# Phonemizer: Text to Phones Transcription for Multiple Languages in Python

**Mathieu Bernard**[1] **and Hadrien Titeux**[1]

**1** LSCP/ENS/CNRS/EHESS/Inria/PSL Research University, Paris, France

## Summary

Phones are elementary sounds the speech is made of, on which syllables and words are built. The transcription of texts from their orthographic form into a phonetic alphabet is an important requirement in various applications related to speech and language processing, for instance for text to speech systems. Phonemizer is a Python package addressing precisely this issue: it transcribes a text from its orthographic representation into a phonetic one. The package is user-friendly and exposes a single `phonemize` function, also available as a command-line interface. It supports about a hundred different languages and provides end-user functionalities such as punctuation preservation, phones accentuation, tokenization at phone/syllable/word levels, as well as parallel processing of large input texts.

## Statement of Need

Whereas the high-level features introduced above are implemented directly by `phonemizer`, the phonetic transcription itself is delegated to third party backends, wrapped in an homogoneous interface by the package. The default backend used by `phonemizer` is eSpeak (Dunn & Vitolins, 2019), a text to speech software built on linguistic expertise and hand written transcription rules. It transcribes text into the International Phonetic Alphabet and supports more than a hundred languages. Using MBROLA voices (Tits & Vitolins, 2019), available for 35 languages, the eSpeak backend transcribes text in the SAMPA computer readable phonetic alphabet. Festival (Black et al., 2014) is another text to speech software used as a backend for `phonemizer`. It is available for American English only, and uses a non standard phoneset for transcription, but this backend is the only one to meet the requirement of some applications by preserving syllable boundaries. The third `phonemizer` backend is Segments (Forkel et al., 2019), a Python package providing Unicode Standard tokenization routines and orthography segmentation. It relies on a grapheme to phone mapping to generate the transcription. This backend is mostly usefull for low-resource languages, for which users with linguistic expertise can write their own mappings. Six languages are provided as exemples with `phonemizer`: Chintang, Cree, Inuktitut, Japanese, Sesotho and Yucatec.

Text to phones transcription is a critical need in different applications related to natural language and speech processing. So far, the `phonemizer` package has been used in the preprocessing pipeline of various deep learning text to speech systems (Ideas Engineering, 2021; Mozilla, 2021; Watanabe et al., 2018). It has also been used as a preprocessing step in word segmentation studies regarding the role of speech prosody in segmentability (Ludusan et al., 2017) and the psychology of child development (Bernard et al., 2020; Cristia et al., 2019). A phonetic transcription generated by the package was used to evaluate a phone discrimination task for the Zero Speech Challenge 2017 (Dunbar et al., 2017). Finally, the `phonemizer` is very suitable to prepare datasets for their use with the Kaldi speech recognition

41 toolkit (Povey et al., 2011), where a phonetic transcription of text is a requirement for various
42 algorithms.

# Acknowledgements

# References

53 Bernard, M., Thiolliere, R., Saksida, A., Loukatou, G. R., Larsen, E., Johnson, M., Fibla, L.,
54 Dupoux, E., Daland, R., Cao, X. N., & Alejandrina, C. (2020). WordSeg: Standardizing
55 unsupervised word form segmentation from text. *Behavior Research Methods*, *52*(1),
56 264–278.

57 Black, A. W., Clark, R., Richmond, K., Yamagishi, J., Oura, K., & King, S. (2014). *The*
58 *festival speech synthesis system* (Version 2.4). CSTR, University of Edinburgh. https:
59 //www.cstr.ed.ac.uk/projects/festival

60 Cristia, A., Dupoux, E., Ratner, N. B., & Soderstrom, M. (2019). Segmentability differences
61 between child-directed and adult-directed speech: A systematic test with an ecologically
62 valid corpus. *Open Mind*, *3*, 13–22.

63 Dunbar, E., Cao, X. N., Benjumea, J., Karadayi, J., Bernard, M., Besacier, L., Anguera, X.,
64 & Dupoux, E. (2017). The zero resource speech challenge 2017. *2017 IEEE Automatic*
65 *Speech Recognition and Understanding Workshop (ASRU)*, 323–330.

66 Dunn, R. H., & Vitolins, V. (2019). eSpeak NG speech synthetizer. In *GitHub repository*
67 (Version 1.50). GitHub. https://github.com/espeak-ng/espeak-ng

68 Forkel, R., Moran, S., List, J.-M., Greenhill, S. J., Ashby, L., Gorman, K., & Kaiping, G.
69 (2019). *Cldf/segments: Unicode standard tokenization* (Version v2.1.3). Zenodo. https:
70 //doi.org/10.5281/zenodo.3549784

71 Ideas Engineering. (2021). Non-autoregressive transformer based neural network for text-to-
72 speech. In *GitHub repository*. GitHub. https://github.com/as-ideas/TransformerTTS

73 Ludusan, B., Mazuka, R., Bernard, M., Cristia, A., & Dupoux, E. (2017). The role of
74 prosody and speech register in word segmentation: A computational modelling perspective.
75 *Proceedings of the Association for Computational Linguistics*, 178–183.

76 Mozilla. (2021). Deep learning for text to speech. In *GitHub repository*. GitHub. https:
77 //github.com/mozilla/TTS

78 Povey, D., Ghoshal, A., Boulianne, G., Burget, L., Glembek, O., Goel, N., Hannemann, M.,
79 Motlicek, P., Qian, Y., Schwarz, P., Silovsky, J., Stemmer, G., & Vesely, K. (2011). The
80 kaldi speech recognition toolkit. *IEEE 2011 Workshop on Automatic Speech Recognition*
81 *and Understanding*.

82 Tits, N., & Vitolins, V. (2019). MBROLA. In *GitHub repostory* (Version 3.3). GitHub.
83 https://github.com/numediart/MBROLA

Watanabe, S., Hori, T., Karita, S., Hayashi, T., Nishitoba, J., Unno, Y., Enrique Yalta Soplin, N., Heymann, J., Wiesner, M., Chen, N., Renduchintala, A., & Ochiai, T. (2018). ESPnet: End-to-end speech processing toolkit. *Proceedings of Interspeech*, 2207–2211. https://doi.org/10.21437/Interspeech.2018-1456