

Natural Language Processing

Introduction to machine translation

Marco Kuhlmann

Department of Computer and Information Science

Machine translation

The screenshot shows the Google Translate interface in a browser window. The address bar displays 'translate.google.com'. The page title is 'Google Translate'. Below the title, there are two tabs: 'Text' (selected) and 'Documents'. The language selection bar shows 'ENGLISH - DETECTED' on the left and 'SWEDISH' on the right, with a bidirectional arrow between them. The source text in English is: 'Machine translation is the task of automatically translating text in one language (the source) into another language (the target)'. The translated text in Swedish is: 'Maskinöversättning är uppgiften att automatiskt översätta text på ett språk (källan) till ett annat språk (målet)'. The interface includes a speaker icon for audio playback, a character count '130/5000', and a 'Send feedback' link at the bottom right.

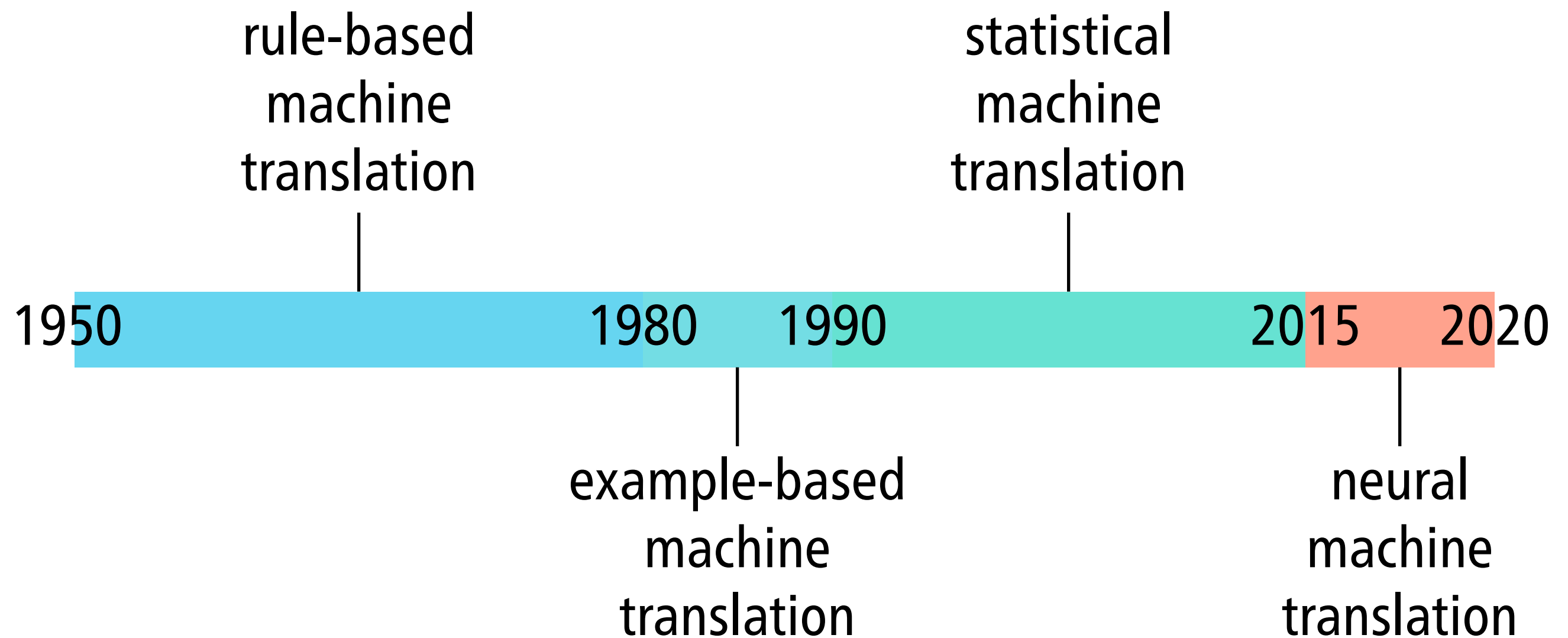
Machine translation is the task of automatically translating text in one language (the source) into another language (the target).

