

Introduction

Translation Corpora: Annotation, Exploitation, Evaluation

Parallel corpora, i.e. collections of originals and their translations, can be used in various ways for the benefit of translation studies, machine translation, linguistics, computational linguistics or simply the human translator.

In computational linguistics, translation corpora have been employed for machine translation but also for term extraction, word sense disambiguation etc. as early as the 1980s (important milestones being Nagao 1984 and Brown et al. 1990). One of the early electronic resources is the Canadian Hansard which was initially used for implementing sentence alignment (Gale and Church 1991), a task that is now a standard feature of applications such as translation memories. Moreover, parallel corpora are used as data basis for multilingual grammar induction, automatic lexicography and many other tasks in information extraction and language processing across different languages.

In translation studies, the focus is more on identifying features that distinguish translations from original texts. From this perspective, the main research interest lies in the detection of patterns of (inevitable) modifications introduced by the translator(s) along the way in terms of local solutions, added information or even larger changes in the register of the text. These modifications may be individual to a given translation task or a translation pair but they may also instantiate typical features of translated text that make translations different from non-translated texts in a wide range of linguistic features. The investigation of corpora is an obvious method to empirically detect these distinctive properties of translations and has been employed since the 1990s as witnessed by Baker (1993, 1996), Johansson and Ebeling (1996) and more recently by Hansen (2003), Teich (2003), Mauranen and Kujamäki (2004) and Hansen-Schirra, Neumann, and Steiner (forthcoming). Furthermore, parallel corpora are used as reference works for translation teaching and in professional translation settings since they enable quick and interactive access to translation solutions (e.g. as translation memories).

Exchange between the translation studies and the computational linguistics communities has traditionally not been very intense. Among other things, this is reflected by the different views on parallel corpora. While computational linguistics does not always strictly pay attention to the translation direction (e.g. when translation rules are extracted from (sub)corpora which actually only consist of translations), translation studies are amongst other things concerned with exactly comparing source and target texts (e.g. to draw conclusions on interference and standardization effects). However, there has recently been more exchange between the two fields – especially when it comes to the annotation of parallel corpora. This special issue brings together the different research perspectives. Its contributions

show – from both perspectives – how the communities have come to interact in recent years.

With issues of the creation of large parallel data collections including multiple annotation and alignment largely solved, the exploitation of these collections remains a bottleneck. In order to effectively use annotated and aligned parallel corpora, the interaction of the different disciplines involved addresses the following issues:

- *Query tools:* We can expect basic computer literacy from researchers nowadays. However, the gap between writing query or evaluation scripts and program usability is immense. One way to address this is by building web query interfaces. Yet in general, what are the claims and possibilities for creating interfaces that address a broader public of researchers using multiply annotated and aligned corpora? An additional ongoing question is the most efficient storage form: are data base formats superior to other formats?
- *Information extraction strategies:* The quality of the information extracted by a query heavily depends on the quality of the annotation of the underlying corpus, i.e. on precision and recall of annotation and alignment. Furthermore, the question arises how we can ensure high precision and recall of queries (while possibly keeping query construction efficient). What are the strategies to compose queries which produce high-quality results? How can the query software contribute to this goal?
- *Corpus quality:* Several criteria for corpus quality have been developed (e.g. in the context of standardization initiatives). Quality can be influenced before compilation by ensuring the balance of the corpus (in terms of register and sample size), its representativeness etc. Also, inter-annotator agreement and – to a lesser extent – intra-annotator agreement are an issue. But, how can we make corpora thus created fit for automatic exploitation? This involves issues such as data format validity throughout the corpus, robust (if not 100% correct) processing with corpus tools/APIs and the like. What are relevant criteria and how can they be addressed?
- *Corpus maintenance:* Beyond the validity of the data format, maintenance of consistent data collections is a more complex task, particularly if the data collection is continually expanded. A change of the annotation scheme entails adjustments in the existing annotation. Questions to this end include whether automatic adjustment is possible and how it can be achieved. Maintenance may also involve compatibility with and adaptations to new data formats. How can we ensure sustainability of the data formats?

A Colloquium held at the Corpus Linguistics 2009 Conference at the University of Liverpool was concerned with the interface between the requirements of linguists and translation studies working with parallel corpora and computational linguists providing the tools and exploiting the corpora for their purposes. In this sense, it was closely related to and a continuation of the workshop “Multilingual Corpora: Linguistic Requirements and Technical Perspectives” held at the Corpus Linguistics

2003 Conference at Lancaster University (see Neumann and Hansen-Schirra 2003). The present special issue is a collection of contributions arising out of this Colloquium. In what follows we outline the contributions responding to some of the questions posed above.

