

Revisiting Syllables in Language Modelling and their Application on Low-Resource Machine Translation

Arturo Oncevay^{ε,χ} Kervy Dante Rivas Rojas^{ρ,α}

Liz Karen Chavez Sanchez^χ Roberto Zariquiey^{ρ,χ}

^εSchool of Informatics, University of Edinburgh, Scotland

^ρPontificia Universidad Católica del Perú (^αIA-PUCP | ^χChana Field Station), Peru

a.oncevay@ed.ac.uk, rzariquiey@pucp.edu.pe

Abstract

Language modelling and machine translation tasks mostly use subword or character inputs, but syllables are seldom used. Syllables provide shorter sequences than characters, require less-specialised extracting rules than morphemes, and their segmentation is not impacted by the corpus size. In this study, we first explore the potential of syllables for open-vocabulary language modelling in 21 languages. We use rule-based syllabification methods for six languages and address the rest with hyphenation, which works as a syllabification proxy. With a comparable perplexity, we show that syllables outperform characters and other subwords. Moreover, we study the importance of syllables on neural machine translation for a non-related and low-resource language-pair (Spanish–Shipibo–Konibo). In pairwise and multilingual systems, syllables outperform unsupervised subwords, and further morphological segmentation methods, when translating into a highly synthetic language with a transparent orthography (Shipibo–Konibo). Finally, we perform some human evaluation, and discuss limitations and opportunities.

1 Introduction

In language modelling (LM), we learn distributions over sequences of words, subwords or characters, and the last two can allow an open-vocabulary generation (Sutskever et al., 2011). We rely on subword segmentation as a widespread approach to generate rare subword units (Sennrich et al., 2016). However, the lack of a representative corpus, in terms of the word vocabulary, can constrain the unsupervised segmentation (e.g. with scarce monolingual texts (Joshi et al., 2020)). As an alternative, we could use character-level modelling, since it also has access to subword information (Kim et al., 2016), but we face long-term dependency issues and require longer training time to converge. Similar issues are extended to other generation tasks, such as machine translation (MT).

In this context, we focus on syllables, which are speech units: “A syl-la-ble con-tains a sin-gle vowel u-nit”. syllables can be defined as a group of segments that is pronounced as a single articulatory movement. Syllables are fundamental phonological units since they participate in important word prosodic patterns, such as stress assignment. In this sense, syllables are more linguistically relevant units than characters, and behave as a mapping function to reduce the length of the sequence with a larger “alphabet” or syllabary. Their extraction can be rule-based and corpus-independent, but data-driven methods or hyphenation using dictionaries can approximate them as well.

We assess whether syllables are useful for encoding and/or decoding a diverse set of languages on two generation tasks. First, for LM, we study 21 languages, to cover different levels of orthographic depth, which is the degree of grapheme-phoneme correspondence (Borgwaldt et al., 2005) and a factor that can increase complexity to syllabification (Marjou, 2021).¹ Whereas for MT, we focus on the distant and low-resource language-pair of Spanish–Shipibo–Konibo. We choose Shipibo–Konibo² because it is an endangered language with scarce textual corpora, which limits unsupervised segmentation methods, and has a transparent orthography, which could be beneficial to syllabification. Also, we consider multilingual MT systems, as they outperformed pairwise systems for the chosen language pair (Mager et al., 2021).

2 Related work

The closest LM study to ours is from Mikolov et al. (2012) for subword-grained prediction in English, where they used syllables as a proxy to split words with low frequency, reduce the vocabulary and compress the model size. Besides, syllable-aware LM

¹E.g., English has a deep orthography (weak correspondence), whereas Finnish is transparent (Ziegler et al., 2010).

²See Appendix A for more details about the language.

was addressed by Assylbekov et al. (2017) for English, German, French, Czech, Spanish and Russian, and by Yu et al. (2017) for Korean. However, in both cases, the syllables were composed with convolutional filters into word-level representations for closed-vocabulary generation. Besides, for subword-aware open-vocabulary LM, Blevins and Zettlemoyer (2019) incorporated morphological supervision with a multi-task objective.

For syllable-based MT, there are mostly studies for related paired languages, such as Indic languages (in statistical MT without subword-based baselines: Kunchukuttan and Bhattacharyya (2016)), Tibetan–Chinese (Lai et al., 2018), and Myanmar–Rakhine (Myint Oo et al., 2019). Instead, Spanish–Shipibo-Konibo is a non-related language-pair. The only distant pair was English–Myanmar (ShweSin et al., 2019), but they did not compare it with unsupervised subwords. Neither of these studies analysed multilingual settings.

3 Open-vocabulary language modelling with a comparable perplexity

Open-vocabulary output We generate the same input unit (e.g. characters, syllables or other subwords) as an open-vocabulary LM task, where there is no prediction of an “unknown” or out-of-vocabulary word-level token (Sutskever et al., 2011). We thereby differ from previous works, and refrain from composing the syllable representations into words to evaluate only word-level perplexity.

Character-level perplexity For a fair comparison across all granularities, we evaluate all results with character-level perplexity:

$$\text{ppl}^c = \exp(\mathcal{L}_{\text{LM}}(\mathbf{s}) \cdot \frac{|\mathbf{s}^{\text{seg}}| + 1}{|\mathbf{s}^c| + 1}) \quad (1)$$

where $\mathcal{L}_{\text{LM}}(\mathbf{s})$ is the cross entropy of a string \mathbf{s} computed by the neural LM, and $|\mathbf{s}^{\text{seg}}|$ and $|\mathbf{s}^c|$ refer to the length of \mathbf{s} in the chosen segmentation and character-level units, respectively (Mielke, 2019). The extra unit considers the end of the sequence.

