

---

# Revisiting Syllable-aware Language Modelling

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Language modelling (LM) is regularly analysed at word, sub-word or character  
2 level inputs, and this study reconsiders syllable units for the task. Rule-based  
3 syllabification typically requires less specialised knowledge than identifying mor-  
4 phemes, and the process can naturally work for low-resource cases, as we do not  
5 need an unsupervised model to extract sub-words. In this paper, we compare differ-  
6 ent granularities from characters to words in an open-vocabulary LM task, where  
7 syllables mostly outperform the rest of them for both English and Shipibo-Konibo  
8 languages. Thereafter, we obtain similarly positive results for syllable-level neural  
9 machine translation (NMT) with Spanish too. [All authors identify as Latinx]

## 10 1 Introduction

11 Previous work on syllable-aware LM in English failed to beat character-level models [2]; however, we  
12 propose to assess the task under two new settings. First, we could employ a plain-vanilla architecture,  
13 without additional composition functions, to analyse an open-vocabulary scenario with syllables [3].  
14 Second, English has a weak correspondence between graphemes (written symbols) and phonemes  
15 (speech units), so we might include an study case with less-ambiguous splits. Therefore, we revisit  
16 syllable-aware LM by using simple recurrent neural networks [8] for open-vocabulary generation [15],  
17 and by also assessing a more phonetic language with a recent alphabetisation (Shipibo-Konibo [1]).  
18 We thereupon explore the syllables effect in another generation task such as NMT.

## 19 2 Methodology and Results

20 We evaluate syllables against words, Byte Pair Encoding [BPE, 14] sub-words, and characters, with a  
21 comparable perplexity [10] in LM; and character [18] and word level [13] metrics in NMT.

### 22 2.1 Languages and Datasets

23 For LM in English (*eng*), we use well-known datasets: Penn Treebank [PTB, 7] and WikiText-2 [9].  
24 In the case of Shipibo-Konibo (*shp*), a low-resource and native language from Peru, we process the  
25 monolingual side of three parallel corpora aligned with Spanish (*spa*) [4]. For one of them, named  
26 Flashcards, we align *eng* sentences from the original *eng-spa* corpus used for its creation [16]. For  
27 comparison purposes, we also analyse the new *eng* monolingual text of Flashcards in LM. Afterwards,  
28 we study the NMT case only with the Flashcards dataset in both *shp-eng* and *shp-spa* language-pairs.  
29 We segment syllables in *shp* with rules [1] and with a dictionary-based method [5] for *eng* and *spa*.

30 Table 1-a-b describes the data for LM. We observe a vast amount of syllable types in the *eng* datasets,  
31 in contrast to *shp*, where syllables are closer to characters than to other granularities. Moreover, the  
32 Flashcards segmentation reveals the perplexing nature of *eng* syllables. For the LM task on *shp*, the  
33 significantly low amount of unique syllables could be interpreted as modelling a language with a  
34 larger alphabet (more characters types) and a smaller average length of token.

Dataset		Ⓐ Split size			Ⓑ # types				Ⓒ $ppl^c \downarrow$			
		Train	Valid	Test	Word	Syl	Char	BPE	Word	Syl	Char	BPE (*)
<i>eng</i>	PTB	887.0k	70.3k	78.6k	10.0k	6.1k	48	4.7k	2.36	<b>2.11</b>	2.52	2.42 (5k)
	WikiText-2	103.2M	217.6k	245.5k	33.2k	19.5k	274	1.3k	2.62	<b>2.15</b>	2.72	2.63 (1k)
	Flashcards	14.7k	1.4k	1.7k	2.5k	2.4k	63	2.3k	<b>2.12</b>	2.24	3.01	2.62 (3k)
<i>shp</i>	Flashcards	12.1k	2.1k	1.4k	2.6k	193	30	1.8k	2.70	<b>2.39</b>	2.64	3.30 (3k)
	Religious	82.4k	9.4k	10.2k	11.1k	331	26	1.0k	3.01	<b>2.37</b>	2.48	2.92 (1k)
	Educational	32.0k	3.6k	4.1k	4.0k	258	32	2.8k	2.65	<b>2.16</b>	2.29	2.77 (3k)

Table 1: (a) Split size in tokens; (b) Number of types per segmentation in Train; (c)  $ppl^c$  on Test for LM. For BPE, we show the best score given various merges between 1k–5k with a 1k-step.

