
Neural Cross-Lingual Tokenization: A Morphologically-Aware Approach to Multilingual Subword Segmentation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Current multilingual language models rely on frequency-based tokenization
2 methods like Byte Pair Encoding (BPE) and SentencePiece, which treat
3 languages as mere symbol sequences without leveraging linguistic structure.
4 While these approaches have achieved remarkable success through
5 scale, they demonstrate significant limitations for morphologically complex
6 languages, leading to over-segmentation, poor cross-lingual transfer, and
7 suboptimal compression ratios. We propose a neural multilingual tokeniza-
8 tion framework that learns morphologically-aware subword representations
9 through cross-lingual consistency constraints and dynamic programming-
10 based segmentation. Our approach addresses fundamental limitations of
11 existing tokenizers by learning semantic equivalences between morphemes
12 across languages, rather than relying solely on frequency statistics. We
13 want to demonstrate that this method achieves superior compression ra-
14 tios and cross-lingual transfer performance, particularly for agglutinative
15 and morphologically rich languages that constitute over 40% of the world’s
16 linguistic diversity.

17 **Keywords:** multilingual NLP, tokenization, morphology, cross-lingual transfer, subword seg-
18 mentation

19 **1 Introduction**

20 Tokenization serves as the fundamental preprocessing step in modern natural language pro-
21 cessing pipelines, directly influencing model efficiency, linguistic fidelity, and cross-lingual
22 transfer capabilities. The dominant paradigm in multilingual language models employs
23 frequency-based subword tokenization methods such as Byte Pair Encoding (BPE) Sen-
24 nrich et al. [2015] and SentencePiece Kudo and Richardson [2018], which have enabled the
25 development of powerful multilingual models like mBERT Devlin et al. [2019] and XLM-R
26 Conneau et al. [2019].

27 However, these statistical approaches treat natural language as mere symbol sequences,
28 ignoring the rich morphological and semantic structure that characterizes human language.
29 This limitation becomes particularly pronounced for morphologically complex languages,
30 where a single root can generate hundreds of inflected forms through agglutination, fusion, or
31 other morphological processes. Recent empirical evidence suggests that current multilingual
32 models exhibit systematic performance gaps between morphologically simple languages (e.g.,
33 English) and complex ones (e.g., Turkish, Finnish, Arabic) Jabbar [2023].

34 This paper introduces a neural multilingual tokenization framework that addresses these
35 fundamental limitations through three key innovations:

- 36 1. Morphologically-aware representation learning that captures cross-lingual morpho-
37 logical equivalences
- 38 2. Convex optimization-based segmentation that optimizes for both compression and
39 linguistic coherence
- 40 3. Typological clustering that enables zero-shot support for rare languages through
41 linguistic prototype transfer

42 Our approach moves beyond frequency-based heuristics to learn principled subword seg-
43 mentations that respect morphological boundaries while maximizing cross-lingual transfer
44 potential.

45 2 Related Work and Current Limitations

46 2.1 Current State-of-the-Art

47 The current multilingual NLP landscape is dominated by large-scale transformer models
48 trained with frequency-based tokenization:

- 49 • AYA (Cohere AI) Abagyan et al. [2025]: Massive multilingual tokenization based
50 on byte-pair-encoding.¹
- 51 • XLM-R Conneau et al. [2019]: Trained on 2.5TB of CommonCrawl data across
52 100 languages, achieving significant improvements over mBERT on cross-lingual
53 benchmarks (+14.6% on XNLI, +13% on MLQA)
- 54 • mT5 Xue et al. [2020]: A multilingual text-to-text transformer covering 101 lan-
55 guages
- 56 • PaLM Chowdhery et al. [2023]: 540B parameter model demonstrating emergent
57 multilingual capabilities

58 dynamic tokenizers

59 theoretical findings

60 2.2 Multi-lingual datasets

61 FLORES-200, mC4, AYA, HPLT

62 2.3 Multi-lingual evaluation datasets

63 MasakhaNER 2.0, XTREME, XTREME-UP, afaik AYA has some instruct-tuning datasets
64 too, COSMMIC, HERCULE, Global MMLU (key dataset), MaXIFE.

65 2.4 Quantitative Evidence of Current Limitations

66 2.4.1 Morphological Segmentation Quality

67 Recent analysis reveals systematic failures in how current tokenizers handle morphologically
68 complex languages:

69 Turkish Analysis Kaya and Tantuğ [2024]: XLM-R SentencePiece creates morphologically
70 meaningless subwords. For example, Turkish “-bandela” is segmented into “ba”, “#ndel”,
71 “#a”, which lacks morphological coherence. Over-segmentation leads to 2-3× longer se-
72 quences compared to morphologically-aware alternatives.

¹Cohere has many of the most recent papers on multilingual tokenization:
<https://cohere.com/research/papers>

73 Arabic Challenges Tawfik et al. [2019]: Root-pattern morphology is poorly captured by
74 frequency-based methods, with dialectal variations creating inconsistent tokenization within
75 the same semantic space.

