# Automatic Post-Editing for the DiscoMT Pronoun Translation Task

**Liane Guillou**

School of Informatics
University of Edinburgh
Scotland, United Kingdom
L.K.Guillou@sms.ed.ac.uk

## Abstract

This paper describes an automated post-editing submission to the DiscoMT 2015 shared task on pronoun translation. Post-editing is achieved by applying pronoun-specific rules to the output of an English-to-French phrase-based SMT system.

## 1 Introduction

The shared task (Hardmeier et al., 2015) focusses on the translation of the English pronouns "it" and "they" into French. While they both serve multiple functions in English, the most significant is as *anaphoric* pronouns, referring back to an entity previously mentioned in the discourse, known as the *antecedent*.

When translated into French, anaphoric pronouns must agree with their antecedent in terms of both number and grammatical gender. Therefore, selecting the correct pronoun in French relies on knowing the number and gender of the antecedent. This presents a problem for current state-of-the-art Statistical Machine Translation (SMT) systems which translate sentences in isolation.

*Inter-sentential* anaphoric pronouns, i.e. those that occur in a different sentence to their antecedent, will be translated with no knowledge of their antecedent. Pronoun-antecedent agreement therefore cannot be guaranteed. Even *intra-sentential* pronouns, i.e. those that occur in the same sentence as their antecedent, may lack sufficient local context to ensure agreement.

The English pronoun "it" may also be used as a pleonastic or event pronoun. *Pleonastic* pronouns such as the "it" in "**it** is raining" or the "il" in "**il** pleut" do not refer to anything but are required by syntax to fill the subject-position slot. *Event* pronouns may refer to a verb, verb phrase or even an entire clause or sentence. The pronoun "they" may also serve as a *generic* pronoun, as in "**They** say

it always rains in Scotland" – here "they" does not refer to a specific person or group. For each pronoun type, translations into French must meet different requirements.

This paper presents an automatic post-editing approach which applies two pronoun-specific rules to the output of an English-to-French phrase-based SMT system. One rule handles anaphoric pronouns and the other handles non-anaphoric (i.e. event and pleonastic) pronouns.

The advantage of a post-editing approach is that the translations of both pronouns and their antecedents (for anaphoric pronouns) are already known. There is therefore no need to keep track of this information within the decoder. Instead, the problem becomes one of identifying incorrectly translated pronouns and amending them based on information extracted from the source-language text. The aim is to leverage knowledge about the target-language and through this maximise the number of changes that will improve the pronoun translations, whilst also attempting to minimise those that may have a detrimental effect.

The post-editing rules make use of information automatically obtained from the source-language text. The risk of doing this is that inaccurate information could lead to incorrect translations. As post-editing takes place after translation, the decoder and language model can no longer be relied upon to recover from bad decisions. However, due to the simplicity of the approach and encouraging results from Weiner (2014) for the English-German pair, post-editing is worth exploring.

## 2 Post-editing Overview

Using the ParCor corpus (Guillou et al., 2014) annotations as a model, automated tools are applied to the full text of each (sentence-split) source-language document in the dataset to extract the following information: anaphoric vs. non-anaphoric pronouns, subject vs. object position and the an-
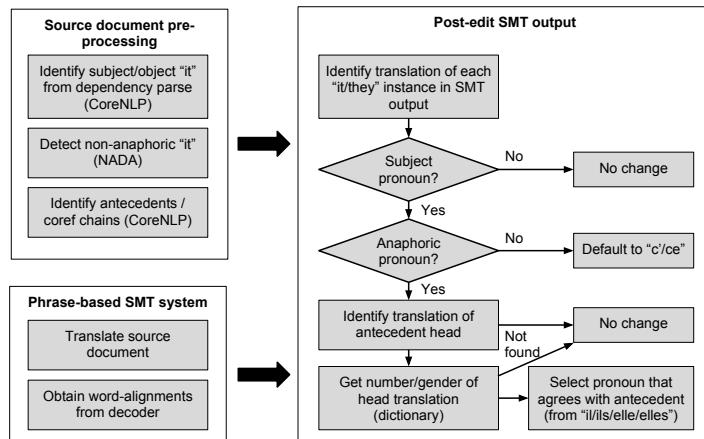
Figure 1: The post-editing process

| Data | Description | Parallel Sentences | Monolingual Sentences |
|---|---|---|---|
| Training | TED, Europarl, News Commentary | 2,372,666 | |
| Tuning | dev2010 + tst2011 | 1,705 | |
| Development test | tst2010 | 1,664 | |
| Development test | tst2012 | 1,124 | |
| Language model | TED, Europarl, News Commentary and News | | 33,869,133 |

