

Morphology Informed Selections for Subword Vocabulary Size

Anonymous ACL submission

Abstract

001 Currently, guidance around selection of an opti-
002 mal or appropriate subword vocabulary size is
003 incomplete and confusing at best. Using a mea-
004 sure of subword-morpheme overlap, our analy-
005 sis shows that one can find a "sweet spot" for
006 a morphology informed subword vocabulary
007 size. This sweet spot exhibits some variation
008 with respect to text complexity and the morpho-
009 logical characteristics of a language. However,
010 it is relatively constant with respect to corpus
011 size.

1 Introduction

012 It is now a best practice in neural machine transla-
013 tion (NMT) to encode input data using a subword
014 (e.g., byte-pair encoding, or BPE) vocabulary (Sen-
015 nrich et al., 2016; Tan et al., 2020). This encoding
016 enables open-vocabulary translation and limits the
017 size of a vocabulary corresponding to a large cor-
018 pus of text. BPE subword methods, for example,
019 iteratively replace the most frequent pair of charac-
020 ter or character sequences in a corpus with a single
021 new character sequence to generate a fixed size
022 vocabulary of subwords capable of tokenizing the
023 corpus. A practitioner can thus specify the size of
024 the subword vocabulary.

025 Currently, guidance around selection of an opti-
026 mal or appropriate subword vocabulary size is in-
027 complete and confusing at best. Certain researchers
028 propose simple heuristics based on NMT experi-
029 ments in certain select languages (Gowda and May,
030 2020). Others recommend performing a sweep
031 over subword vocabulary sizes (Ding et al., 2019)
032 or other computationally intensive trial and error
033 methods to select subword vocabulary size. Still
034 others suggest that specific numbers of BPE merges
035 exhibit similarities across languages, which could
036 motivate consistent choices for subword vocabulary
037 size in a multilingual context (Gutierrez-Vasques
038 et al., 2021).

040 In this work, we add another perspective and
041 attempt to bring some clarity to the selection of
042 subword vocabulary sizes for NMT and other Nat-
043 ural Language Processing (NLP) experiments. The
044 usage of subword vocabularies is most often moti-
045 vated by a desire to enable open vocabulary meth-
046 ods while, at the same time, limiting vocabulary
047 size for the purpose of corpus tokenization. Thus,
048 capturing the true vocabulary of the corpus, or
049 rather morphology of the language, is still the
050 end goal of such approaches. In this paper, we
051 show that the overlap between subwords and mor-
052 phemes follows a predictable pattern as a function
053 of subword vocabulary size (at least for corpora
054 in a certain domain and of a certain size). That
055 is, a subword vocabulary size can be predictably
056 selected based on the criteria that the subword vo-
057 cabulary should have a maximum overlap with a
058 corresponding language morphology. We calculate
059 such overlaps for 27 languages to motivate practi-
060 tioners to consider and experiment with morphol-
061 ogy informed selections of subword vocabulary
062 size.

2 Related Work

063 Various attempts have been made to identify opti-
064 mal subword vocabulary sizes. Gowda and May
065 (2020) perform systematic NMT experiments on
066 four different target languages. In these experi-
067 ments they use a range of BPE vocabulary between
068 500 and 64K types. They finally make a recom-
069 mendation for using a simple heuristic to identify
070 the near-optimal vocabulary size, which is where a
071 mean sentence length measure is small (low num-
072 bers of subwords per sentence) and the frequency
073 of subwords in the corpus at the 95% class rank
074 is 100 or higher. Although this gives some clear
075 guidance, such an approach relies on the ability to
076 segment text into sentences (which is not always
077 an easy task for many languages in the world).

078 In other NMT experiment informed research,
079

080 [Denkowski and Neubig \(2017\)](#) make a general
081 recommendation of 32K BPE types for NMT sys-
082 tems with a secondary recommendation of 16K for
083 system with less than 1 million parallel sentences.
084 [Ding et al. \(2019\)](#), on the other hand, conduct ex-
085 periments with 5 different NMT architectures on
086 4 language pairs and come to the conclusion that
087 a sweep over BPE merge operations from 0-4K or
088 even 0-32K types is useful. This shows how re-
089 sults from NMT-based studies can vary (or even
090 contradict). Further, the authors are not aware of
091 any such works that use language morphology to
092 motivate the selection of subword vocabulary size.

093 In a different vein of research, [Gutierrez-Vasques](#)
094 [et al. \(2021\)](#) utilize information theory as a tool to
095 explore BPE merges. At each merge, they ana-
096 lyze Shannon entropy across 47 languages. This
097 entropy across subword distributions shows a lack
098 of variability, which suggests that a language that
099 is complex at the word level is not as complex at
100 the subword level. However, the "turning point"
101 highlighted in [Gutierrez-Vasques et al. \(2021\)](#) is at
102 around 200 BPE merges, which is significantly less
103 than the number of BPE merges generally recom-
104 mended in practice for NMT and other NLP exper-
105 iments. This discrepancy, along with the inconsis-
106 tency of NMT-based studies of subword vocabulary
107 size, begs the question: is there a linguistically in-
108 formed way to provide guidance on vocabulary size
109 selections that is more consistent with published
110 NMT studies?

