

# Perbandingan Algoritma Sentencepiece BPE dan Unigram Pada Tokenisasi Artikel Bahasa Indonesia

Triwidyastuti Jamaluddin<sup>1</sup>, Moch Arif Bijaksana<sup>2</sup>, Ibnu Asror<sup>3</sup>

<sup>1,2,3</sup>Fakultas Informatika, Universitas Telkom, Bandung

<sup>1</sup>triwidyastuti@student.telkomuniversity.ac.id, <sup>2</sup>arifbijaksana@telkomuniversity.ac.id,

<sup>3</sup>iasror@telkomuniversity.ac.id

---

## Abstrak

Tokenisasi merupakan sebuah konsep yang mencakup proses sederhana dimana urutan teks dipecah menjadi bagian-bagian yang lebih kecil atau token dan kemudian dimasukkan sebagai input ke dalam model Natural language processing (NLP), atau proses model yang lebih kompleks seperti menerapkan pengetahuan dunia Deep Learning (DL). Tokenisasi akan lebih rumit ketika berhadapan dengan kasus semua kata dikelompokkan menjadi satu token atau tanpa pemisah dan kesalahan dalam tipografi. Paper ini mengusulkan model unsupervised tokenization menggunakan subword unit tokenizer dan detokenizer representasi oleh neural network, implementasi algoritma Byte Pair Encoding (BPE) dan Unigram Language Model. Selain itu, mengeksplorasi sentencePiece, model segmentasi pada kalimat dapat dilatih tanpa spasi. Eksperimen menggunakan bahasa Indonesia menghasilkan akurasi 54.6% dan 87.0% untuk Byte Pair Encoding (BPE) dan Unigram Language Model, masing-masing.

**Kata kunci:** *SentencePiece, Subword Tokenizer, Byte Pair Encoding (BPE), Unigram Language Model.*