Table 1: Baseline training, tuning and development data.

tecedent of each anaphoric pronoun. This information is then leveraged by two post-editing rules; one for anaphoric pronouns and one for non-anaphoric pronouns. These rules are automatically applied to the 1-best output of the baseline SMT system described in Section 3. The process for extracting source-language information and application of the post-editing rules is outlined in Figure 1 and described in Sections 4 and 5.

## 3 Baseline Machine Translation System

The baseline system used to produce the SMT output is of a similar design to that provided as part of the shared task resources. It is a phrase-based system built using the Moses toolkit (Koehn et al., 2007) and trained/tuned using only the pre-processed (tokenised, lower-cased) parallel data provided for the shared task. Training, tuning and (development) test data are described in Table 1.

Word alignments are computed using Giza++ with *grow-diag-final-and* symmetrization, and with sentences restricted to 80 tokens or fewer (as Giza++ produces more robust alignments for shorter sentences). The maximum phrase length is set to 7. As memory and disk space are not a concern, sig-test filtering which prunes unlikely phrase pairs from the phrase table, is not used in

training the baseline system. Tuning is performed using MERT (Och, 2003) with an N-best list of 200, and using the dev2010+tst2011 data.

The language model is a 5-gram KenLM (Heafield, 2011) model, trained using lmplz, with modified Kneser-Ney smoothing and no pruning. The memory optimisations that were made for the shared task baseline[1] are not replicated as they are not required. The language model uses the *probing data structure*; the fastest and default data structure for KenLM, it makes use of a hash table to store the language model n-grams.

By restricting the training data to sentences of 80 or fewer tokens, the baseline SMT system is trained on 27,481 fewer parallel sentences than the shared task baseline. There are no other differences in the data used; for tuning, development-testing or language model construction.

The baseline SMT system scores nearly one BLEU point higher than the shared task baseline for the IWSLT 2010 (34.57 vs. 33.86) and 2012 (41.07 vs. 40.06) test sets. BLEU scores were calculated using the case-insensitive, multi-bleu perl script provided in the Moses toolkit.

The decoder is set to output word alignments, which are used later for automatic post-editing.

---

[1] Provided as part of the shared task resources

## 4 Extracting Source-language Information

Guided by the ParCor annotation scheme, the following is extracted from the source-language text:

- Position: subject or object ("it" only)

- Function: anaphoric or non-anaphoric (i.e. pleonastic / event, for "it" only)

- Antecedent: for anaphoric pronouns only

The first step is to identify whether the pronoun appears in subject or object position. The pronoun "it" may be used in either position, unlike "they" which is always a subject-position pronoun. When translating into French it is necessary to ensure that each instance of "it" is correctly translated, with different French pronouns used depending on the position that the pronoun fills. Instances of "it" are categorised as being either subject- or object-position pronouns using the dependency parser provided as part of the Stanford CoreNLP tool[2]. Subject-position pronouns are those that participate in an *nsubj* or *nsubjpass* dependency relation.

The next step is to determine the function of each instance of "it". NADA (Bergsma and Yarowsky, 2011) is used as it considers the entire sentence, unlike the pleonastic sieve in the Stanford coreference resolution system (Lee et al., 2011), which uses only fixed expressions to identify pleonastic "it". Instances of "it" with a NADA probability below a specified threshold are treated as non-anaphoric, and those above, as anaphoric. Here, a non-anaphoric pronoun is either an event or pleonastic pronoun; a finer distinction cannot be made using currently available tools. The NADA threshold is set to 0.41 (see Section 6).

For instances of "it" identified as anaphoric, and all instances of "they", the pronoun's nearest non-pronominal antecedent is extracted using the coreference resolution system (Raghunathan et al., 2010; Lee et al., 2011) provided in the Stanford CoreNLP tool[3]. To avoid falsely identifying coreference chains across document boundaries, the source-language text is split into documents prior to coreference resolution. Full coreference chains are retained in case the nearest antecedent is not translated by the baseline SMT system.

