



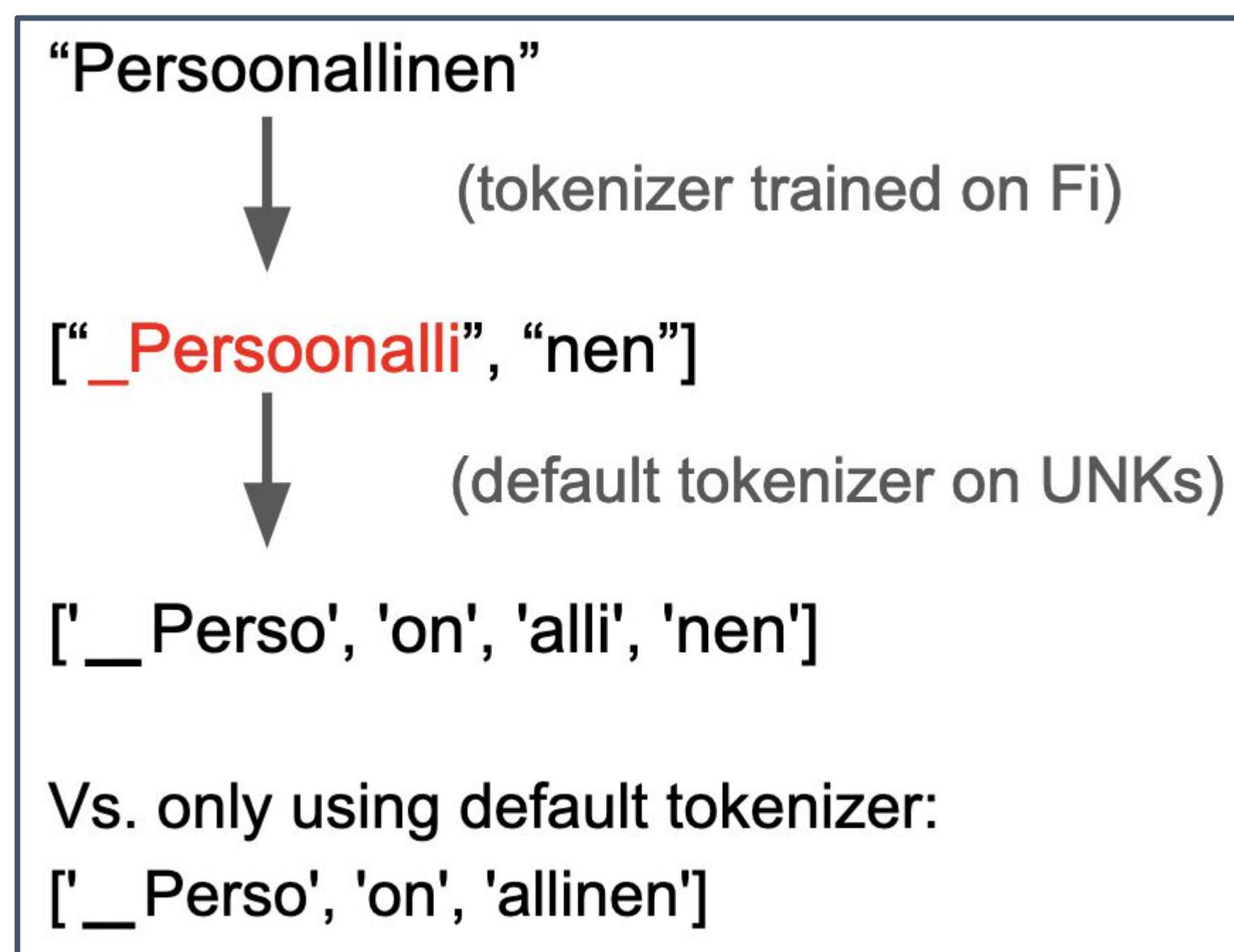
## BACKGROUND

### BYTE-PAIR ENCODING

- What's wrong with BPE?
  - Byte-pair encoding (BPE) replaces frequently occurring byte pairs with a new byte to create token vocabulary.
  - Balance between information of word-level tokenization and flexibility of character-level tokenization.
  - Related work shows BPE is not perfect:
    - BPE causes language imbalance.
    - Tokenizing spaces separately improves accuracy.
- How can we improve?
  - Language-specific BPE: Train BPE on dataset of balanced languages.
  - Morpheme tokenizer: Segment words into tokens based on morphemes.
    - Known improvement for Korean NMT [1].