The volume sets off with a focus on annotation, alignment and query on the syntactic level: Volk, Marek and Samuelsson discuss a trilingual parallel treebank, the Stockholm Multilingual Treebank SMULTRON. The ultimate purpose of the resource is its exploitation for machine translation, a typical application scenario for parallel treebanks. Interestingly, the resource only consists of translations in the three languages English, German and Swedish. The authors discuss solutions for some important questions in querying the treebank, thus focussing on an issue in working with parallel corpora that typically only arises at a later stage of corpus construction but that is not trivial at all.

In their contribution, Vintar and Fišer discuss the exploitation of multilingual resources – and translations in particular – for a monolingual computational linguistic task, the construction and enrichment of the Slovene WordNet. They turn the problem of a lesser-studied language into an advantage in drawing on the rich body of translations existing for Slovene. At various stages of their work, parallel corpora are used to disambiguate word senses with the help of translations – making use of a typical feature of translation, namely settling on one interpretation of ambiguous items in the source text –, as well as to extract a bilingual lexicon of word-aligned items in order to enrich the resource with domain-specific lexical items. Vintar and Fišer show how monolingual resources can be successfully exploited with the help of parallel corpora that contain the required information.

Fantinuoli's contribution demonstrates an even more practice-oriented exploitation of corpora, both monolingual and parallel. Fantinuoli describes the design of a software, InterpretBank, which supports conference interpreters in all stages of their work. Based on Baroni and Bernardini's (2004) BootCat mechanism, it harvests the web for domain-specific documents given a set of search terms, performs term extraction on them and uses additional resources, e.g. Wikipedia or bilingual online dictionaries, to propose definitions, translations, collocations and keyword-in-context information. All available modules, for harvesting, management and retrieval, are adapted to the specific needs of interpreters, reducing the time needed for preparation and allowing for efficient retrieval while interpreting. A pilot module adds the possibility to include parallel resources, e.g. translation memories or the OPUS corpora, in the preparation phase.

The contribution by Čulo, Hansen-Schirra, Maksymski and Neumann returns to a more theory-oriented topic. It discusses the analysis of the bilingual CroCo Corpus, a richly annotated and aligned corpus of English and German translations and originals, with respect to a translation-specific research question. It exemplifies the exploitation of a resource that comes close to a parallel treebank for a research question that has a long history in translation studies, namely the study of shifts (e.g. Vinay and Darbelnet 1958, Catford 1965 etc.). The goal of this contribution is a

heuristic identification of shifts in translation that can then be interpreted as properties of translations. While the main aim of the study is to advance empirical knowledge in the field of translation studies, it also has some clear implications for computational handling of translation shifts for instance in machine translation

The translation-related research question investigated by Čulo et al. sets the scene for the final paper in this special issue: Alves and Vale introduce an innovative approach to adopting a corpus perspective on psycholinguistic research into the translation process. The authors describe LITTERAE, a computer tool that allows annotating linear representations of the process of producing a translation of a source text. They then go on to discuss quantitative findings yielded with LITTERAE which suggest certain patterns in target text production. The paper represents a highly interesting way of reducing the gap between corpus-based and process-oriented investigations of translations. It thus rounds off this special issue with a perspective beyond corpus linguistics.

Summarizing, the articles in this special issue address a number of the issues discussed above: Vintar and Fišer are concerned with information extraction from various multilingual resources, whereas Čulo et al. exemplify the linguistic interpretation of parallel data on the basis of a heuristic information extraction procedure. Information extraction as well as its interpretation is also exemplified in Alves and Vale. Questions of querying are also a major concern of Volk et al. They also discuss corpus quality, in particular annotation quality. This, in turn, is also addressed by Padó. The only area of interest not covered by one of the contributions is maintenance of continually expanding resources. This is an area addressed by work in the area of sustainability of corpora, for instance in the framework of the European CLARIN project¹ and similar national initiatives (e.g. Rehm et al. 2009).