NADA and CoreNLP were run on tokenised, but not lower-cased data, in order to ensure parser

[2]Stanford CoreNLP version 3.3.1 `http://nlp.stanford.edu/software/corenlp.shtml`
[3]Considers pronoun-antecedent distances $\leq 3$ sentences

accuracy. The tokenisation and sentence segmentation is the same as that used in the pre-processed data distributed for the shared task. The CoreNLP tool was run with the following annotators: *tokenize*, *ssplit*, *pos*, *lemma*, *ner*, *parse* and *dcoref*. The following parameters were set to true: *tokenize.whitespace* and *ssplit.eolonly*.

## 5 Automatic Post-Editing Rules

Automatic post-editing is applied to the 1-best output of the baseline SMT system described in Section 3. The process makes use of information extracted from the source-language text (Section 4) and the word alignments output by the decoder.

For each source-language pronoun, one of two post-editing rules is applied, depending on whether the pronoun is identified as anaphoric or non-anaphoric. The rules are outlined in Figure 1 and described in detail in the following sections.

### 5.1 Anaphoric Rule

This rule is applied to all instances of "they" and subject-position "it" that are identified as anaphoric, both inter- and intra-sentential. *Cataphoric* pronouns, where the pronoun appears before its antecedent, are very rare (Guillou et al., 2014) and are ignored for the sake of simplicity. Instances of object-position "it" are excluded as the focus of the shared task is on subject-position pronouns only. Target-language pronoun forms are predicted using the projected translation of the head of the nearest non-pronominal antecedent.

**On the source-language side:**

1. Identify the nearest non-pronominal antecedent

2. Identify the antecedent head word (provided by CoreNLP for each antecedent)

3. Using word alignments output by the decoder, project source-language pronoun and antecedent head positions to the SMT output

**On the target-language side (SMT output):**

4. If no antecedent can be found for the pronoun, do not attempt to amend its translation. (It may be non-anaphoric but not detected by NADA)

5. For all other pronouns, use the word alignments to identify the translations of the pronoun and antecedent head

6. Extract the number and gender of the antecedent head translation via a dictionary of

French nouns extracted from the Lefff (Sagot, 2010) and augmented by entries from dict.cc[4]

7. If the antecedent head word is aligned to multiple words in the translation select the rightmost noun (should be the head in most cases)

8. If the antecedent head translation **is a noun**[5]:

    (a) Predict "elle" for feminine, singular; "il" for masculine, singular

    (b) Predict "elles" for feminine, plural; "ils" for masculine, plural

    (c) If the antecedent is split-reference of the format **N and N**, split it into two nouns. If both are feminine, predict "elles", otherwise predict "ils"

9. If the antecedent head translation **is not a noun** (i.e. not in the dictionary) or is not translated:

    (a) Traverse further back through the coreference chain and repeat from *step 5*

    (b) If the antecedent head is not translated, apply a default value. If the source-language pronoun is translated as a pronoun, but not "il/elle" (for "it") or "ils/elles" (for "they"), predict "il" for "it" and "ils" for "they". If the pronoun is not translated, do nothing as the SMT system may have correctly learned to drop a pronoun

10. If the pronoun in the SMT output and the predicted translation disagree, the post-editing rule replaces the translation in the SMT output with the predicted value

This method allows for the prediction of a plural pronoun for cases where an English singular noun is translated into French using a plural noun. For example, "vacation" is singular in English but may be translated as "vacances" (plural) in French.

### 5.2 Non-Anaphoric Rule

This rule is applied to instances of subject-position "it" that are identified as non-anaphoric, i.e. those with a NADA probability below the specified threshold. It does not apply to instances of "they".

The first step is to identify the translation of the pronoun (using the word alignments). The translation that should appear in the post-edited SMT output is then predicted.

---

[4]`www.dict.cc`
[5]If the word is hyphenated and not in the dictionary, look up the right-most part, which should be the head

**1) Translation is an event/pleonastic pronoun:** As NADA does not appear to distinguish event and pleonastic pronouns (i.e. both are considered equally non-anaphoric; see Section 6) it is not straightforward to predict a correct translation for non-anaphoric "it". The French pronoun "ce" may function as both an event and a pleonastic pronoun, but "il" is used only as a pleonastic pronoun. All instances of "it" translated as "ce/c'/il" are left as they are in the SMT output. Changing them may do more harm than good and would be performed in an uninformed manner. The hope is that these pronouns, or at least the pleonastic ones, may be correctly translated using local context.