## METHODOLOGY

### LANGUAGE SPECIFIC TOKENIZER

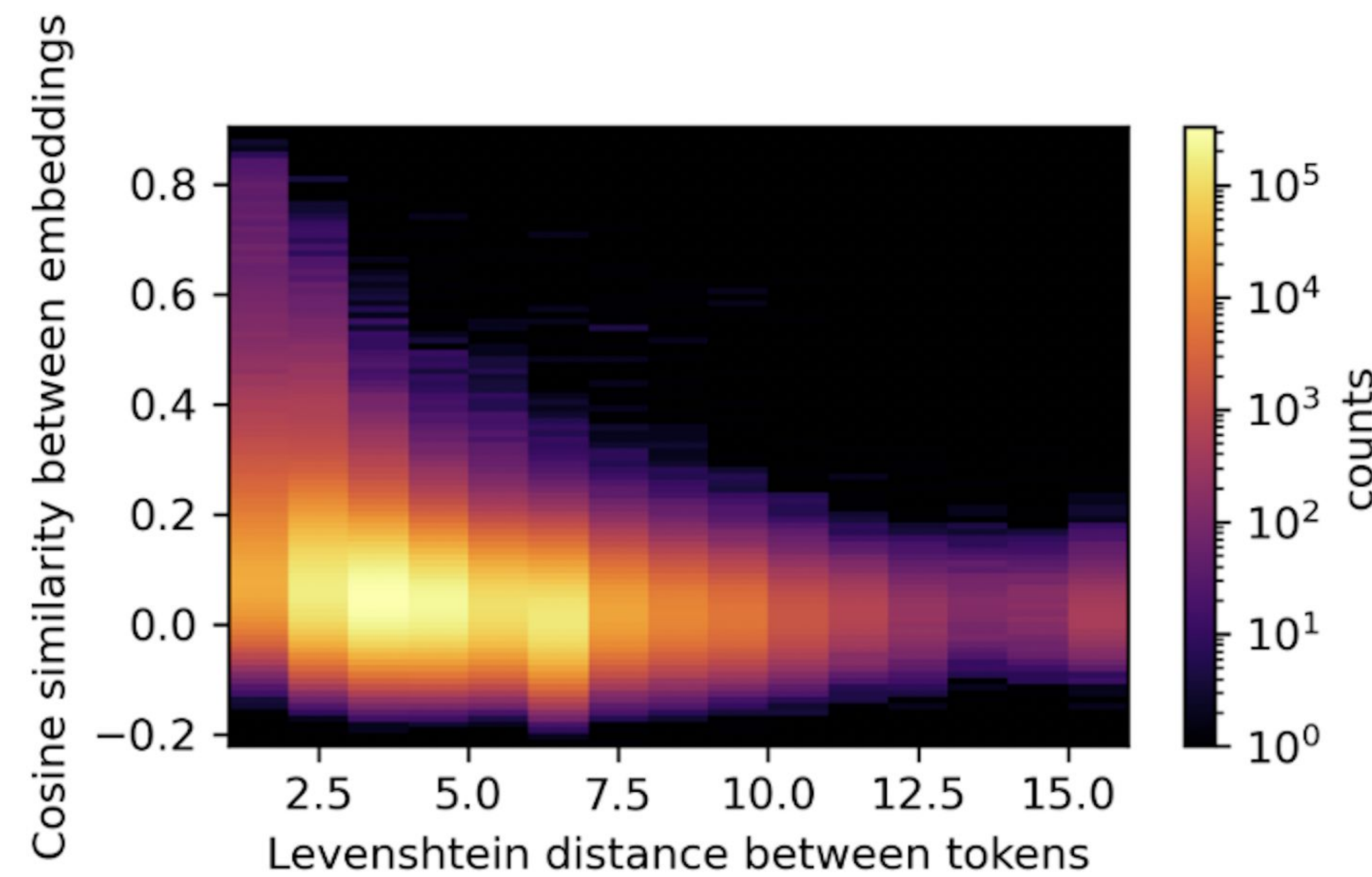


### MORPHEME TOKENIZER

- Greedy algorithm:
  - Tokenize the word according to tokenization provided in MorphyNet [2].
  - Split the word into the stem and the most frequent potential affix.
  - Use default tokenizer on UNKs.