Acknowledgements

We believe that this volume provides a good overview of some important issues in working with parallel corpora, not only focussing on computational issues but also giving insight into the linguistic analysis of translations. If this succeeds this is not least thanks to the efforts the reviewers put into providing feedback to the authors and thus ensuring the quality of this issue. The reviewers were: Sabine Bartsch (Technische Universität Darmstadt), Stefan Evert (University of Osnabrück), Johann Haller (IAI, Saarbrücken), Kerstin Kunz (Saarland University, Saarbrücken), Anke Lüdeling (Humboldt-Universität Berlin), Reinhardt Rapp (University of Mainz, Gernersheim), Josef Schmied (Chemnitz University of Technology), Erich Steiner (Saarland University, Saarbrücken), Elke Teich (Saarland University, Saarbrücken), Mihaela Vela (German Research Center for Artificial Intelligence (DFKI),

¹ <http://www.clarin.eu/external/>, last visited 9 March 2010

Saarbrücken), Andreas Witt (Institute for the German Language (IDS), Mannheim). We are also grateful to the authors for their contributions and collaboration.

References

- Baker, Mona. 1993. Corpus Linguistics and Translation Studies. Implications and Applications. In *Text and technology. In honour of John Sinclair*, ed. Mona Baker, Gill Francis, and Elena Tognini-Bonelli, 233-250. Amsterdam, Philadelphia: Benjamins.
- . 1996. Corpus-based translation studies: The challenges that lie ahead. In *Terminology, LSP and Translation. Studies in Language Engineering in Honour of Juan C. Sager*, ed. Harold Somers, 175-186. Amsterdam: Benjamins.
- Baroni, Marco, and Silvia Bernardini. 2004. BootCaT: Bootstrapping Corpora and Terms from the Web. *Proceedings of LREC2004*, 1313-1316. Lisbon: ELDA.
- Brown, Peter F., John Cocke, Stephen A. Della Pietra, Vincent J. Della Pietra, Fredrick Jelinek, John D. Lafferty, Robert L. Mercer, and Paul S. Roossin. 1990. A Statistical Approach to Machine Translation. *Computational Linguistics* 16, no. 2: 79-85.
- Catford, John C. 1965. *A linguistic theory of translation. an essay in applied linguistics*. Oxford: Oxford University Press.
- Gale, William A., and Kenneth W. Church. 1991. Identifying Word Correspondences in Parallel Texts. In *Speech and Natural Language. Proceedings of a Workshop Held at Pacific Grove, California, February 19-22, 1991*, 152-157. Pacific Grove, California: Morgan Kaufmann.
- Hansen, Silvia. 2003. *The Nature of Translated Text. An Interdisciplinary Methodology for the Investigation of the Specific Properties of Translations*. Saarbrücken: DFKI/Universität des Saarlandes.
- Hansen-Schirra, Silvia, Stella Neumann, and Erich Steiner. forthcoming. *Cross-linguistic Corpora for the Study of Translations - Insights from the language pair English-German*. Berlin: de Gruyter.
- Johansson, Stig, and Jarle Ebeling. 1996. Exploring the English-Norwegian Parallel Corpus. In *Synchronic Corpus Linguistics*, ed. Carol E. Percy, Charles F. Meyer, and Ian Lancashire, 3-15. Amsterdam: Rodopi.
- Mauranen, Anna, and Pekka Kujamäki, eds. 2004. *Translation universals. Do they exist?* Amsterdam, Philadelphia: Benjamins.
- Nagao, Makoto. 1984. A framework of a mechanical translation between Japanese and English by analogy principle. In *Proceedings of the international NATO symposium on Artificial and human intelligence, Lyon, France, 173-180*. New York: Elsevier.
- Neumann, Stella, and Silvia Hansen-Schirra, eds. 2003. *Proceedings of the Workshop on Multilingual Corpora: Linguistic Requirements and Technical Perspectives. Corpus Linguistics Conference 2003, 27 March 2003*. Lancaster. <http://www.coli.uni-saarland.de/conf/muco03/Proceedings.htm>.

- Rehm, Georg, Oliver Schonefeld, Andreas Witt, Erhard Hinrichs, and Marga Reis. 2009. Sustainability of annotated resources in linguistics: A web-platform for exploring, querying, and distributing linguistic corpora and other resources. *Literary & Linguistic Computing* 24, no. 2: 193-210.
- Teich, Elke. 2003. *Cross-linguistic Variation in System and Text. A Methodology for the Investigation of Translations and Comparable Texts*. Berlin and New York: Walter de Gruyter.
- Vinay, Jean-Paul, and Jean Darbelnet. 1958. *Stylistique comparée du français et de l'anglais. Méthode de traduction*. Paris: Didier.