

Concept Alignment for Multilingual Machine Translation

04.07.2021

Arianna Masciolini

Context

- ❖ GF is well suited for **domain-specific MT** systems, where precision is more important than coverage, as it provides strong guarantees of grammatical correctness
- ❖ in such systems, **lexical exactness** is as important as grammaticality
 - ❖ need for high-quality **translation lexica** preserving semantics *and* morphological correctness

The problem

- ❖ manually building a translation lexicon
 - ❖ is time consuming
 - ❖ requires significant linguistic knowledge
- ❖ desire to **automate** this process at least in part
 - ❖ possible when **example parallel data** are available

A parallel corpus

Alice thought she might as well wait, as she had nothing else to do, and perhaps after all it might tell her something worth hearing.

For some minutes it puffed away without speaking, but at last it unfolded its arms, took the hookah out of its mouth again, and said, 'So you think you're changed, do you?'

'I'm afraid I am, sir,' said Alice; 'I can't remember things as I used--and I don't keep the same size for ten minutes together!'

Alice pensò che poteva aspettare, perchè non aveva niente di meglio da fare, e perchè forse il Bruco avrebbe potuto dirle qualche cosa d'importante.

Per qualche istante il Bruco fumò in silenzio, finalmente sciolse le braccia, si tolse la pipa di bocca e disse: — E così, tu credi di essere cambiata?

— Ho paura di sì, signore, — rispose Alice. — Non posso ricordarmi le cose bene come una volta, e non rimango della stessa statura neppure per lo spazio di dieci minuti!

From Lewis Carroll, *Alice's adventures in Wonderland*. Parallel text at paralleltext.io

Alignment

Word alignment:

Alice thought she might as well wait, as she had
nothing else to do, and perhaps after all it might tell
her something worth hearing.

Alice pensò che poteva aspettare, perchè non aveva
niente di meglio da fare, e perchè forse il Bruco
avrebbe potuto dirle qualche cosa d'importante.

Phrase alignment:

Alice thought she might as well wait, as she had
nothing else to do, and perhaps after all it might tell
her something worth hearing.

Alice pensò che poteva aspettare, perchè non aveva
niente di meglio da fare, e perchè forse il Bruco
avrebbe potuto dirle qualche cosa d'importante.

Statistical approaches

Standard approaches are statistical (IBM models).

❖ **pros:**

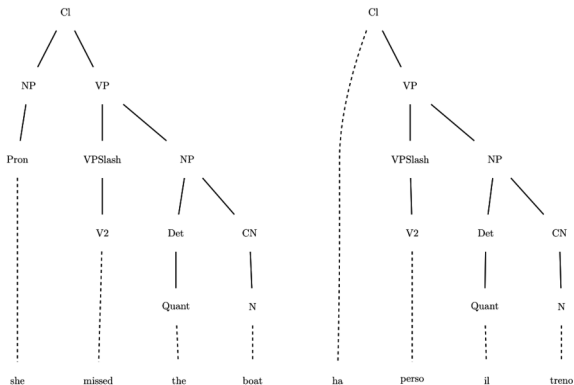
- ❖ easy to use
- ❖ can handle noisy data
- ❖ fast on large corpora

❖ **cons:**

- ❖ *require* large amounts of raw data
- ❖ correspondences between strings → no morphological info
- ❖ “fixed” level of abstraction (word, phrase or sentence)

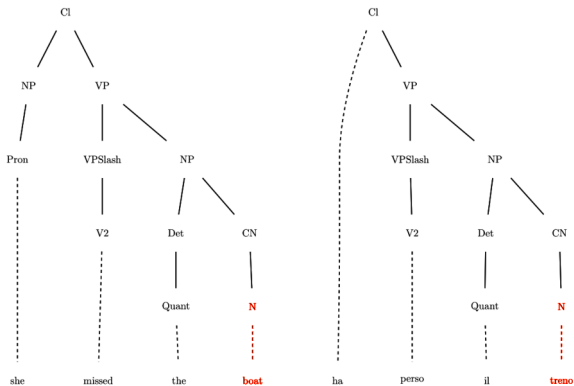
Syntax-based approaches I

Alternative: tree-to-tree alignment.



