

Bridging Boundaries: Combining Phonetic and Orthographic Information to Improve Automated Syllabification Performance

Gus Lathouwers, Wieke Harmsen, Catia Cucchiarini, Helmer Strik

Radboud University

Syllabification concerns the task of dividing words into syllables. Due to many exceptions and subword pattern interactions, training an algorithm to perform syllabification with high accuracy remains a challenge. Different syllabification algorithms have been put forth over the past few decades in the literature, both language-specific and language-independent. Syllabification algorithms can be applied to orthographic representations of words, as well as phonetic representations, and may aid in applications such as text-to-speech and spelling correction software. Given that research on Dutch syllabification algorithms is generally outdated or algorithms are not tailored to Dutch-specific language features, our research set out to apply modern deep-learning techniques for improved syllabification performance. Previously, syllabification algorithms have been applied to phonetic wordsets (e.g., Krantz et al., 2019), and orthographic wordsets (e.g., Trogkanis & Elkan, 2010); yet the two approaches have not been combined to complement each other.

A new deep-learning model was developed that combines orthographic and phonetic information from two independently trained neural nets into a unified deep-learning model using attention mechanisms. Results show that the integration of phonetic in addition to orthographic information in the deep learning model yields improvements. The mean word accuracy of 99.65% is a 0.10% improvement in comparison with the model trained solely on orthographic data, and a 0.14% improvement in comparison with the best model reported in the literature for Dutch orthographic syllabification (Trogkanis & Elkan, 2010). A similar approach using a transformer model applied to the English language achieved a 97.49% word accuracy, representing a 1.18% improvement over the orthographic-only model.

The outcome of the current research indicates that combining phonetic and orthographic information leads to increased accuracy on word processing tasks such as syllabification.

References

- Trogkanis, N., & Elkan, C. (2010). Conditional random fields for word hyphenation. In J. Hajič, S. Carberry, S. Clark, & J. Nivre (Eds.), *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics* (pp. 366–374). Association for Computational Linguistics.
- Krantz, J., Dulin, M., & De Palma, P. (2019). Language-Agnostic Syllabification with Neural Sequence Labeling. *International Conference on Machine Learning and Applications*.