	BLEU $\uparrow$						CharacTER $\downarrow$					
	Word	Syl	Char	BPE 5k–10k–15k			Word	Syl	Char	BPE 5k–10k–15k		
<i>shp-eng</i>	16.26	18.38	<b>19.60</b>	16.90	15.65	16.21	63.86	<b>53.57</b>	54.25	56.20	58.71	57.51
<i>eng-shp</i>	16.35	<b>19.70</b>	17.32	16.61	16.80	17.17	57.07	<b>51.76</b>	53.40	55.61	55.80	56.91
<i>shp-spa</i>	8.91	<b>13.20</b>	10.62	8.68	8.76	9.14	68.37	<b>55.33</b>	58.98	62.20	63.00	64.61
<i>spa-shp</i>	9.76	<b>14.78</b>	13.39	11.62	11.62	12.12	65.79	<b>55.05</b>	55.24	62.48	63.42	62.88

Table 2: NMT results at word (BLEU) and character level (CharacTER) on the Flashcards dataset.

## 2.2 Language Modelling with a Comparable Perplexity

For a fair comparison across all granularities, we evaluate all results with character-level perplexity:  $ppl^c = \exp(L \cdot (s^{seg} + 1)/(s^c + 1))$ , where  $L$  is the cross-entropy loss of a string  $s$  computed by a neural LM, and  $s^{seg}$  and  $s^c$  refer to the length of  $s$  in the chosen segmentation and character level units, respectively [10]. Furthermore, we generate the same input unit as an open-vocabulary task, where there is no prediction of an “unknown” token [15], with an exception at word-level in PTB. We thereby differ from previous work [2], and refrain from composing the syllable representations into words to evaluate only word-level perplexity. Following other open-vocabulary LM studies [12, 11], we use a low-compute version of an LSTM neural network, named Average SGD Weight-Dropped [8], with a smaller embedding size (300 units) for faster training. Additionally, we use the SentencePiece [6] segmentation format in both characters and syllables, and the original one for BPE [14].

Table 1-c shows that syllables mostly result in better perplexities than the remaining granularities in LM, even for a low-phonetic language as *eng*, and with a very competitive score when they do not achieve the best one. Moreover, syllables outperform the rest in the open-vocabulary scope (excluding words). Beating characters implies a gain in time processing as well, given the shorter sequences of syllables. Other settings that could be further explored are working unsupervised morphemes or morphological-aided supervision [17], and constraining the BPE-vocabulary size to the number of syllable types [6].

## 2.3 Syllables for Neural Machine Translation (NMT)

To further explore the value of syllables, we build *eng-shp* and *spa-shp* NMT models with all granularities as inputs-outputs. Each model uses a two-layer LSTM encoder-decoder with a hidden layer of 512, an embedding size of 300, and joint BPE with various merges. Table 2 presents the BLEU [13] and CharacTER [18] scores, where syllables predominantly stand out again, with an exception against characters in *shp-eng* at the word level metric. This result reinforces the initial concern for the phonetic ambiguity of *eng*, as we infer an inherent difficulty to reconstruct a word with generated syllable sequences. We also hypothesise that the joint BPE models do not provide the best scores given the potentially small word and sub-word overlapping of *shp* with both *eng* and *spa*.

## 3 Conclusion

Our results suggest that syllables might be valuable for both open-vocabulary LM and NMT tasks, where they behave positively even for a poor phonemically-spelt language. Syllables do not have an embedded meaning; however, the required effort for their segmentation could be advantageous concerning other morphological-aware or unsupervised-driven methods. Finally, we are currently working on exploring a multilingual scope and the hybrid-LM scenario [12] with syllables.

