

# A Japanese Dataset and Efficient Multilingual LLM-Based Methods for Lexical Simplification and Lexical Complexity Prediction

Adam Nohejl<sup>1</sup> Akio Hayakawa<sup>2</sup> Yusuke Ide<sup>1</sup> Taro Watanabe<sup>1</sup>

<sup>1</sup>Nara Institute of Science and Technology

<sup>2</sup>Universitat Pompeu Fabra

{nohejl.adam.mt3, ide.yusuke.ja6, taro}@is.naist.jp

akio.hayakawa@upf.edu

## Publication Information

Volume 32, issue 4, pp. 1129–1188.

DOI: <https://doi.org/10.5715/jnlp.32.1129>

## Abstract

Lexical simplification (LS) is the task of making text easier to understand by replacing complex words with simpler equivalents. LS involves the subtask of lexical complexity prediction (LCP). We present MultiLS-Japanese, the first unified LS and LCP dataset targeting non-native Japanese speakers, and one of the ten language-specific MultiLS datasets. We propose methods for LS and LCP based on large language models (LLMs) that outperform existing LLM-based methods on 7 and 8 of the 10 MultiLS languages, respectively, while using only a fraction of their computational cost. Our methods rely on a single prompt across languages and introduce a novel calibrated token-probability scoring technique, G-SCALE, for LCP. Our ablations confirmed the benefits of G-SCALE and of concrete wording in the LLM prompt. We made the MultiLS-Japanese dataset available online under a CC-BY-SA license, including detailed metadata.<sup>1)</sup>

---

1) The complete dataset is made available at <https://huggingface.co/datasets/naist-nlp/multils-japanese>