111 3 Methodology

112 To provide linguistically informed guidance on the
113 selection of subword vocabulary sizes, we analyze
114 the overlap of subwords and morphemes for a range
115 of subword vocabulary sizes and for a variety of
116 languages. We look for a "sweet spot" where the
117 overlap between subwords and morphemes is a
118 maximum.

119 For each language, we pre-process the available
120 data to remove blank lines and to lower case all
121 characters. We then use the morphological analysis
122 implemented in the Python polyglot library ([Vir-
123 pioja et al., 2013](#)) to obtain morphemes for each
124 word contained in the corpus. Although this use of
125 polyglot limits the method to the 135 supported lan-
126 guages, vocabulary sizes for languages supported
127 by polyglot are likely a good starting point for vo-
128 cabulary sizes in related language experiments.

129 For a range of vocabulary sizes from 0 to 8000,

130 we first train a unigram subword model using that
131 vocabulary size. SentencePiece¹ is used for all
132 experiments, and all experiments define the same
133 random seed to maintain reproducibility. Next, we
134 encode each word in the corpus to get the unique
135 subwords corresponding to the word. These unique
136 subwords are compared with the corresponding
137 morphemes for that work to determine the percent-
138 age of these subwords that are also morphemes.
139 This percentage is what we define as the overlap
140 between subwords and morphemes. Because we
141 are doing this at a word level, we then aggregate
142 the overlap metrics for all words in a corpus to get
143 the average overlap for a given vocabulary size.

144 4 Experiments

145 In order to evaluate the behavior of subword-
146 morpheme overlap in many languages, we use data
147 from the JHU Bible Corpus ([McCarthy et al., 2020](#)).
148 We filtered all of the text files in the corpus down
149 to those that were: (i) supported by polyglot’s mor-
150 phological analysis; and (ii) including the full text
151 of the Bible. This resulted in 63 full Bibles in
152 27 languages including Latin and Cyrillic writing
153 system scripts. We ran the analysis detailed in
154 Section 3 on each of these full Bible files. The
155 reported maximum overlaps and vocabulary sizes
156 at the maximum overlaps are the averages for the
157 set of full Bibles in each respective language.

158 We also wanted to analyze the influence of cor-
159 pus size on the subword-morpheme overlap. To this
160 end, we took a single full Bible from 4 languages
161 (Polish [pol], French [fra], Vietnamese [vie], and
162 Romanian [ron]) and split the Bible into random
163 collections of Bible verses ranging in length from
164 100 verses to 30,000 verses. We then ran the anal-
165 ysis detailed in Section 3 on each of these collec-
166 tions. The result is a morphology optimal subword
167 vocabulary size as a function of the number of sam-
168 ples/lines in the respective corpus.

169 Finally, to analyze the influence of text com-
170 plexity and morphological extremes, we ran the
171 analysis detailed in Section 3 on two various spe-
172 cific pairings of Bible texts. The first of these pairs
173 was the Spanish [spa] Nueva Versión Internacional
174 (NVI) Bible and the Spanish Nueva Versión Inter-
175 nacional Simplificada (NVIs). The NVIs is a "sim-
176 plified" version of the NVI, and, thus, this pairing
177 should demonstrate how subword-morpheme over-
178 lap is influenced by text complexity. The second

¹<https://github.com/google/sentencepiece>

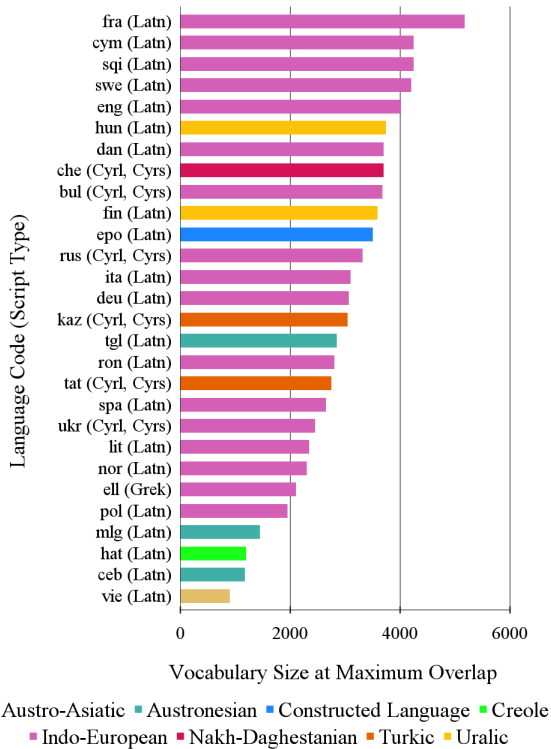


Figure 1: Subword vocabulary sizes for various languages at the maximum overlap of subwords and morphemes. Writing system scripts are shown in parentheses and language families are shown via color.

of these pairs was a Quechua [quc] Bible and the Bible in Basic English (BBE). This basic English text (BBE) should differ significantly, on a morphological level, as compared to the agglutinating language of Quechua.