Maskinöversättning är uppgiften att automatiskt översätta text på ett språk (källan) till ett annat språk (målet).

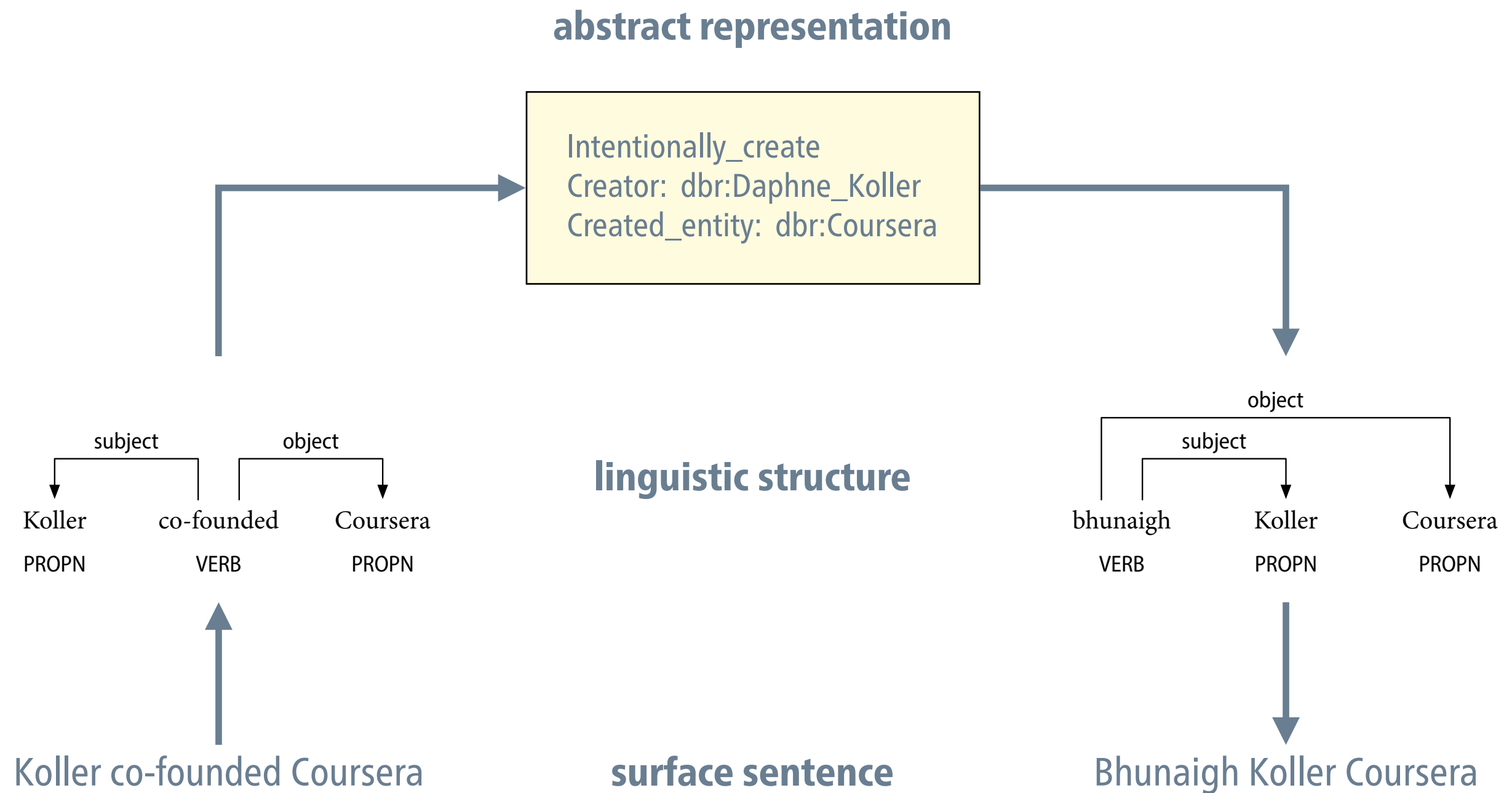
130/5000

Send feedback

A timeline of machine translation



Interlingual machine translation



Noisy Channel Model

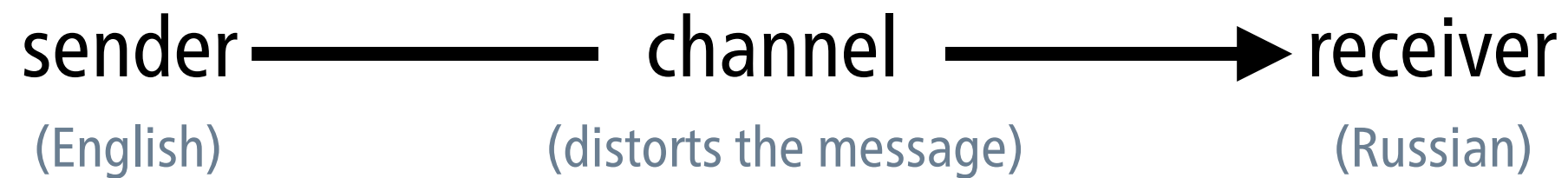


Image source



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.'

