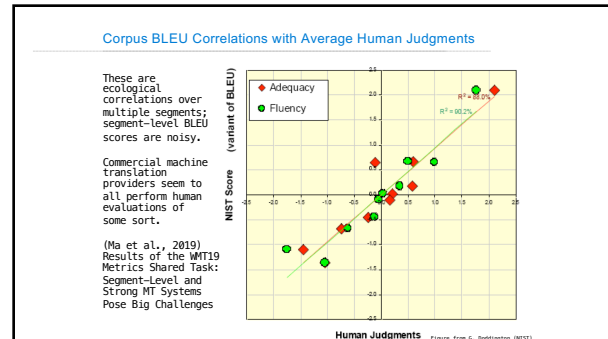
Annotations:

- If "of the" appears twice in hypothesis but only at most once in a reference, then only the first is "correct"
- "Clipped" precision of n-gram tokens
- Brevity penalty only matters if the hypothesis corpus is shorter than the shortest reference.
- BLUE is a mean of clipped precisions, scaled down by the brevity penalty.

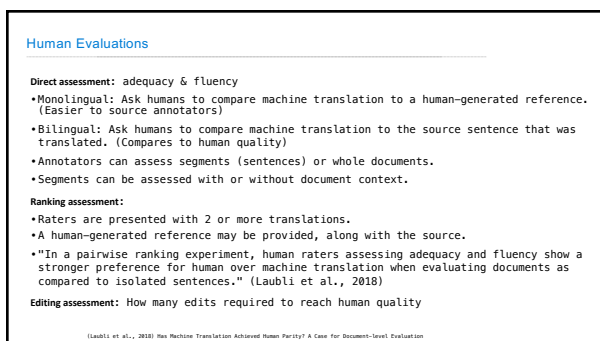
8



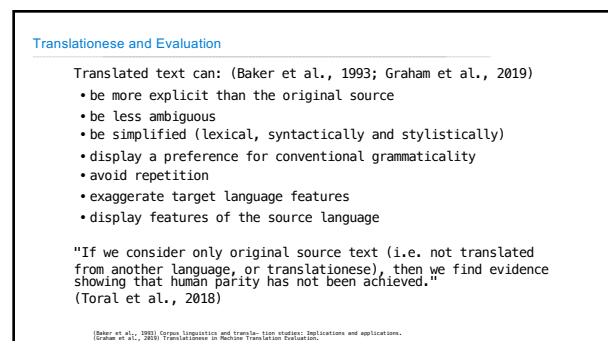
9



10



11



12

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 Chinese: all systems are outperformed by the human translator;
- × English to Czech: 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 Gujarati: all systems are outperformed by the human translator;
- × English to Kazakh: all systems are outperformed by the human translator;
- × English to Lithuanian: all systems are outperformed by the human translator;
- ✓ English to Russian: Facebook-FAIR is tied with human performance.

(Barrault et al., 2019). *Proceedings of the 2019 Conference on Machine Translation (WMT19)*.

13

Statistical Machine Translation
(1990 - 2015)

14

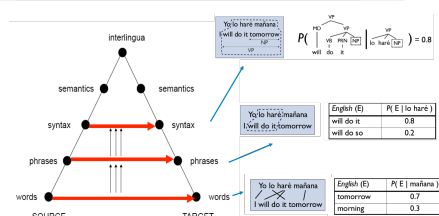


When I look at an article in Russian, I say: "This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode."

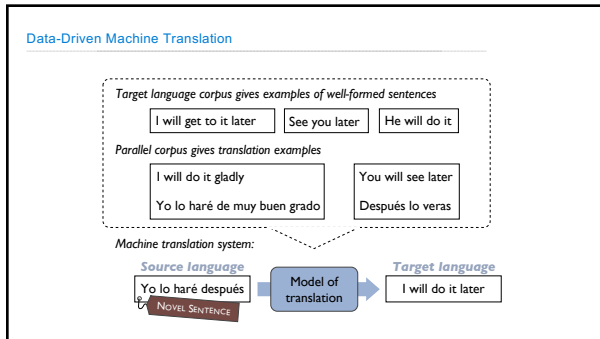
Warren Weaver (1949)

15

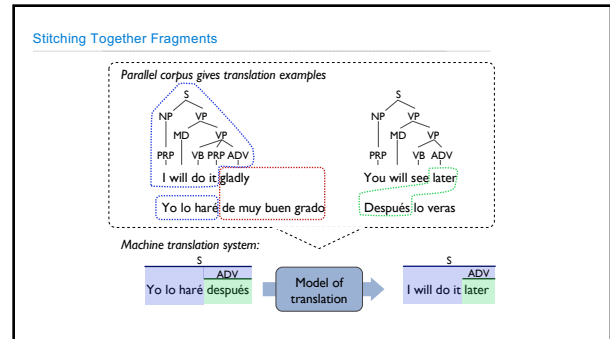
Levels of Transfer: Vauquois Triangle (1968)



