

Word Alignment



IT Systems Engineering | Universität Potsdam

Dr. Mariana Neves
(adapted from the original slides
of Prof. Philipp Koehn)

November 23rd, 2015

- Further discussion on word alignment, such as problems and quality measurement
- Present a method on word alignment based on the IBM models

Given a sentence pair, which words correspond to each other?

	michael	geht	davon	aus	,	dass	er	im	haus	bleibt
michael	■									
assumes		■	■	■						
that						■				
he							■			
will										■
stay										■
in								■		
the								■		
house									■	

- It does not need to be one-by-one.
- Words can have multiple or no alignment points.

	john	wohnt	hier	nicht
john	■			
does		?		?
not				■
live		■		
here			■	

Is the English word **does** aligned to the German **wohnt** (verb) or **nicht** (negation) or neither?

	john	biss	ins	grass
john	■			
kicked		■	■	■
the		■	■	■
bucket		■	■	■

How do the idioms **kicked the bucket** and **biss ins grass** match up?
Outside this exceptional context, **bucket** is never a good translation for **grass**

Word Alignment?

	john	biss	ins	grass
john	■			
kicked		■	■	■
the		■	■	■
bucket		■	■	■

The better solution here is a phrasal alignment!

- Sure alignments:
 - John to John
- Possible alignments:
 - kicked to biss
 - the to im
 - bucket to Grass

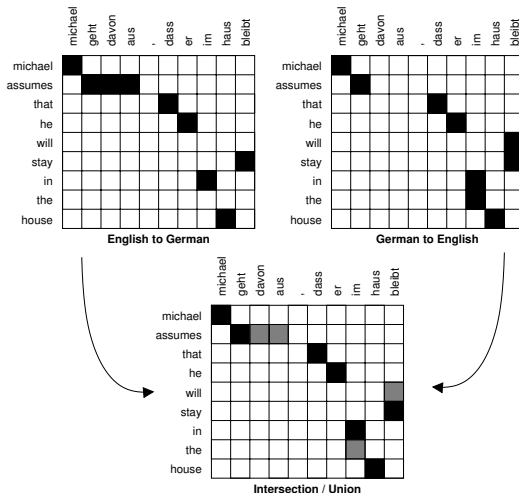
- Manually align corpus with *sure* (S) and *possible* (P) alignment points ($S \subseteq P$)
- Alignment Error Rate (AER): common metric for evaluation word alignments

$$\text{AER}(S, P; A) = 1 - \frac{|A \cap S| + |A \cap P|}{|A| + |S|}$$

- $\text{AER} = 0$: alignment A matches all sure, any possible alignment points

- IBM Models create a **many-to-one** mapping
 - words are aligned using an alignment function
 - a function may return the same value for different input (one-to-many mapping)
 - a function can not return multiple values for one input (no many-to-one mapping)
- Real word alignments have **many-to-many** mappings

Symmetrizing Word Alignments



- Intersection of GIZA++ bidirectional alignments

- The **intersection** usually contains good alignment points (high precision), but not all of them.
- The **union** usually contains most of the desired align points (high recall), but also faulty points.

- We want to explore the space between the two extremes:
 - Take the all alignment points in the intersection (reliable).
 - Add some of the points from the union (neighboring candidates), incrementally.

grow-diag-final(e2f,f2e)

- 1: neighboring = $\{(-1,0),(0,-1),(1,0),(0,1),(-1,-1),(-1,1),(1,-1),(1,1)\}$
- 2: alignment $A = \text{intersect}(e2f,f2e)$; **grow-diag**(); **final**(e2f); **final**(f2e);

grow-diag()

- 1: **while** new points added **do**
- 2: **for all** English word $e \in [1\dots e_n]$, foreign word $f \in [1\dots f_n]$, $(e, f) \in A$ **do**
- 3: **for all** neighboring alignment points $(e_{\text{new}}, f_{\text{new}})$ **do**
- 4: **if** $(e_{\text{new}}$ unaligned OR f_{new} unaligned) AND $(e_{\text{new}}, f_{\text{new}}) \in \text{union}(e2f,f2e)$ **then**
- 5: add $(e_{\text{new}}, f_{\text{new}})$ to A
- 6: **end if**
- 7: **end for**
- 8: **end for**
- 9: **end while**

final()

- 1: **for all** English word $e_{\text{new}} \in [1\dots e_n]$, foreign word $f_{\text{new}} \in [1\dots f_n]$ **do**
- 2: **if** $(e_{\text{new}}$ unaligned OR f_{new} unaligned) AND $(e_{\text{new}}, f_{\text{new}}) \in \text{union}(e2f,f2e)$ **then**
- 3: add $(e_{\text{new}}, f_{\text{new}})$ to A
- 4: **end if**
- 5: **end for**

- Statistical Machine Translation, Philipp Koehn (section 4.5).