**2) Translation is another pronoun**: If an instance of "it" is translated as a pronoun outwith the set "ce/c'/il", it will be corrected to the default "ce" (or "c'" if the next word in the SMT output starts with a vowel or silent "h"). The French pronouns "ce/c'/cela/ça" may be used as neutral pronouns, referring to *events*/actions/states or general classes of people/things, and "il/ce/c'/cela/ça" may be used as impersonal pronouns, marking the subject position but not referring to an entity in the text, i.e. *pleonastically* (Hawkins et al., 2001). "ce/c'/cela/ça" may all be used as either pleonastic or event pronouns. "ce" is selected as the default as it occurs most frequently in the training data, suggesting common usage. There are some cases in which only "il" should be used as the impersonal pronoun, such as expressions of time. These are not easy to detect and are therefore ignored.

**3) Translation is not a pronoun**: If an instance of "it" is translated using something other than a pronoun, it is not amended. This may also indicate that the pronoun has been dropped.

**4) No translation**: There is no provision for handling cases where a pleonastic or event pronoun may in fact be required but was dropped in the SMT output. I am not aware of any tools that can separate pleonastic and event instances of "it" for English and inserting a pronoun might not be the correct thing to do in all cases.

If the pronoun in the SMT output and the predicted translation disagree, the post-editing rule replaces the translation in the SMT output with the predicted value.

### 6 Setting the NADA Threshold

NADA returns a probability between 0 and 1, and the decision as to whether an instance of "it" is

anaphoric can be made by thresholding this probability. The NADA documentation suggests a general threshold value of 0.5; for probabilities over this value the pronoun is said to be referential (i.e. anaphoric) and for those below this value, that it is non-referential. However, different threshold values may be appropriate for different genres[6].

The TED-specific NADA threshold was set using the manual ParCor (Guillou et al., 2014) annotations over the TED Talks portion of the corpus. NADA was run over the English TED Talks in ParCor and the probabilities it assigned for each instance of "it" were compared with the pronoun type labels (i.e. anaphoric/pleonastic/event).

There are 61 instances of "it" marked as pleonastic in the ParCor annotations. Looking at *all* 133 instances of "it" in the ParCor TED Talks for which their NADA probabilities fall below 0.5, there are a mixture of pleonastic, event, and "anaphoric with no explicit antecedent" pronouns. These could acceptably be treated as non-referential. However, there are also a number of anaphoric pronouns that fall into this range and it would be unacceptable to treat these as non-referential. Setting the threshold is therefore a trade-off between precision and recall. Whatever threshold is set, there will be both false positives and false negatives. At a threshold of $\leq 0.41$, 37 (60.66%) of pronouns marked as pleonastic in ParCor are correctly identified and 24 (39.34%) are not. 37 pronouns marked in ParCor as event pronouns and 35 anaphoric pronouns (of which 4 have no explicit antecedent) are also (incorrectly) identified as non-referential.

## 7 Post-Editing Statistics

The shared task test set contains 307 instances of "they" and 809 instances of "it". Automated pre-processing of the source-language texts identifies 581 instances of "it" as subject-position pronouns and 228 as object-position pronouns (for which no change will be made). Of the 888 instances of "it" and "they" identified as subject-position pronouns, the translation of 316 are changed in the SMT output by the post-editing rules. 303 changes are applied to pronouns identified as anaphoric (36 "they" and 267 "it") and 13 to pronouns identified as non-anaphoric. The pronoun changes are summarised in Table 2. 10 pronouns were not trans-

lated by the baseline SMT system, and as such, were not considered for amendment.

| Pronoun type | Form | Before | After | Count |
|---|---|---|---|---|
| Non-anaphoric | it | ç | ce/c' | 7 |
| Non-anaphoric | it | cela | ce/c' | 3 |
| Non-anaphoric | it | elle | ce/c' | 1 |
| Non-anaphoric | it | le | ce/c' | 1 |
| Non-anaphoric | it | on | ce/c' | 1 |
| Anaphoric | it | il | ils | 3 |
| Anaphoric | it | il | elle | 51 |
| Anaphoric | it | il | elles | 3 |
| Anaphoric | it | elle | il | 17 |
| Anaphoric | it | elle | ils | 1 |
| Anaphoric | it | le/l' | il | 3 |
| Anaphoric | it | on | il | 1 |
| Anaphoric | it | ç | il | 10 |
| Anaphoric | it | ç | ils | 2 |
| Anaphoric | it | ç | elle | 5 |
| Anaphoric | it | cela | il | 6 |
| Anaphoric | it | cela | elle | 3 |
| Anaphoric | it | cela | elles | 1 |
| Anaphoric | it | ce/c' | il | 84 |
| Anaphoric | it | ce/c' | ils | 5 |
| Anaphoric | it | ce/c' | elle | 68 |
| Anaphoric | it | ce/c' | elles | 4 |
| Anaphoric | they | ils | elles | 32 |
| Anaphoric | they | elles | ils | 4 |
| **Total** | | | | 316 |