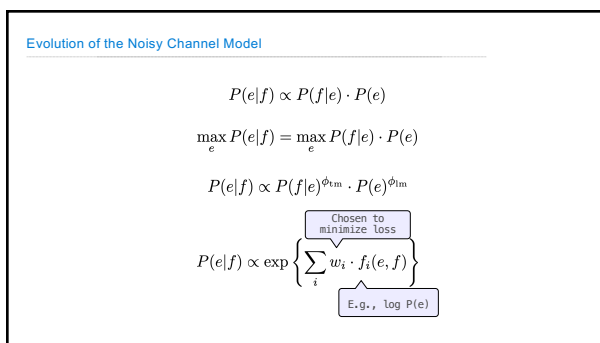
16



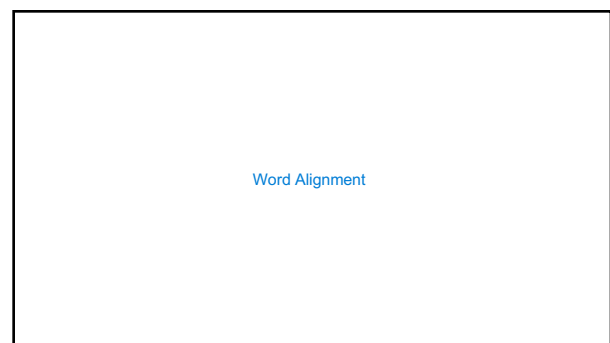
17



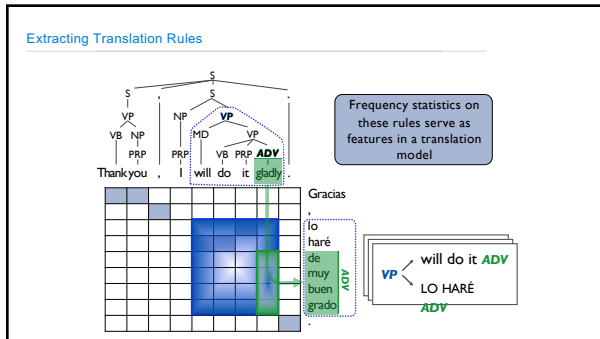
18



19



20



21

Counting Aligned Phrases

d'assister à la reunion et ||| to attend the meeting and
 assister à la reunion ||| attend the meeting
 la reunion and ||| the meeting and
 nous ||| we
 ...

- Relative frequencies are the most important features in a phrase-based or syntax-based model.
- Scoring a phrase under a lexical model is the second most important feature.
- Estimation does not involve choosing among segmentations of a sentence into phrases.

Slide by Greg Barrett

22

Interlude: Lexical Translation Models

23

HMM Alignment Model

24

Alignment Link Posteriors

$$c(e|f; \mathbf{e}, \mathbf{f}) = \sum_a \underbrace{p(a|\mathbf{e}, \mathbf{f}) \sum_{j=1}^{\ell_e} \delta(e, e_j) \delta(f, f_{a(j)})}_{\text{Non-zero for any alignment vector (for sentence pair } \mathbf{e}, \mathbf{f}) \text{ that has word } e \text{ aligned to word } f}}$$

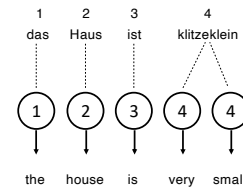
$$c(e|f; \mathbf{e}, \mathbf{f}) = \sum_i \sum_j \delta(e, e_j) \cdot \delta(f, f_i) \cdot P(a(j) = i | \mathbf{e}, \mathbf{f})$$

$$= \sum_i \sum_j \delta(e, e_j) \cdot \delta(f, f_i) \cdot \sum_a \underbrace{P(a|\mathbf{e}, \mathbf{f}) \cdot \delta(a(j), i)}_{\text{Non-zero for any alignment vector (for sentence pair } \mathbf{e}, \mathbf{f}) \text{ that has position } j \text{ aligned to position } i}}$$

25

Model 1 Posteriors

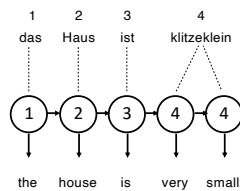
$$P(a(j) = i | \mathbf{e}, \mathbf{f}) = \frac{t(e_j | f_i)}{\sum_{i'} t(e_j | f_{i'})}$$



26

HMM Alignment Model

$$P(a, \mathbf{e} | \mathbf{f}) \propto \prod_j P(e_j | f_{a(j)}) \cdot P(a(j) | a(j-1))$$



(Egert, Stephan, Hermann Ney, and Christoph Eliezer, 2006) "HMM-based word alignment in statistical translation." In: "Proceedings of the 2006 Conference on Empirical Methods in Natural Language Processing (EMNLP 2006)".

27

HMM Alignment Model Posteriors

$$P(a(j) = i | \mathbf{e}, \mathbf{f}) = \sum_a \underbrace{P(a|\mathbf{e}, \mathbf{f}) \cdot \delta(a(j), i)}_{\text{Non-zero for alignments where } j \text{ is aligned to } i}}$$

$$= \sum_a \frac{P(a, \mathbf{e} | \mathbf{f}) \cdot \delta(a(j), i)}{P(\mathbf{e} | \mathbf{f})}$$

$$= \frac{\alpha_j(i) \cdot \beta_j(i)}{P(\mathbf{e} | \mathbf{f})} \quad \text{Forward-Backward algorithm}$$

$$\alpha_j(i) = \sum_{i'} P(a(j) = i | a(j-1) = i') \cdot P(e_j | f_i) \cdot \alpha_{j-1}(i')$$

$$\beta_j(i) = \sum_{i''} P(a(j+1) = i'' | a(j) = i) \cdot P(e_{j+1} | f_{i''}) \cdot \beta_{j+1}(i'')$$

$$\alpha_j(i) = P(e_1, e_2, \dots, e_j, a(j) = i | \mathbf{f})$$

$$\beta_j(i) = P(e_{j+1}, e_{j+2}, \dots, e_\ell | a(j) = i, \mathbf{f})$$

28

Interlude: Phrase-Based Models

29

What's Next?

Neural models: attention and the transformer architecture

Tricks of the trade: back-translation, knowledge distillation, subword models, and coverage vectors

30