76 Japanese

77 2.4.2 Performance Gaps Across Language Families

78 Empirical evidence demonstrates systematic performance degradation for morphologically
79 complex languages:

Language Type	XNLI Accuracy	Relative Performance
English (Fusional)	91.3%	Baseline
Turkish (Agglutinative)	79.1%	-13.4%
Finnish (Agglutinative)	81.2%	-11.1%
Arabic (Root-pattern)	77.8%	-14.8%

Table 1: XLM-R performance across language types showing consistent degradation for morphologically complex languages.

80 Key findings Jabbar [2023]: Consistent 10-15% performance degradation for agglutinative
81 languages, with the gap persisting even with increased training data. Morphological align-
82 ment shows potential for 2-8% improvements.

83 2.4.3 Compression Efficiency Limitations

84 Current tokenizers demonstrate suboptimal compression for morphologically rich languages:

Language	Tokens/Word	Overhead
English	1.3	Baseline
Turkish	2.1	+61%
Finnish	1.9	+46%
Arabic	2.3	+77%

Table 2: Compression inefficiency directly translates to 30-70% more computational cost.

85 2.5 Fundamental Limitations of Frequency-Based Approaches

86 2.5.1 Lack of Semantic Awareness

87 BPE and SentencePiece optimize for:

$$\max \sum frequency(subword_i) \quad (1)$$

88 This objective ignores:

- 89 • Morphological boundaries
- 90 • Cross-lingual semantic equivalences
- 91 • Compositional meaning structure
- 92 • Corpora distribution (often towards English)

93 2.5.2 Language-Specific Bias

94 Training on concatenated multilingual corpora introduces systematic biases where high-
95 resource languages dominate vocabulary, morphologically simple languages receive dispro-
96 portionate representation, and cross-lingual patterns emerge accidentally rather than by
97 design.

- 98 2.5.3 Extra papers
- 99 1. <https://arxiv.org/pdf/2301.10472>
- 100 2. <https://arxiv.org/pdf/2403.10691>
- 101 3. <https://arxiv.org/pdf/2402.06619>
- 102 4. <https://cohere.com/research/papers>
- 103 5. <https://oldi.org/>
- 104 6. <https://arxiv.org/pdf/2305.11938>
- 105 7. <https://arxiv.org/abs/2506.15372>
- 106 8. <https://arxiv.org/abs/2405.07883>
- 107 9. <https://arxiv.org/pdf/2411.18553>
- 108 10. <https://arxiv.org/pdf/2412.15210>
- 109 11. <https://arxiv.org/pdf/2505.19599>
- 110 12. <https://arxiv.org/pdf/2410.13394>
- 111 13. <https://arxiv.org/pdf/2412.03304>
- 112 14. <https://arxiv.org/pdf/2503.10267>
- 113 15. <https://arxiv.org/pdf/2506.01776>
- 114 16. <https://arxiv.org/pdf/2003.11080>

115 References

- 116 Diana Abagyan, Alejandro R Salamanca, Andres Felipe Cruz-Salinas, Kris Cao, Hangyu
117 Lin, Acyr Locatelli, Marzieh Fadaee, Ahmet Üstün, and Sara Hooker. One tokenizer to
118 rule them all: Emergent language plasticity via multilingual tokenizers. arXiv preprint
119 arXiv:2506.10766, 2025.
- 120 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra,
121 Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann,
122 et al. Palm: Scaling language modeling with pathways. Journal of Machine Learning
123 Research, 24(240):1–113, 2023.
- 124 Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume
125 Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin
126 Stoyanov. Unsupervised cross-lingual representation learning at scale. arXiv preprint
127 arXiv:1911.02116, 2019.
- 128 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of
129 deep bidirectional transformers for language understanding. In Proceedings of the 2019
130 conference of the North American chapter of the association for computational linguistics:
131 human language technologies, volume 1 (long and short papers), pages 4171–4186, 2019.
- 132 Haris Jabbar. Morphpiece: A linguistic tokenizer for large language models. arXiv preprint
133 arXiv:2307.07262, 2023.
- 134 Yiğit Bekir Kaya and A Cüneyd Tantuğ. Effect of tokenization granularity for turkish large
135 language models. Intelligent Systems with Applications, 21:200335, 2024.
- 136 Taku Kudo and John Richardson. Sentencepiece: A simple and language indepen-
137 dent subword tokenizer and detokenizer for neural text processing. arXiv preprint
138 arXiv:1808.06226, 2018.
- 139 Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare
140 words with subword units. arXiv preprint arXiv:1508.07909, 2015.
- 141 Ahmed Tawfik, Mahitab Emam, Khaled Essam, Robert Nabil, and Hany Hassan Awadalla.
142 Morphology-aware word-segmentation in dialectal arabic adaptation of neural machine
143 translation. In Proceedings of the Fourth Arabic Natural Language Processing Workshop,
144 pages 11–17, 2019.

¹⁴⁵ Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant,
¹⁴⁶ Aditya Barua, and Colin Raffel. mt5: A massively multilingual pre-trained text-to-text
¹⁴⁷ transformer. arXiv preprint arXiv:2010.11934, 2020.