5 Results and Discussion

Figure 1 shows the subword vocabulary size at the maximum overlap between subwords and morphemes for the 27 languages we considered. This vocabulary size ranges from 900 on the low end (for Vietnamese [vie]) to 5,175 at the high end (for French [fra]). The analysis finds higher vocabulary sizes for corpora with a large number of morphemes (like French, with 4300+ morphemes in the corpus) and lower vocabulary sizes for corpora with a small number of morphemes (like Vietnamese, with only 864 morphemes found in the corpus). In fact, the lowest vocabulary size, for Vietnamese, is consistent with the fact that Vietnamese is known to be an extreme in Austro-Asiatic languages in that it has very little morphology (Noyer, 1998).

To see how the subword-morpheme overlap changes as a function of subword vocabulary size,

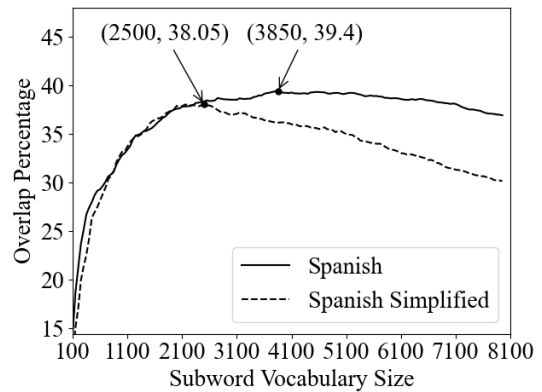


Figure 2: The percentage overlap between subwords and morphemes as a function of subword vocabulary size for the Nueva Versión Internacional Bible ("Spanish" in the figure) and the Nueva Versión Internacional Simplificada Bible ("Spanish Simplified" in the figure).

see Figure 2 and Figure 3. All of these curves indicate a general trend: the subword-morpheme overlap rises gradually to a peak and then starts to decrease. This trend makes sense in terms of the subword merges. Starting from characters at the low end of vocabulary size, the overlap with rise gradually as these characters and sets of characters are merged into morphemes. However, eventually the merges will start merging two morphemes or a non-morpheme set of characters with a morpheme. These latter merges with decrease the overall overlap.

Of course the shape of the subword-morpheme overlap curve will change as a function of both text complexity (influencing the total number of morphemes) and morphological characteristics of a language (influencing the number of subwords per word and per sentence). Figure 2 shows that the overlap curve for a simplified Spanish corpus (the NVIs) peaks earlier than the corresponding non-simplified version. Figure 3 shows how the long words of the agglutinating language of Quechua cause the overlap curve to peak earlier and decrease more rapidly than the curve for the English BBE translation. This is because Quechua has both fewer total morphemes and longer words than English.

Finally, Figure 4 shows how the vocabulary size at the maximum subword-morpheme overlap changes as a function of corpus size. This variability over corpus size is shown for a subset of the languages represented in Figure 1. One can see that the vocabulary size rises quickly to a relatively

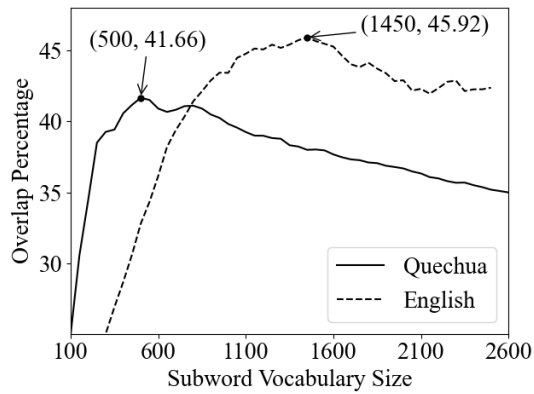


Figure 3: The percentage overlap between subwords and morphemes as a function of subword vocabulary size for the Quechua [quc] and the Bible in Basic English (BBE).

