

"Neural Machine Translation of Rare Words with Subword Units"

challenge → translating rare & unknown words in NMT systems. Standard NMT models typically operates on a fixed 30K to 50K words, → giving way to open vocabulary probably & facing OOV (out of vocab)

BPE algorithm → to tokenize for better translation
↳ many "words classes" can be translated via smaller units than full words.

BPE. Algorithm

- ① Initialization → symbol vocabulary initialized with character vocab. Special end of word symbol "." added to each word.
- ② Frequency Calculation → All adjacent symbol pairs in the training data are counted. Weighted according to frequency.
- ③ Identify most frequent pair of symbols.
- ④ Merging → each occurrence of the most frequent pair is replaced with a new symbol which represents the merged sequence of characters.
- ⑤ Each merge operation produces a new symbol → added to symbol vocabulary.
- ⑥ Steps 2-5 repeated iteratively.

for better efficiency \rightarrow pairs crossing word boundaries
are not considered. algorithm runs on dictionary
with word frequencies.