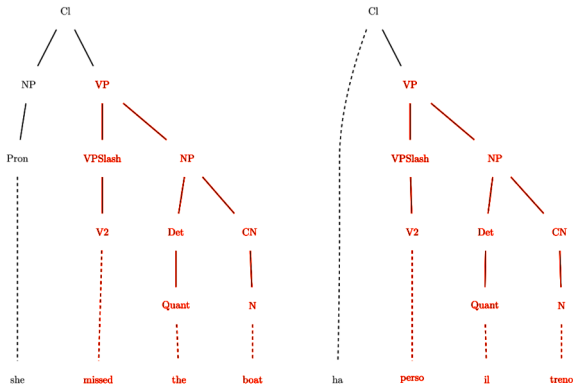
Syntax-based approaches II

Word alignment



Syntax-based approaches III

Phrase alignment



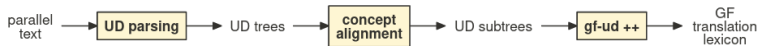
Comparison

statistical	syntax-based
require large amounts of data	work consistently well even on individual sentence pairs
works with raw data	requires the data to be analyzed
correspondences between strings	correspondences between grammatical objects
“fixed” level of abstraction (word or phrase)	all levels of abstraction → concept alignment

Why not just use GF?

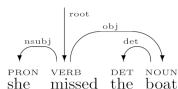
- ❖ quality of the analysis is crucial
 - ❖ lack of robust GF parsers
- ❖ dependency trees are an easier target for a parser
 - ❖ neural parsers such as **UDPipe**

Overview



1. parse parallel data to UD trees
2. search for aligned UD subtrees
3. convert them to GF trees and then grammar rules

UD trees



text = she missed the boat

1 she she PRON _ _ 2 nsubj _ _

2 missed miss VERB _ _ 0 root _ _

3 the the DET _ _ 4 det _ _

4 boat boat NOUN _ _ 2 obj _ _

2 missed miss VERB _ _ 0 root _ _

1 she she PRON _ _ 2 nsubj _ _

4 boat boat NOUN _ _ 2 obj _ _

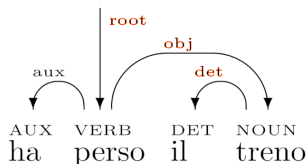
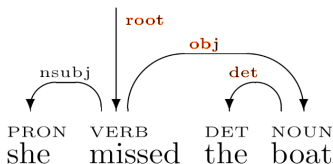
3 the the DET _ _ 4 det _ _

Graphical, CoNNL-U and Rose Tree representation of the same UD tree.

- ❑ dependency-labelled links between words (head-dependent pairs)
- ❑ POS tags
- ❑ ...

Extracting concepts

Matching dependency labels

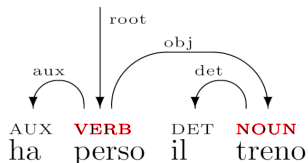
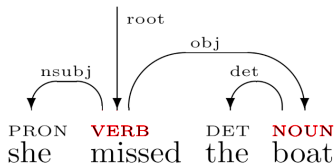


- ❑ $\langle \textit{she missed the boat, ha perso il treno} \rangle$
- ❑ $\langle \textit{missed the boat, perso il treno} \rangle$
- ❑ $*\langle \textit{the boat, il treno} \rangle$
- ❑ $\langle \textit{the, il} \rangle$

Aligning heads of matching trees

- ❑ $\langle \textit{the boat, il treno} \rangle \rightarrow * \langle \textit{boat, treno} \rangle$
- ❑ $\langle \textit{missed the boat, perso il treno} \rangle \rightarrow \langle \textit{missed, ha perso} \rangle$
(including the auxiliary)

Using POS tags



- ❖ more reliable **ignoring function words**
- ❖ in this case, basically same results as when matching labels
- ❖ can increase recall when labels do not coincide
- ❖ can increase precision if used **in conjunction with labels**

Translation divergences

Divergence: systematic cross-linguistic distinction.

- ❖ categorial
 - ❖ ⟨*Gioara listens **distractedly**, Gioara lyssnar **distraherad***⟩
 - ❖ ⟨*Herbert completed his **doctoral** thesis, Herbert ha completato la sua tesi **di dottorato***⟩
- ❖ conflational
 - ❖ ⟨*Filippo is interested in **game development**, Filippo är intresserad av **spelutveckling***⟩
- ❖ structural
 - ❖ ⟨*I called **Francesco**, Ho telefonato a **Francesco***⟩
- ❖ head swapping
 - ❖ ⟨*Anna **usually** goes for walks, Anna **brukar** promenera*⟩
- ❖ thematic
 - ❖ ⟨***Yana** likes **books**, A **Yana** piacciono **i libri***⟩

Reusing known alignments

- allows using CA in conjunction with statistical tools
- iterative application

Searching for specific patterns

- ❑ gf-ud pattern matching allows looking for specific syntactic patterns
- ❑ possible generalization via pattern replacement

