

Automatic Induction of Transfer Rules using Aligned Bilingual Corpora

Yvette Graham

National Center for Language Technology
School of Computing
Dublin City University

May 31 2006

BiText Project

Project Aims

- ▶ Automatically generate a number of aligned bilingual corpora (bi-texts)
- ▶ Acquire existing bilingual corpora.
- ▶ Automatically annotate the bi-texts with LFG c-structures and f-structures.
- ▶ Automatically induce transfer rules from the corpora.

Traditionally, transfer rules for machine translation have been hand-coded - takes time and resources.

Automatically inducing such transfer rules from large bilingual corpora could provide an efficient way of acquiring such transfer rules.

Transfer Rules for MT

SL sentence parsed and f-structure achieved.

Transfer Rule database is searched for a matching rule

- ▶ The SL sentence f-structure must match the LHS of the transfer rule.
- ▶ If multiple rules match, then the most specific and most likely one is selected.

Recursively apply transfer rules to any parts of the f-structure not taken care of by the initial rule.

Perhaps adding a default rule, for cases in which there is no matching transfer rule for an given f-structure.

Adding this information to the RHS of the initial rule produces the TL f-structure.

This f-structure is then used to generate the TL sentence.

Some advantages of this approach to MT

The f-structure notation allows us to infer desirable transfer rules for f-structure relations that correspond to non-adjacent words and thus would not be discovered in the string-based world.

For example:

▶ **Es gibt ...**
There is ...
(lit. *It gives ...*)

▶ **Es scheint ... zu geben.**
There seems [...] to be.
(lit. *It seems ... to give.*)

- - Riezler, Maxwell (2006)

Transfer rules allow for mapping between different linguistic types.

For example, extraction of the phrase pair:

▶ **zutiefst dankbar** ⇒ **a deep appreciation**

Filtering of f-structure phrases based on consistency of linguistic types would find this mapping invalid because of the incompatibility of types A and N for adjective *dankbar* and nominal *appreciation*.

- - Riezler, Maxwell (2006)

Automatically Inducing Transfer Rules

- ▶ F-structure annotated bilingual corpora are used to extract transfer rules for a given language pair.
- ▶ Relevant sentences and their translations are extracted from such corpora and analysed.
- ▶ The similarities between the SL f-structures are examined
⇒ generalisation about the structure of *the sentence to be translated* and form the LHS of the transfer rule.
- ▶ The similarities amongst the TL f-structures are examined
⇒ a generalisation about the structure of *the translation of the sentence* and forms the RHS of the rule.
- ▶ This results in a transfer rule meaning:
given a SL sentence whose f-structure matches the LHS of the rule, the f-structure of its translation is likely to be that of the RHS of the rule.

Example: Automatically Inducing Transfer Rule for the verb “*talk*” from English to French

- ▶ Search the English side of all available English/French bi-text corpora for sentences in which the main predicate is *talk*.
- ▶ Retrieve both the English sentences and their matching French sentences.
- ▶ Examine the f-structures of the sentences looking for regularities amongst the English f-structures and then similarities amongst the f-structures of their translations.
- ▶ For Example: the search of the corpora could return the following three sentence pairs
 1. The man talks to the dog. *L'homme parle au chien.*
 2. He talks to Mary all day. *Il parle à Marie toute la journée.*
 3. The president talks to the nurse. *Le président parle à l'infirmière.*

F-Structures for Sentence Pair 1

The man talks to the dog.

$$\left[\begin{array}{l} \text{SUBJ } 1: \left[\begin{array}{l} \text{PRED 'man'} \\ \text{DEF +} \\ \text{NUM sg} \end{array} \right] \\ \text{PRED 'talk(SUBJ, OBL OBJ)'} \\ \text{OBL } 2: \left[\begin{array}{l} \text{PRED 'to(OBJ)'} \\ \text{OBJ } 3: \left[\begin{array}{l} \text{PRED 'dog'} \\ \text{NUM sg} \\ \text{DEF +} \end{array} \right] \end{array} \right] \end{array} \right]$$

L'homme parle au chien.