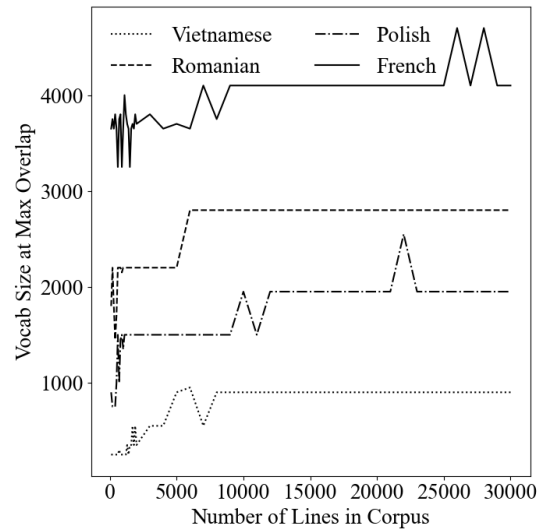


Figure 4: The "sweet spot" vocabulary size (at maximum subword-morpheme overlap) as a function of corpus size for 4 languages (Polish [pol], French [fra], Vietnamese [vie], and Romanian [ron]). The vocabulary size at the maximum overlap stays relatively constant (i.e., it plateaus) after the corpus size reaches around 10-15K samples. From vocabulary sizes from 100 to 2000 we use an interval of 100, and from 2000+ we sample with an interval size of 1000.

constant value as a function of corpus size. This suggests that, if a practitioner finds a morphology informed choice of subword vocabulary size (at least for unigram subwords), the choice of vocabulary size can be re-used for experiments with a variety of corpus sizes. In fact, a practitioner could look at the language represented in Figure 1, find a language similar to the language they are using in their experiments, and select a vocabulary size similar to that of the related languages. Related languages can be found using language classifications, such as those in the Ethnologue (David M. Eberhard et al., 2021). Such a process may be a good starting point for vocabulary size selections.

6 Conclusions and Future Work

Using a measure of subword-morpheme overlap, our analysis shows that one can find a "sweet spot" for a morphology informed subword vocabulary size. For Bible data, this vocabulary size shows little variation with corpus sizes greater than 15,000 samples, although it does exhibit some variation with respect to text complexity and general morphological characteristics of a language. We acknowledge that the results presented here are very limited in terms of domain (the Bible) and this kind of subword-morpheme analysis may produce different results in other domains or with different corpus sizes. In any event, we submit that such a morphology informed analysis could serve as a starting point for vocabulary size in NMT or other NLP experiments. In future work, we would like to more fully explore data from other domains and the variation in downstream NLP task performance

with morphology informed vocabulary sizes.

References

- David M. Eberhard, Gary F. Simons, and Charles D. Fennig. 2021. *Ethnologue: Languages of the World*, twenty-fourth edition. SIL International, Dallas, Texas.
- Michael Denkowski and Graham Neubig. 2017. [Stronger baselines for trustable results in neural machine translation](#). In *Proceedings of the First Workshop on Neural Machine Translation*, pages 18–27, Vancouver. Association for Computational Linguistics.
- Shuoyang Ding, Adithya Renduchintala, and Kevin Duh. 2019. A call for prudent choice of subword merge operations in neural machine translation. In *MTSummit*.
- Thamme Gowda and Jonathan May. 2020. Finding the optimal vocabulary size for neural machine translation. *arXiv preprint arXiv:2004.02334*.
- Ximena Gutierrez-Vasques, Christian Bentz, Olga Sozinova, and Tanja Samardzic. 2021. [From characters to words: the turning point of BPE merges](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 3454–3468, Online. Association for Computational Linguistics.

294 Arya D. McCarthy, Rachel Wicks, Dylan Lewis, Aaron
295 Mueller, Winston Wu, Oliver Adams, Garrett Nicolai,
296 Matt Post, and David Yarowsky. 2020. [The Johns](#)
297 [Hopkins University Bible corpus: 1600+ tongues](#)
298 [for typological exploration](#). In *Proceedings of the*
299 *12th Language Resources and Evaluation Confer-*
300 *ence*, pages 2884–2892, Marseille, France. European
301 Language Resources Association.

302 Rolf Noyer. 1998. Vietnamese ‘morphology’ and the def-
303 inition of word. *University of Pennsylvania Working*
304 *Papers in Linguistics*, 5(2):5.

305 Rico Sennrich, Barry Haddow, and Alexandra Birch.
306 2016. Neural machine translation of rare words with
307 subword units. *ArXiv*, abs/1508.07909.

308 Zhixing Tan, Shuo Wang, Zonghan Yang, Gang Chen,
309 Xuancheng Huang, Maosong Sun, and Yang Liu.
310 2020. Neural machine translation: A review of meth-
311 ods, resources, and tools. *AI Open*, 1:5–21.

312 Sami Virpioja, Peter Smit, and Stig-Arne Grönroos and
313 Mikko Kurimo. 2013. Morfessor 2.0: Python imple-
314 mentation and extensions for morfessor baseline. In
315 *Aalto University publication series*. Department of
316 Signal Processing and Acoustics, Aalto University.