68 **References**

- 69 [1] Carlo Alva and Arturo Oncevay. Spell-checking based on syllabification and character-level  
70 graphs for a Peruvian agglutinative language. In *Proc. of the First Workshop on Subword and*  
71 *Character Level Models in NLP*, pages 109–116. Association for Comp. Linguistics, 2017.
- 72 [2] Zhenisbek Assylbekov, Rustem Takhanov, Bagdat Myrzakhmetov, and Jonathan N. Washington.  
73 Syllable-aware neural language models: A failure to beat character-aware ones. In *Proceedings*  
74 *of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1866–  
75 1872. Association for Computational Linguistics, 2017.
- 76 [3] Ryan Cotterell, Sebastian J. Mielke, Jason Eisner, and Brian Roark. Are all languages equally  
77 hard to language-model? In *Proceedings of the 2018 Conference of the NAACL-HLT, Vol. 2*  
78 *(Short Papers)*, pages 536–541. Association for Computational Linguistics, 2018.
- 79 [4] Héctor E. Gómez Montoya, Kervy D. Rivas Rojas, and Arturo Oncevay. A continuous im-  
80 provement framework of machine translation for Shipibo-konibo. In *Proceedings of the 2nd*  
81 *Workshop on Technologies for MT of Low Resource Languages*, pages 17–23. EAMT, 2019.
- 82 [5] Kozea. Pyphen. Available in: <https://pyphen.org/>, 2013.
- 83 [6] Taku Kudo and John Richardson. SentencePiece: A simple and language independent subword  
84 tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference*  
85 *on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71.  
86 Association for Computational Linguistics, 2018.
- 87 [7] Mitchell P. Marcus, Beatrice Santorini, and Mary A. Marcinkiewicz. Building a large annotated  
88 corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2):313–330, 1993.
- 89 [8] Stephen Merity, Nitish Shirish Keskar, and Richard Socher. Regularizing and optimizing LSTM  
90 language models. In *International Conference on Learning Representations*, 2018.
- 91 [9] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture  
92 models. *arXiv preprint arXiv:1609.07843*, 2016.
- 93 [10] Sebastian J. Mielke. Can you compare perplexity across different segmentations? Available in:  
94 <http://sjmielke.com/comparing-perplexities.htm>, 2019.
- 95 [11] Sebastian J. Mielke, Ryan Cotterell, Kyle Gorman, Brian Roark, and Jason Eisner. What kind  
96 of language is hard to language-model? In *Proc. of the 57th Annual Meeting of the Association*  
97 *for Computational Linguistics*, pages 4975–4989. Association for Comp. Linguistics, 2019.
- 98 [12] Sebastian J. Mielke and Jason Eisner. Spell once, summon anywhere: A two-level open-  
99 vocabulary language model. In *The Thirty-Third AAAI Conf. on Artificial Intelligence*, 2019.
- 100 [13] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic  
101 evaluation of machine translation. In *Proc. of the 40th Annual Meeting of the Association for*  
102 *Computational Linguistics*, pages 311–318. Association for Computational Linguistics, 2002.
- 103 [14] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare  
104 words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for*  
105 *Computational Linguistics*, pages 1715–1725. Association for Computational Linguistics, 2016.
- 106 [15] Ilya Sutskever, James Martens, and Geoffrey Hinton. Generating text with recurrent neural  
107 networks. In *Proceedings of the 28th International Conference on International Conference on*  
108 *Machine Learning*, pages 1017–1024, 2011.
- 109 [16] Tatoeba. Tab-delimited bilingual sentence pairs. Available in: <http://www.manythings.org/anki/>, 2010. These are selected sentence pairs from the Tatoeba Project.
- 110
- 111 [17] Clara Vania and Adam Lopez. From characters to words to in between: Do we capture  
112 morphology? In *Proceedings of the 55th Annual Meeting of the Association for Computational*  
113 *Linguistics*, pages 2016–2027. Association for Computational Linguistics, 2017.
- 114 [18] Weiyue Wang, Jan-Thorsten Peter, Hendrik Rosendahl, and Hermann Ney. CharacTer: Transla-  
115 tion edit rate on character level. In *Proceedings of the First Conference on Machine Translation:*  
116 *Vol. 2, Shared Task Papers*, pages 505–510. Association for Computational Linguistics, 2016.