



Special issue: selected papers from LREC 2016

Nancy Ide¹ · Nicoletta Calzolari²

Published online: 5 June 2019
© Springer Nature B.V. 2019

This special issue of *Language Resources and Evaluation* includes papers selected from among those presented at LREC on May 25–27, 2016 in Potorož, Slovenia. The selected papers were recommended by LREC Scientific Committee members during the reviewing process and have been expanded and elaborated for inclusion in the special issue as full-length journal articles.

Three of the four included papers describe a diverse array of language resources. “The DialogBank: Dialogues with Interoperable Annotations” describes a gold standard corpus annotated for dialogue acts, using an ISO scheme derived from years of standards development in the field along with two other schemes, all of which are shown to be interoperable. The paper, which substantially expands the original LREC paper (Bunt et al. 2016), provides a comprehensive description of the development of the corpus and its annotation schemes, together with an assessment of the difficulties encountered in the process of development and current limitations. A similarly comprehensive resource overview is provided in “SenseDefs: A Multilingual Corpus of Semantically Annotated Textual Definitions”, which elaborates the description of a large-scale corpus in multiple languages comprised of sense-annotated definitions using the BabelNet unified sense inventory that was originally presented at LREC 2016 (Camacho-Collados et al. 2016). “The Spoken Wikipedia Corpus Collection: Harvesting, Alignment and Exploitation” describes in detail the conversion of data from Spoken Wikipedia to an aligned corpus in

✉ Nancy Ide
ide@cs.vassar.edu
Nicoletta Calzolari
glottolo@ilc.cnr.it

¹ Department of Computer Science, Vassar College, Poughkeepsie, USA

² Istituto di Linguistica Computazionale “A. Zampolli” - CNR, Pisa, Italy

English/German/Dutch, also first presented at LREC 2016 (Köhn et al. 2016). The authors provide an open-source software pipeline that downloads, extracts, normalizes, and aligns Spoken Wikipedia data, so that additional languages can be easily added to the resource. The fourth paper, “Applying Data Mining and Machine Learning Techniques for Sentiment Shifter Identification”, presents a methodology rather than a resource. It significantly extends the original LREC paper (Noferesti and Shamsfard 2016) by presenting and evaluating three methods for identifying words and expressions that affect text polarity in order to improve the accuracy of opinion mining systems, including a machine learning-based algorithm and two weighted association rule mining algorithms.

Among the three major resources presented in this *LRE* special issue, it is notable that all are freely available and, in particular, that the creation and evaluation of these resources is described in sufficient detail to fully understand their fundamental linguistic properties as well as their format and contents. Comprehensive descriptions of resource creation have not always been readily available for reference by others, which has ramifications for the *replicability/reproducibility* of results of studies that rely on specific resources. As a step toward remedying this situation, *LRE* has actively fostered open availability of data and code as a means to enable research replicability and reproducibility, which has become an important issue for all scientific disciplines in the past few years.

Following the success of the 2016 LREC 4REAL Workshop on Replicability and Reproducibility of Research,¹ *LRE* introduced an additional category for submissions to the journal that will be featured in a special section (see “Replicability and reproducibility of research results for human language technology: Introducing an LRE special section” (Branco et al. 2017). In addition, the *LRE* review form (and, for the 2020 conference, the LREC review form) now asks for an assessment of the replicability of results as well as resource creation and evaluation procedures, along with an indication of open availability of data and code. In this way we hope to influence practices in the field around open availability of code and data as well as the comprehensive reporting of research results and, especially, procedures for resource creation, evaluation, and use. Because much of the discussion concerning replicability and reproducibility of scientific research focuses on transparency of methods, the *LRE* community has an important duty to promote best practices for the description and availability of language resources, which differ significantly in their nature from resources used in other fields (e.g., genomics and other biomedical areas) in the complexity of preparation procedures, formats, and added content. These resources are the foundation of work in the field of human language technology, but it has been shown that a lack of even basic processing information about resource preparation can make replication impossible (see, for example, Fokkens-Zwirello et al. 2013).

LRE and its sister conference LREC continue to be the premier venues for the publication of work on all aspects of resource development, use, and evaluation. The papers in this special issue are representative of the breadth and quality of work

¹ Proceedings available at http://www.lrec-conf.org/proceedings/lrec2016/workshops/LREC2016Workshop-4REAL_Proceedings.pdf.

in the field, which, in part through community efforts to identify best practices and encourage reproducible results, continues to grow and thrive.

Nancy Ide, Nicoletta Calzolari
LRE Co-editors-in-chief

References

- Branco, A., Cohen, K. B., Vossen, P., Ide, N., & Calzolari, N. (2017). Replicability and reproducibility of research results for human language technology: Introducing an LRE special section. *Language Resources and Evaluation*, 51(1), 1–5. <https://doi.org/10.1007/s10579-017-9380-0>.
- Bunt, H., Petukhova, V., Malchanau, A., Wijnhoven, K., & Fang, A. (2016). The dialogbank. In N. Calzolari, K. Choukri, T. Declerck, S. Goggi, M. Grobelnik, B. Maegaard, J. Mariani, H. Mazo, A. Moreno, J. Odiijk & S. Piperidis (Eds.) *Proceedings of the tenth international conference on Language Resources and Evaluation (LREC 2016)*, European Language Resources Association (ELRA), Paris, France.
- Camacho-Collados, J., Bovi, C.D., Raganato, A., & Navigli, R. (2016). A large-scale multilingual disambiguation of glosses. In N. Calzolari, K. Choukri, T. Declerck, S. Goggi, M. Grobelnik, B. Maegaard, J. Mariani, H. Mazo, A. Moreno, J. Odiijk & S. Piperidis (Eds.) *Proceedings of the tenth international conference on Language Resources and Evaluation (LREC 2016)*, European Language Resources Association (ELRA), Paris, France.
- Fokkens-Zwirello, A., van Erp, M., Postma, M., Pedersen, T., Vossen, P., & Freire, N. (2013). Offspring from reproduction problems: What replication failure teaches us. In P. Fung & M. Poesio (Eds.) *Proceedings of the 51st annual meeting of the association for computational linguistics, association for computational linguistics* (pp. 1691–1701), runner up Best Paper Award.
- Köhn, A., Stegen, F., & Baumann, T. (2016). Mining the spoken wikipedia for speech data and beyond. In N. Calzolari, K. Choukri, T. Declerck, S. Goggi, M. Grobelnik, B. Maegaard, J. Mariani, H. Mazo, A. Moreno, J. Odiijk & S. Piperidis (Eds.) *Proceedings of the tenth international conference on Language Resources and Evaluation (LREC 2016)*, European Language Resources Association (ELRA), Paris, France.
- Noferesti, S., & Shamsfard, M. (2016). Using data mining techniques for sentiment shifter identification. In N. Calzolari, K. Choukri, T. Declerck, S. Goggi, M. Grobelnik, B. Maegaard, J. Mariani, H. Mazo, A. Moreno, J. Odiijk & S. Piperidis (Eds.) *Proceedings of the tenth international conference on Language Resources and Evaluation (LREC 2016)*, European Language Resources Association (ELRA), Paris, France.

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.