Example predication patterns:

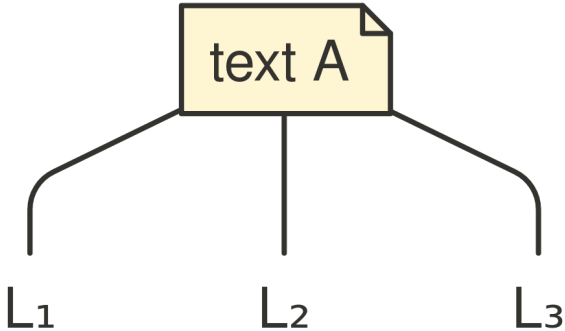
- ❑ $\langle \textit{she missed the boat, ha perso il treno} \rangle \rightarrow \langle [\textit{subj}] \textit{ missed} [\textit{obj}], \textit{ ha perso} [\textit{obj}] \rangle$
- ❑ $\langle \textit{she told you that, hon berättade det för dig} \rangle \rightarrow \langle [\textit{subj}] \textit{ told} [\textit{iobj}] [\textit{obj}], [\textit{subj}] \textit{ berättade} [\textit{obj}] \textit{ för} [\textit{obl}] \rangle$

Propagating concepts to a new language

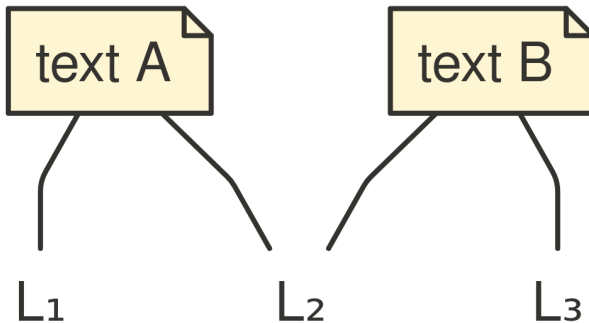
Concept Propagation

- ❖ So far, we focused on how to identify correspondences in bilingual parallel texts (***Concept Extraction***)
- ❖ what happens when we need to handle a third language?
 - ❖ ***Concept Propagation***: finding the expression corresponding to a known concept in a new language

Scenario 1



Scenario 2



Generating grammar rules

Requirements

- ❑ aligned UD trees
- ❑ dependency configurations for gf-ud
- ❑ **morphological dictionaries**
- ❑ **extraction grammar**

Morphological dictionaries

Purely morphological unilingual dictionaries.

Example:

```
...  
lin morphologic_A =  
    mkAMost "morphologic" "morphologicly" ;  
lin morphological_A =  
    mkAMost "morphological" "morphologically" ;  
lin morphology_N =  
    mkN "morphology" "morphologies" ;  
...
```

Extraction grammar

Defines the syntactic categories and functions to build lexical entries.

Example (prepositional NPs):

PrepNP : Prep -> NP -> PP # case head

Lexical rules

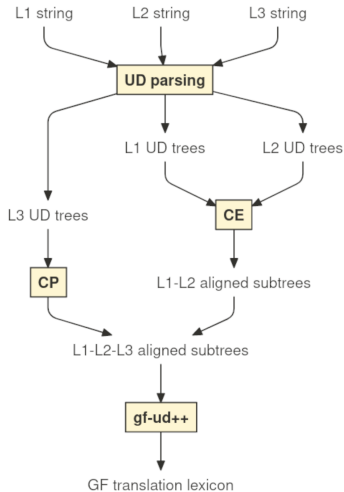
Abstract:

```
fun in_the_field__inom_området_PP : PP ;
```

English concrete:

```
lin in_the_field__inom_område_PP =  
  PrepNP in_Prep (DetCN the_Det (UseN field_N))
```

Detailed view



Refining the generated lexicon

Postprocessing tools:

- ▣ interactive selection
- ▣ CoNNL-U synoptic viewer

Summary

- ❑ (parsing)
- ❑ Concept Extraction
- ❑ Concept Propagation
- ❑ GF lexicon generation
- ❑ postprocessing

Links

Links to everything mentioned in this talk, and more:

- ❑ **An overview of the IBM models**
- ❑ `fast_align`
- ❑ **the UD standard**
- ❑ **UDPipe**
- ❑ **B. J. Dorr's paper on translation divergences**
- ❑ `gf-ud`
- ❑ **the concept-alignment repo**
- ❑ **my thesis report** where everything is explained in detail but not everything is up to date
- ❑ **the paper on CE** I wrote together with Aarne
- ❑ **the CoNNL-U synoptic viewer**

Questions?