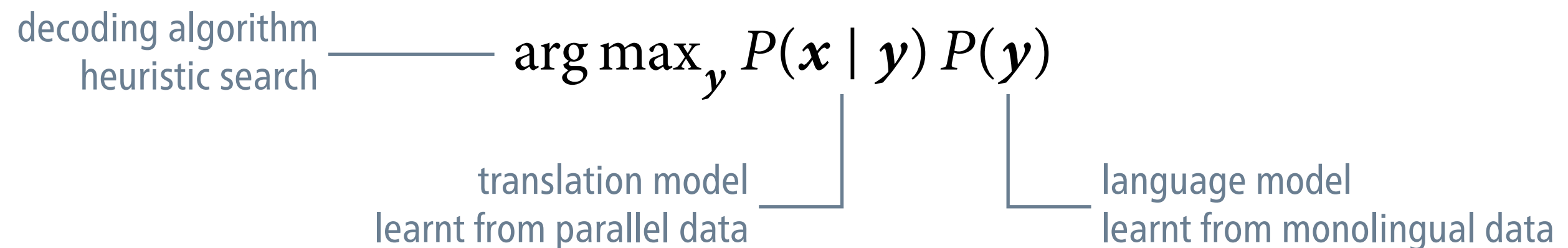
Warren Weaver (1894–1978)

Statistical machine translation (SMT)

- Formulate machine translation as an optimisation task: Given a source sentence \mathbf{x} , find the most probable target sentence \mathbf{y} :

$$\arg \max_{\mathbf{y}} P(\mathbf{y} | \mathbf{x})$$

- Use Bayes' rule to decompose the probability model into two components that can be learned separately:



Parallel corpora

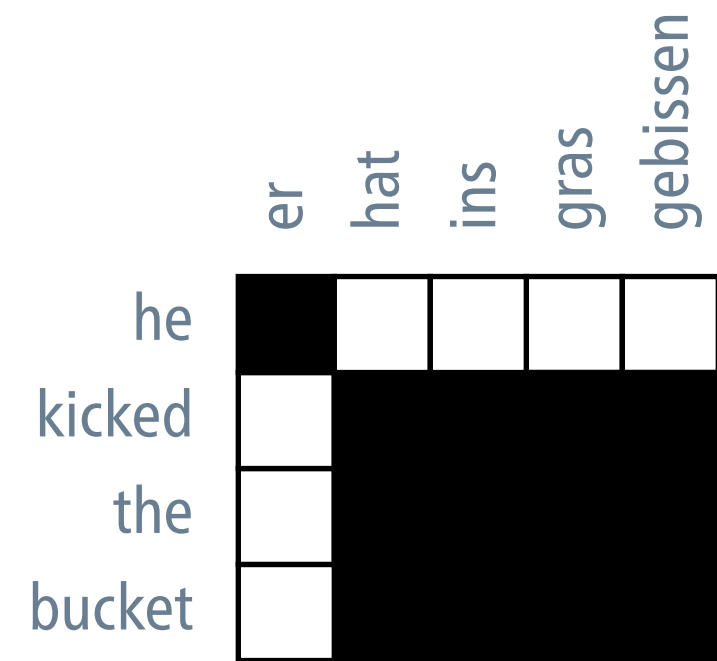
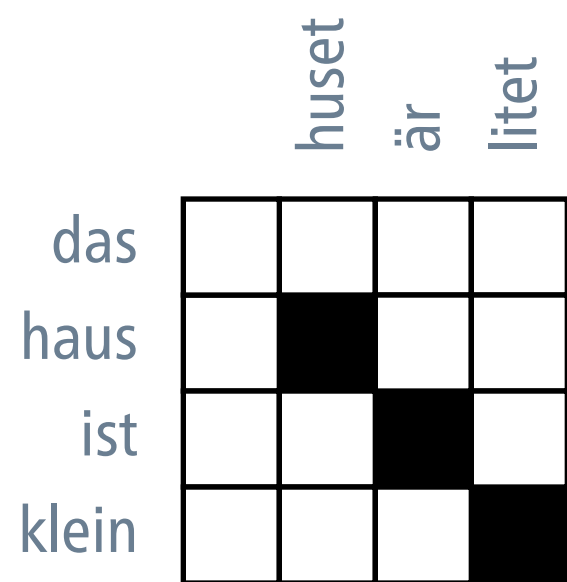
- Canadian Hansard (English–French); extracted from the proceedings of the Canadian Parliament.
- Europarl (21 languages); 30.32 M parallel sentences extracted from the proceedings of the European parliament.

[Link to the Europarl website](#)

- OPUS (several languages); growing collection of translated texts, automatically preprocessed and aligned.

[Link to the OPUS website](#)

Word-to-word alignments



IBM Model 1

$$P(\mathbf{x}, a \mid \mathbf{y}) = \frac{\varepsilon}{(|\mathbf{y}| + 1)^{|\mathbf{x}|}} \prod_{j=1}^{|\mathbf{x}|} t(x_j \mid y_{a(j)})$$

normalisation constant

|

ε

|

alignment function

- First in a series of increasingly complex statistical translation models; deals only with lexical (word-to-word) translation.
- Central component: The lexical translation probability t of observing a source word x , given the aligned target word y .

Training statistical machine translation models

- We would like to estimate the lexical translation probabilities from a parallel corpus – but we do not have the alignments.
- We can bootstrap the translation probabilities and alignments in parallel using the Expectation–Maximization (EM) algorithm:
 1. initialise the model parameters randomly
 2. calculate alignments based on the current model parameters
 3. estimate new model parameters from the new alignments
 4. repeat steps 2–3 until convergence

[Brown et al. \(1993\)](#)

Statistical machine translation (SMT)

- Research on statistical machine translation led to significant improvements in the availability and quality of translation.

The first version of Google Translate (2006–2016) was an SMT system.

- However, the best systems were extremely complex and required large amounts of external resources and feature engineering.

Open-source example: [Moses](#)

Evaluation: BLEU (Bilingual Evaluation Understudy)

- BLEU compares the automatic translation of a source sentence to one or several human-created reference translations.
- BLEU combines n -gram precision (for n up to 4) with a brevity penalty for too-short translations.
- BLEU has been criticised for not correlating well with human judgement, and several other evaluation measures exist.