

Word Alignment

Goal: Given a sentence in two languages, align words (tokens).

State of the art: GIZA++ (Och and Ney, 2000):

- Unsupervised, only sentence-parallel texts needed.
- Word alignments formally restricted to a **function**:

$$\text{src token} \mapsto \text{tgt token or NULL}$$

- A cascade of models refining the probability distribution:
 - IBM1: only lexical probabilities: $P(\textit{kočka} = \textit{cat})$
 - IBM2: absolute reordering added (not used in practice now)
 - IBM3: adds fertility: 1 word generates several others
 - IBM4/HMM: to account for relative reordering
- Only many-to-one links created \Rightarrow used twice, in both directions.