Table 2: Automated post-editing changes

The most frequent changes are "c'/ce" → "il" (84), "c'/ce" → "elle" (68), "il" → "elle" (51), and "ils" → "elles" (32). The change "c'/ce" → "il/elle" takes place due to the decision to use gendered translations of all instances of "it" identified as anaphoric (even if "c'/ce" might also have been an acceptable translation). Biases in the training data may account for some of the other changes. For example, the change "ils" → "elles" may result from the common alignment of "they" to "ils" which arises due to the rule in French that "ils" is used unless all of the antecedents are feminine (in which case "elles" is used). This may result in more masculine pronouns requiring replacement with a feminine pronoun than vice versa.

The changes "il" → "elle" and "ils" → "elles" are made to conform with the gender of the translation of the antecedent head of an anaphoric pronoun. The post-editing rules also allow for changes from singular to plural (and vice versa) and from one number and gender to another. For example in translating "it" → "vacation" the anaphoric rule would allow for an instance of "il" (masc. sg.) in the SMT output to be changed to "elles" → "vacances" (fem. pl.).

---

[6]TED Talks are considered out-of-domain. NADA was trained using the Penn Treebank and Google N-Grams corpus

## 8 Results

The official shared task results report a BLEU score of 36.91 for the post-edited SMT output. This score is lower than the official baseline system (37.18), comparable with the UU-Tiedemann system (36.92), and higher than the other competing systems. However, the post-editing system outperformed only two of the five competing systems in terms of the *accuracy* measures, suggesting that BLEU is a poor measure of pronoun translation performance. The *accuracy with OTHER* measure reveals that the post-edited SMT output contains correct translations for only 114/210 pronoun instances, according to human judgements.

There is a small decrease of 0.36 BLEU between the baseline system used to provide SMT output and the post-edited version for the test set (38.83 vs. 38.47 respectively, as calculated using case-insensitive multi-bleu[7]).

An examination of the human judgements from the shared task manual evaluation reveals that the post-editing process makes many mistakes. 34 instances were worsened by post-editing and only 9 improved. The remaining instances were neither better nor worse following post-editing. Translation accuracy differs for "it" and "they". For "it" 32 instances are judged to be correct vs. 60 incorrect. The opposite is observed for "they", with 47 instances judged to be correct vs. 14 incorrect. (Instances marked as "other" or "bad translation" cannot be commented upon further and are excluded from the counts). The poor translation of "it" could be due to the method used to identify anaphoric and non-anaphoric instances (no such method was used for "they"), differences in coreference resolution accuracy for "it" and "they", or something else entirely.

## 9 Limitations of Post-Editing

Although specific failures in the baseline SMT system, the external tools and the post-editing rules await detailed analysis, the following possible problems with the external tools should at least be considered: incorrect identification of subject-position "it", of non-anaphoric pronouns and of antecedents. These problems may arise from a mismatch between the TED Talks domain, and the domain of the data that the tools were trained on.

---

[7]The official shared task BLEU scores appear to have been calculated using a different method

As the post-editing rules affect only pronouns, agreement issues may occur. For example, if the baseline SMT system outputs "ils sont partis" ("they[masc] have left") and the post-editing rules amend "ils" to "elles", the verb "partis" should also be amended: "elles sont *parties*" ("they[fem] have left"). Agreement issues could be addressed within a dependency-parser-based post-editing framework such as the Depfix system for Czech (Mareček et al., 2011; Rosa, 2014).

Another limitation is the lack of an available tool for detecting event pronouns. Whilst NADA appears to detect some of these, it is an accidental consequence of its inability to distinguish a pleonastic ("il/ce") from an event pronoun ("ce"). NADA was also shown to perform poorly for TED data (see Section 6).

