

IUCL: Combining Several Information Sources for SemEval Task 5

Alex Rudnick, Levi King, Can Liu, Markus Dickinson, Sandra Kübler

Indiana University
Bloomington, IN, USA

{alexr, leviking, liucan, md7, skuebler}@indiana.edu

Abstract

We describe the Indiana University system for SemEval Task 5. The system is based on combining different information sources to arrive at a final L2 guess, incorporating: phrase tables, an L2 language model, and collocational tendencies. We also consider different data sources.

1 Introduction

The task of translating an L1 fragment occurring in the midst of an L2 sentence is one in which a phrase occurs in an already-rich target language (L2) context; this makes the task quite different from translation in the general case. For this reason, we decided to approach the problem by combining standard Machine Translation technology with target language information, such as contextual relationships. This is broken down into various steps: 1) constructing candidate translations for the L1 fragment, including weights for the likelihood of each translation; 2) scoring candidate translations via a language model of the L2; 3) scoring candidate translations via dependency-driven word similarity measure (Lin, 1998) (which we call *SIM*) and 4) combining the scores from 1)-3) via minimized error rate training (MERT), to arrive at a final solution. Step 1) models transfer knowledge between the L1 and L2; step 2) models facts about the L2 grammar, i.e., what translations fit well into the local context; step 3) models collocational and semantic tendencies of the L2; and step 4) gives different weights to each of the three sources of information. Although we did not finish step 3) in time for the official results, we report it here, as it represents the most novel aspect of the system – namely, the exploitation of the

rich L2 context – and it results in our team’s best system. In general, our system is fully language-independent, with accuracy varying due to the size of data sources and quality of input technology (e.g., syntactic parse accuracy).

2 Data Sources

The data sources serve two major steps of our system: for L2 candidate generation, we use Europarl and BabelNet; and for candidate ranking using L2 context, we use Wikipedia and the Google Books Syntactic N-grams. For the candidate generation, we later experiment with expanding to larger sets of data (section 4.3).

Europarl The Europarl Parallel Corpus (Europarl, v7) (Koehn, 2005) is a corpus of proceedings of the European Parliament, containing 21 European languages with sentence alignments. From this corpus, we build phrase tables for English-Spanish, English-German, French-English, Dutch-English.

BabelNet In the cases where the constructed phrase tables do not contain a translation for a source phrase, we will need to back off to smaller phrases and find candidate translations for these components. To better handle sparsity, we extend look-up using the multilingual dictionary BabelNet, v2.5 (Navigli and Ponzetto, 2012) as a way to look up translation candidates.

Wikipedia For German and Spanish, we use recent Wikipedia dumps, obtained through the Wikipedia Extractor tool.¹ To save time during parsing, sentences longer than 25 words are removed. The remaining sentences are POS-tagged and dependency parsed using Mate Parser with its pre-trained models (Bohnet, 2010; Bohnet and Kuhn, 2012; Seeker and Kuhn, 2013). To keep

MD: check
version
number

LK: Alex
may have the
exact info for
the dumps if
needed

This work is licenced under a Creative Commons Attribution 4.0 International License. Page numbers and proceedings footer are added by the organizers. License details: <http://creativecommons.org/licenses/by/4.0/>

¹http://medialab.di.unipi.it/wiki/Wikipedia_Extractor

our English Wikipedia data set comparable in size to the German and Spanish sets, we choose an older (2006), smaller dump. Long sentences were removed, and the remaining sentences are POS-tagged and dependency parsed using the pre-trained Stanford Parser (Klein and Manning, 2003; de Marneffe et al., 2006). The resulting sizes of the data sets are (roughly): German: 389M words, 28M sentences; Spanish: 147M words, 12M sentences; English: 253M words, 15M sentences. Dependencies extracted from these parsed data sets serve as training for the SIM system described in section 3.3.

Google Books Syntactic N-grams For English, we also trained a SIM system on the arcs dataset of the Google Books Syntactic N-Grams (Goldberg and Orwant, 2013), which has 919M items.

3 Our System

3.1 Constructing Candidate Translations

Given that this is essentially a local translation problem, we use as our starting point a phrase table constructed from parallel text, weighted with probabilities. We use GIZA++ (Och and Ney, 2000) and Moses (Koehn et al., 2007) to construct the phrase tables. The quality of the phrase table depends upon the size of the data—an issue we discuss with larger phrase tables in section 4.3—but in the case of missing phrases from the table, we back off to subphrases in the following way: ...

- some text: split up the source phrase into smaller components and look up each component from Babelnet.
- MD: 1. What is the exact back-off procedure? 2. Do we back off only in the case of missing phrases, or do we consider multiple possibilities for one phrase, even if its already in the table?

3.2 Scoring Candidate Translations via a L2 Language Model

To examine how well a phrase fits into an L2 context, a quick and intuitive method is to use an N-gram language model (LM), built from the L2, and rank the candidate phrases based on their LM scores. With a large vocabulary, the construction and query of an N-gram language model is potentially very time-consuming, and thus we use the KenLM Language Model Toolkit for efficiency (Heafield, 2011).