$$\left[\begin{array}{l} \text{SUBJ } 1: \left[\begin{array}{l} \text{PRED 'homme'} \\ \text{DEF +} \\ \text{NUM sg} \\ \text{GEND masc} \end{array} \right] \\ \text{PRED 'parler(SUBJ, OBJ)'} \\ \text{OBJ } 2: \left[\begin{array}{l} \text{PRED 'chien'} \\ \text{DEF +} \\ \text{NUM sg} \\ \text{GEND masc} \\ \text{PFORM à_} \end{array} \right] \end{array} \right]$$

F-Structures as Equations: Sentence Pair 1

English F-Structure:

{
SUBJ (0, 1)
PRED (0, 'talk<SUBJ, OBL OBJ>')
OBL (0, 2)
PRED (1, 'man')
DEF (1, +)
NUM (1, sg)
PRED (2, 'to<OBJ>')
OBJ (2, 3)
PRED (3, 'dog')
NUM (3, SG)
DEF (3, +)
}

French F-Structure:

{
SUBJ (0, 1)
PRED (0, 'parler<SUBJ, OBJ>')
OBJ (0, 2)
PRED (1, 'homme')
DEF (1, +)
NUM (1, sg)
GEND (1, masc)
PRED (2, 'chien')
DEF (2, +)
NUM (2, sg)
GEND (2, MASC)
PFORM (2, à_)
}

Linking the Content of the Two F-Structures

- ▶ Examine the values of the predicates of the English f-structure
- ▶ Translate each of the values using a bilingual dictionary or word list
- ▶ Search for the translated predicates in the French F-structure
- ▶ Record links between the f-structures

1. Predicate values of the english f-structure

$$\left\{ \begin{array}{l} \text{PRED (0, 'talk(SUBJ, OBL OBJ)')} \\ \text{PRED (1, 'man')} \\ \text{PRED (2, 'to(OBJ)')} \\ \text{PRED (3, 'dog')} \end{array} \right\}$$

2. Dictionary Look-Up

talk \Rightarrow parler

man \Rightarrow homme

to \Rightarrow à

dog \Rightarrow chien

3. Translations in the french f-structure predicate values:

$$\left\{ \begin{array}{l} \text{PRED (0, 'parler(SUBJ, OBJ)')} \\ \text{PRED (1, 'homme')} \\ \text{PRED (2, 'chien')} \end{array} \right\}$$

4. Resulting links between f-structures:

0 \Rightarrow 0 write as variable A

1 \Rightarrow 1 write as variable B

3 \Rightarrow 2 write as variable C

Linking the Content of the F-Structures

English F-Structure:

$$\left\{ \begin{array}{l} \text{SUBJ (A, B)} \\ \text{PRED (A, 'talk(SUBJ, OBL OBJ)')} \\ \text{OBL (A, 2)} \\ \text{PRED (B, 'man')} \\ \text{DEF (B, +)} \\ \text{NUM (B, sg)} \\ \text{PRED (2, 'to(OBJ)')} \\ \text{OBJ (2, C)} \\ \text{PRED (C, 'dog')} \\ \text{NUM (C, sg)} \\ \text{DEF (C, +)} \end{array} \right\}$$

\Rightarrow

French F-Structure:

$$\left\{ \begin{array}{l} \text{SUBJ (A, B)} \\ \text{PRED (A, 'parler(SUBJ, OBJ)')} \\ \text{OBJ (A, C)} \\ \text{PRED (B, 'homme')} \\ \text{DEF (B, +)} \\ \text{NUM (B, sg)} \\ \text{GEND (B, masc)} \\ \text{PRED (C, 'chien')} \\ \text{DEF (C, +)} \\ \text{NUM (C, sg)} \\ \text{GEND (C, MASC)} \\ \text{PFORM (C, à-)} \end{array} \right\}$$

Examine the left side of the rules

LHS for Sentence Pair 2:

LHS for Sentence Pair 1:

SUBJ (A, B)
PRED (A, 'talk(SUBJ, OBL OBJ)')
OBL (A ,2)
PRED (B, 'man')
DEF (B, +)
NUM (B, sg)
PRED (2, 'to(OBJ)')
OBJ (2, C)
PRED (C, 'dog')
NUM (C, SG)
DEF (C, +)