While post-editing rules could potentially be written to insert a pronoun in the SMT output where one is syntactically required in the the target language, or to delete a pronoun for syntactic or stylistic reasons, this was not done in the current system.

The approach may also be difficult to extend to other languages which are less well provisioned in terms of parsers and coreference resolution systems or for which baseline SMT quality is poor.

## 10 Summary and Future Work

The post-editing approach makes use of two pronoun-specific rules applied to the output of a baseline English-to-French phrase-based SMT system. One rule handles anaphoric pronouns, the other handles non-anaphoric pronouns.

Before extending this work to develop new rules or applying the technique to other language pairs, it is important to first understand where the post-editing method performs well and where it performs poorly. A detailed analysis of the post-edits as compared with the human judgements from the manual evaluation would be a logical first step. Limitations of both the external tools and the post-editing rules should be assessed.

### Acknowledgements

# References

Shane Bergsma and David Yarowsky. 2011. NADA: A Robust System for Non-Referential Pronoun Detection. In *Proceedings of DAARC 2011*, pages 12–23.

Liane Guillou, Christian Hardmeier, Aaron Smith, Jörg Tiedemann, and Bonnie Webber. 2014. ParCor 1.0: A Parallel Pronoun-Coreference Corpus to Support Statistical MT. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, Reykjavik, Iceland. European Language Resources Association (ELRA).

Christian Hardmeier, Preslav Nakov, Sara Stymne, Jörg Tiedemann, Yannick Versley, and Mauro Cettolo. 2015. Pronoun-focused MT and cross-lingual pronoun prediction: Findings of the 2015 DiscoMT shared task on pronoun translation. In *Proceedings of the Second Workshop on Discourse in Machine Translation*, Lisbon, Portugal.

Roger Hawkins, Richard Towell, and Marie-Noëlle Lamy. 2001. *French Grammar and Usage*. Hodder Arnold, 2 edition.

Kenneth Heafield. 2011. KenLM: Faster and smaller language model queries. In *Proceedings of the Sixth Workshop on Statistical Machine Translation*.

Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open Source Toolkit for Statistical Machine Translation. In *Proceedings of the 45th Annual Meeting of the ACL on Interactive Poster and Demonstration Sessions*, ACL '07, pages 177–180, Stroudsburg, PA, USA. Association for Computational Linguistics.

Heeyoung Lee, Yves Peirsman, Angel Chang, Nathanael Chambers, Mihai Surdeanu, and Dan Jurafsky. 2011. Stanford's Multi-pass Sieve Coreference Resolution System at the CoNLL-2011 Shared Task. In *Proceedings of the Fifteenth Conference on Computational Natural Language Learning: Shared Task*, CONLL Shared Task '11, pages 28–34, Stroudsburg, PA, USA. Association for Computational Linguistics.

David Mareček, Rudolf Rosa, Petra Galuščáková, and Ondřej Bojar. 2011. Two-step Translation with Grammatical Post-processing. In *Proceedings of the Sixth Workshop on Statistical Machine Translation*, WMT '11, pages 426–432, Stroudsburg, PA, USA. Association for Computational Linguistics.

Franz Josef Och. 2003. Minimum Error Rate Training in Statistical Machine Translation. In *Proceedings of the 41st Annual Meeting on Association for Computational Linguistics - Volume 1*, ACL '03, pages 160–167, Stroudsburg, PA, USA. Association for Computational Linguistics.

Karthik Raghunathan, Heeyoung Lee, Sudarshan Rangarajan, Nathanael Chambers, Mihai Surdeanu, Dan Jurafsky, and Christopher Manning. 2010. A Multi-pass Sieve for Coreference Resolution. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, EMNLP '10, pages 492–501, Stroudsburg, PA, USA. Association for Computational Linguistics.

Rudolf Rosa. 2014. Depfix, a Tool for Automatic Rule-based Post-editing of SMT. *The Prague Bulletin of Mathematical Linguistics*, 102:47–56.

Benoît Sagot. 2010. The Lefff, a Freely Available and Large-coverage Morphological and Syntactic Lexicon for French. In *7th International Conference on Language Resources and Evaluation (LREC 2010)*, Valletta, Malta.

Jochen Weiner. 2014. Pronominal anaphora in machine translation. Master's thesis, Karlsruhe Institute of Technology, Karlsruhe, Germany.