3.3 Scoring Candidate Translations via Dependency-Based Word Similarity

Definition The candidate ranking based on the N-gram language model is based on very shallow information. We can also rank the candidate phrases based on how well each of the components fits into the L2 context based on syntactic information. In this case, the fitness is measured in terms of dependency-based word similarity based on dependency triples consisting of the the head, the dependent, and the dependency label. We slightly adapted the word similarity measure by Lin (1998):

SK: Can, please check comment

$$SIM(w_1, w_2) = \frac{2 \prod c(h, d, l)}{c(h, -, l) + c(-, d, l)} \quad (1)$$

SK: how do w1 and w2 relate to h,d,l?

where $c(h, d, l)$ is the frequency that a particular (*head, dependent, label*) dependency triple occurs in the L2 corpus. $c(h, -, l)$ is the frequency that a word occurs as a head in a dependency labeled l with any dependent. $c(-, d, l)$ is the frequency that a word occurs as a dependent in a dependency labeled l with any head.

The fitness of a phrase is the average word similarity over all its components. For example the fitness of the phrase "eat with chopsticks" would be computed as:

$$fit(\text{eat with chopsticks}) = \frac{SIM(\text{eat}) + SIM(\text{with}) + SIM(\text{chopsticks})}{3} \quad (2)$$

Since we consider the heads and dependents of a target phrase component, these may be situated inside or outside the phrase. Both cases are included in our calculation, thus enabling us to consider a more variable, syntactically determined local context of the phrase. By basing the calculation on a single words's head and dependent, we focus on tightly associated words and thus avoid data sparseness issues.

Back-Off Lexical-based dependency triples suffer from data sparsity, so in addition to computing the lexical fitness of a phrase, we also calculate the POS fitness. For example, the POS fitness of "eat with chopsticks" would be computed as follows:

$$\text{fit}(\text{eat/VBG with/IN chopsticks/NNS}) = \frac{\text{SIM}(\text{VBG}) + \text{SIM}(\text{IN}) + \text{SIM}(\text{NNS})}{3} \quad (3)$$

Storing and Caching The large vocabulary and huge number of combinations of our head/dependent/label” triples poses an efficiency problem when querying the dependency-based word similarity values. Thus, we stored the dependency triples in a database with a Python programming interface (sqlite3) and built database indexes on the frequent query types. However, for frequently searched dependency triples, re-querying the database is still inefficient. Thus we built a query cache to store the recently queries triples. Using the database and cache significantly speeds up our system.

Combining SIM and LM During training time, we obtained dependency-based word similarity statistics for all target languages using corresponding dependency triples. During testing time, the word similarity value is queried and combined with the N-gram language model to rank the candidate phrase.

SK: how???

3.4 Tuning Weights with MERT

- ZMERT (Zaidan, 2009) - exact set-up

4 Experiments

4.1 Official System

(Bird et al., 2009)

4.2 Expansion #1: Experiments with SIM over Dependencies

4.3 Expansion #2: Experiments with Large Phrase Tables

- EU bookshop corpus - MultiUN corpus - you can find so many corpora on Opus!! (Tiedemann, 2012)

5 Conclusion

References

Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural Language Processing with Python*. O’Reilly Media.

Bernd Bohnet and Jonas Kuhn. 2012. The best of both worlds – A graph-based completion model for

transition-based parsers. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics (EACL)*, pages 77–87, Avignon, France.

Bernd Bohnet. 2010. Top accuracy and fast dependency parsing is not a contradiction. In *Proceedings of the 23rd International Conference on Computational Linguistics (COLING)*, pages 89–97, Beijing, China.

Marie-Catherine de Marneffe, Bill MacCartney, and Christopher Manning. 2006. Generating typed dependency parses from phrase structure parses. In *Proceedings of LREC-06*, Genoa, Italy.

Yoav Goldberg and Jon Orwant. 2013. A dataset of syntactic-ngrams over time from a very large corpus of English books. In *Second Joint Conference on Lexical and Computational Semantics (*SEM)*, pages 241–247, Atlanta, GA.

Kenneth Heafield. 2011. KenLM: faster and smaller language model queries. In *Proceedings of the EMNLP 2011 Sixth Workshop on Statistical Machine Translation*, pages 187–197, Edinburgh, Scotland, United Kingdom, July.

Dan Klein and Christopher Manning. 2003. Accurate unlexicalized parsing. In *Proceedings of ACL-2003*, pages 423–430, Sapporo, Japan.

Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, Prague, Czech Republic.

Philipp Koehn. 2005. Europarl: A parallel corpus for statistical machine translation. In *MT Summit*.

Dekang Lin. 1998. An information-theoretic definition of similarity. In *International Conference on Machine Learning (ICML)*, volume 98, pages 296–304.

Roberto Navigli and Simone Paolo Ponzetto. 2012. BabelNet: The automatic construction, evaluation and application of a wide-coverage multilingual semantic network. *Artificial Intelligence*, 193:217–250.

Franz Josef Och and Hermann Ney. 2000. Improved statistical alignment models. In *Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 440–447, Hong Kong.

Wolfgang Seeker and Jonas Kuhn. 2013. Morphological and syntactic case in statistical dependency parsing. *Computational Linguistics*, 39(1):23–55.

Jörg Tiedemann. 2012. Parallel data, tools and interfaces in OPUS. In *Proceedings of the 8th International Conference on Language Resources and Evaluation (LREC)*, pages 2214–2218, Istanbul, Turkey.

Omar F. Zaidan. 2009. Z-MERT: A fully configurable open source tool for minimum error rate training of machine translation systems. *The Prague Bulletin of Mathematical Linguistics*, 91:79–88.