## RESULTS

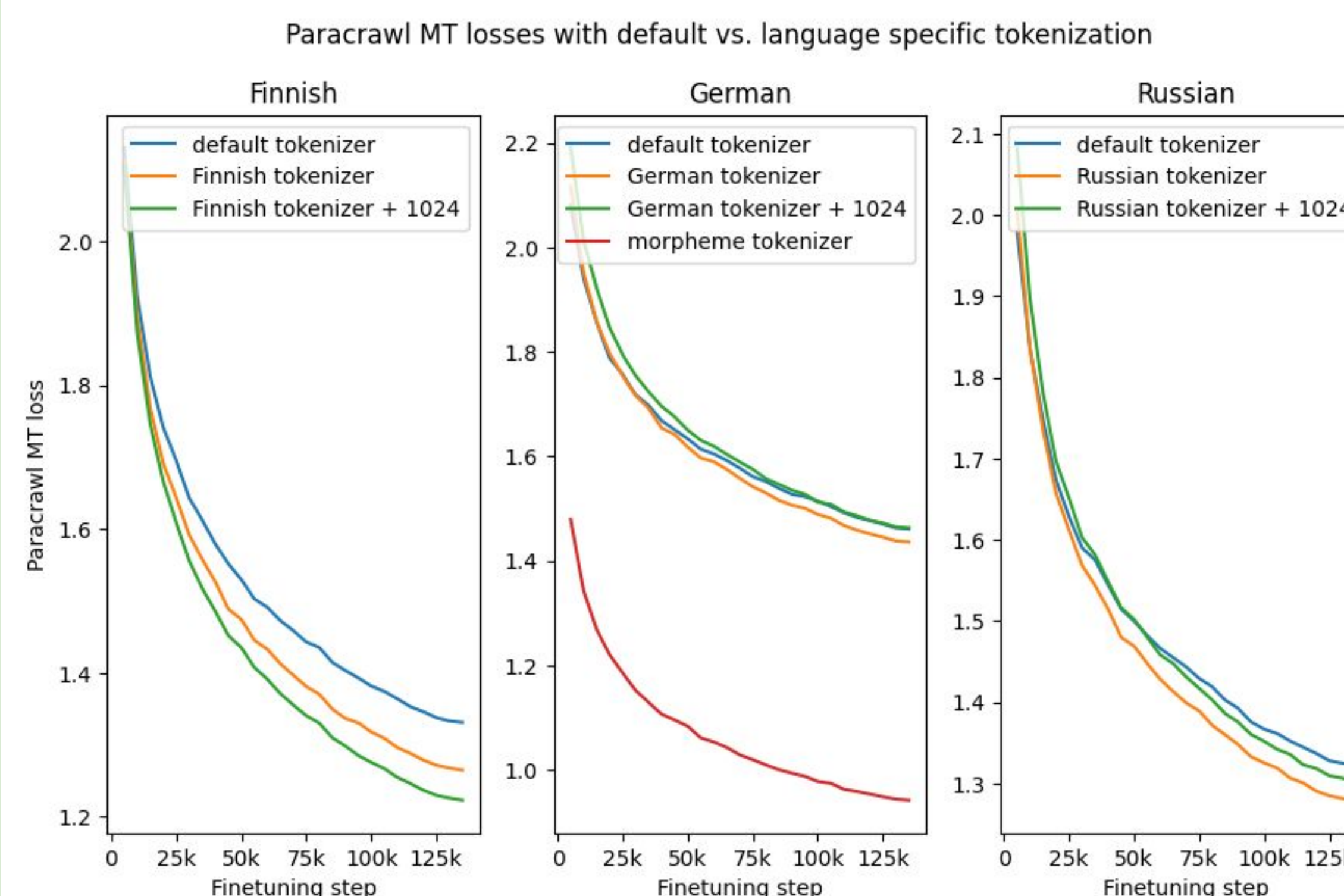
### TOKEN SIMILARITY VS. EDIT DISTANCE



### SUBSTITUTING NEW TOKENIZERS

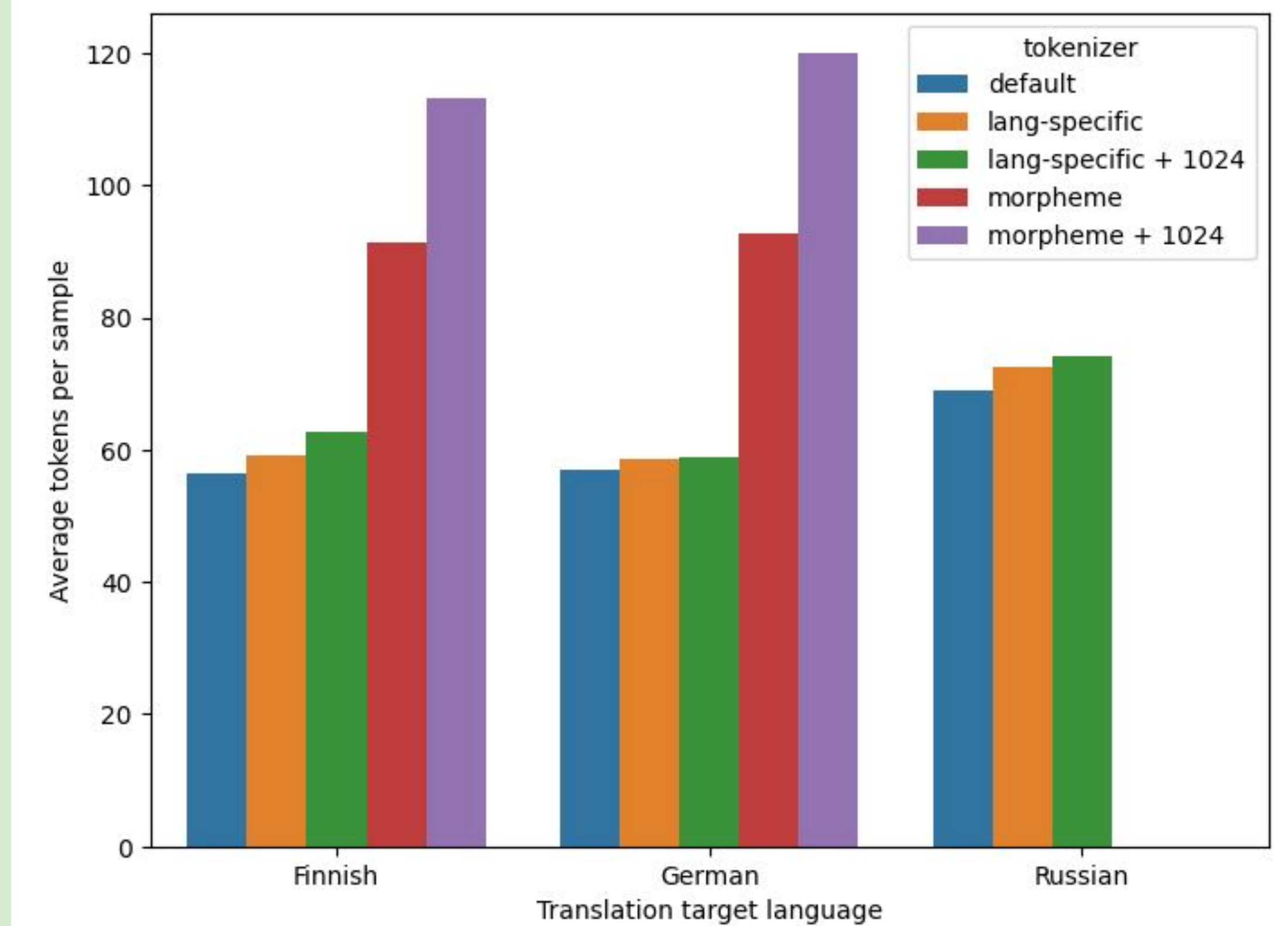
Finetuning Tokenizer	Finnish	German	Russian
Pretrained	20.80 (41.97)	18.77 (39.84)	13.37 (33.28)
Language-specific	20.18 (40.24)	18.68 (39.96)	12.97 (31.60)
Language-specific (Added tokens)	17.80 (38.39)		
Morpheme			N/A
Morpheme (Added tokens)			N/A

BLEU score and chrF++ for each strategy



## RESULTS CONT.

### TOKEN SEGMENTATION RATES



## DISCUSSION AND FUTURE WORK

### DISCUSSION

- Substituting linguistically-informed tokenizers hurt translation performance.
  - BLEU and chrF++ follow segmentation rates better than loss.

### LIMITATIONS

- Unsure whether decreased accuracy from alternative tokenizers is due to (i) distribution shift or (ii) linguistically informed tokenization not helping

### FUTURE WORK

- How to insert new tokens into vocabulary?
- Pre-train LLM from scratch using new tokenization strategies.
- Use a variety of more distantly-related languages.

### CITATIONS

[1] Kyubyong Park, Joohong Lee, Seongbo Jang, and Dawoon Jung. An empirical study of tokenization strategies for various Korean NLP tasks. 2020.  
[2] Khuyagbaatar Batsuren, Gábor Bella, and Fausto Giunchiglia. MorphyNet: a large multilingual database of derivational and inflectional morphology. In Garrett Nicolai, Kyle Gorman, and Ryan Cotterell, editors Proceedings of the 18th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology, pages 39–48. Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.sigmorphon-1.5. URL <https://aclanthology.org/2021.sigmorphon-1.5>.