SUBJ (A, B)
PRED (A, 'talk(SUBJ, OBL OBJ)')
OBL (A ,2)
ADJUNCT (A, D)
PRED (B, 'pro')
CASE (B, nom)
GEND (B, male)
NUM (B, sg)
PRED (2, 'to(OBJ)')
OBJ (2, C)
PRED (C, 'Mary')
PRED (D, 'day')
SPEC (D, 5)
QUANT (5, E)
PRED (E, 'all')

LHS for Sentence Pair 3:

SUBJ (A, B)
PRED (A, 'talk(SUBJ, OBL OBJ)')
OBL (A ,2)
PRED (B, 'president')
DEF (B, +)
NUM (B, sg)
PRED (2, 'to(OBJ)')
OBJ (2, C)
PRED (C, 'nurse')
NUM (C, SG)
DEF (C, +)

The common equations among the set of f-structures gives us the LHS of the transfer rule:

SUBJ (A, B)
PRED (A, 'talk(SUBJ, OBL OBJ)')
OBL (A ,2)
PRED (2, 'to(OBJ)')
OBJ (2, C)

Examine the right side of the rules

RHS for Sentence Pair 1:

subj (A, B)
pred (A, 'parler<subj, obj>')
obj (A, C)
PRED (B, 'homme')
DEF (B, +)
NUM (B, sg)
GEN (B, masc)
PRED (C, 'chien')
DEF (C, +)
NUM (C, sg)
GEN (C, MASC)
PFORM (C, à_)

RHS for Sentence Pair 2:

subj (A, B)
pred (A, 'parler<subj, obj>')
obj (A, C)
ADJUNCT (A, D)
PRED (B, 'pro')
CASE (B, nom)
NUM (B, sg)
GEN (B, masc)
PRED (C, 'Marie')
PRED (D, 'journee')
QUANT (D, E)
SPEC (D, 5)
PRED (E, 'tout')
DET (5, 6)
PRED (6, 'le')

RHS for Sentence Pair 3:

subj (A, B)
pred (A, 'parler<subj, obj>')
obj (A, C)
PRED (B, 'president')
DEF (B, +)
NUM (B, sg)
GEN (B, masc)
PRED (C, 'infirmiere')
DEF (C, +)
NUM (C, sg)
GEN (C, FEM)
PFORM (C, à_)

The common equations among the set of f-structures gives us the RHS of the transfer rule:

subj (A, B)
pred (A, 'parler<subj, obj>')
obj (A, C)

Resulting Transfer Rule

$$\left\{ \begin{array}{l} \text{SUBJ (A, B)} \\ \text{PRED (A, 'talk<SUBJ, OBL OBJ>')} \\ \text{OBL (A ,2)} \\ \text{PRED (2, 'to<OBJ>')} \\ \text{OBJ (2, C)} \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} \text{SUBJ (A, B)} \\ \text{PRED (A, 'parler<SUBJ, OBJ>')} \\ \text{OBJ (A, C)} \end{array} \right\}$$

Finally, replacing any remaining f-structure labels with variables:

$$\left\{ \begin{array}{l} \text{SUBJ (A, B)} \\ \text{PRED (A, 'talk<SUBJ, OBL OBJ>')} \\ \text{OBL (A ,D)} \\ \text{PRED (D, 'to<OBJ>')} \\ \text{OBJ (D, C)} \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} \text{SUBJ (A, B)} \\ \text{PRED (A, 'parler<SUBJ, OBJ>')} \\ \text{OBJ (A, C)} \end{array} \right\}$$

Current Project Status

Work to date:

- ▶ Compiled a list of available bilingual corpora.
- ▶ Developed a small trilingual corpus of sentences in German, English and French.
- ▶ Used XLE to parse the sentences producing their f-structures in prolog format.
- ▶ Manually compiled word-lists for the example sentences.

Currently working on:

- ▶ Implementing and testing different algorithms to extract transfer rules from the f-structures of aligned bilingual sentences for a given language pair.

Some Future Work

Further develop the induction algorithm

Use the Human Centre Corpus to induce transfer rules and as validation for the induction algorithm

- ▶ French English (900 sent. approx)

Scale things up further to use a larger f-structure annotated corpus

Use the induced grammars developed by members of the LFG group to annotate existing bi-texts

Use these annotated bitexts to induce transfer rules