3.1 Experimental setup

Languages and datasets Corpora are listed in Table 4 in Appendix B. We first choose WikiText-2-raw (en_w ; Merity et al., 2016), which contains around two million word-level tokens extracted from Wikipedia articles in English. Furthermore, we employ 20 Universal Dependencies (UD; Nivre et al., 2020) treebanks, similarly to Blevins

and Zettlemoyer (2019).³ Finally, we include the Shipibo-Konibo (shp) side of the parallel corpora provided by the AmericasNLP shared task on MT (Mager et al., 2021), which is also used in §4.

Syllable segmentation (SYL) For splitting syllables in different languages, we used rule-based syllabification tools for English, Spanish, Russian, Finnish, Turkish and Shipibo-Konibo, and dictionary-based hyphenation tools for all languages except the ones mentioned above. We list the tools in Appendix C.

Segmentation baselines Besides characters (CHAR) and the annotated morphemes in the UD treebanks (MORPH), we consider Polyglot (POLY)⁴, which includes models for unsupervised morpheme segmentation trained with Morfessor (Virpioja et al., 2013). Moreover, we employ an unsupervised subword segmentation baseline of Byte Pair Encoding (BPE; Sennrich et al., 2016)⁵ with different vocabulary sizes from 2,500 to 10,000 tokens, with 2,500 steps. We also fix the parameter to the syllabary size. Appendix C includes details about the segmentation format.

Model and training Following other open-vocabulary LM studies (Mielke and Eisner, 2019; Mielke et al., 2019), we use a low-compute version of an LSTM neural network, named Average SGD Weight-Dropped (Merity et al., 2018). See the hyperparameter details in Appendix E.

3.2 Results and discussion

Table 1 shows the ppl^c values for the different levels of segmentation we considered in the study, where we did not tune the neural LM model for a specific segmentation. We observe that syllables always result in better perplexities than other granularities, even for deep orthography languages such as English or French. The results obtained by the BPE baselines are relatively poor as well, and they could not beat characters in any dataset, even though we searched for an optimal vocabulary size for the BPE algorithm. The advantage of using syllables is that we do not need to tune a

³The languages are chosen given the availability of an open-source syllabification or hyphenation tool. We prefer to use the UD treebanks, instead of other well-known datasets for language modelling (e.g. Multilingual Wikipedia Corpus (Kawakami et al., 2017)), because they provide morphological annotations, which are fundamental for this study.

⁴polyglot-nlp.com

⁵We use: <https://github.com/huggingface/tokenizers>

	CHAR	MORPH	POLY	SYL	BPE _{best}
en _w *	2.48 ±0.0	-	2.8 ±0.0	1.96 ±0.0	2.91 ±0.0
bg	3.56 ±0.03	4.09 ±0.05	4.69 ±0.01	2.87 ±0.0	5.19 ±0.01
ca	2.84 ±0.0	3.11 ±0.02	3.26 ±0.01	2.21 ±0.0	3.31 ±0.0
cs	3.32 ±0.0	3.11 ±0.01	4.18 ±0.01	2.66 ±0.0	4.24 ±0.0
da	4.25 ±0.01	4.42 ±0.04	5.6 ±0.0	3.1 ±0.01	6.21 ±0.03
de	3.5 ±0.04	3.36 ±0.08	3.79 ±0.0	2.48 ±0.0	3.86 ±0.02
en*	4.11 ±0.01	4.39 ±0.08	5.67 ±0.01	2.82 ±0.07	5.65 ±0.04
es*	3.16 ±0.01	3.71 ±0.04	3.95 ±0.01	2.51 ±0.0	3.98 ±0.0
fi*	3.77 ±0.01	4.05 ±0.12	4.76 ±0.01	3.1 ±0.0	5.27 ±0.01
fr	3.09 ±0.01	3.67 ±0.02	3.82 ±0.01	2.3 ±0.01	3.87 ±0.01
hr	3.52 ±0.02	3.92 ±0.01	4.34 ±0.0	2.8 ±0.0	4.52 ±0.02
it	2.8 ±0.0	3.19 ±0.0	3.43 ±0.01	2.27 ±0.01	3.61 ±0.0
lv	4.55 ±0.02	5.31 ±0.0	6.82 ±0.02	3.59 ±0.0	7.19 ±0.0
nl	3.83 ±0.05	3.69 ±0.1	4.44 ±0.01	2.76 ±0.01	4.83 ±0.01
pl	4.03 ±0.01	4.77 ±0.22	5.96 ±0.04	3.19 ±0.0	5.99 ±0.0
pt	3.31 ±0.01	3.46 ±0.03	4.03 ±0.01	2.56 ±0.0	4.24 ±0.01
ro	3.4 ±0.02	3.89 ±0.04	4.25 ±0.01	2.72 ±0.0	4.71 ±0.01
ru*	3.28 ±0.0	2.93 ±0.01	4.05 ±0.0	2.69 ±0.01	4.04 ±0.0
sk	6.16 ±0.05	5.1 ±0.07	7.61 ±0.08	4.62 ±0.01	10.51 ±0.03
tr*	4.16 ±0.05	4.86 ±0.05	6.41 ±0.07	3.66 ±0.03	6.98 ±0.1
uk	4.92 ±0.02	6.45 ±0.11	8.11 ±0.03	4.24 ±0.02	9.23 ±0.02
shp*	4.48±0.01	-	-	2.15 ±0.02	3.50±0.03

Table 1: Character-level perplexity (\downarrow) in test. We show the mean and standard deviation for three runs with different seeds. BPE shows the best result from models with different vocabulary sizes. SYL presents the syllabification-based result if it is available (*), or the hyphenation otherwise.

hyper-parameter to extract a different set of subword pieces.

As a significant outcome, we note that syllables did not fail to beat characters, at least in an open-vocabulary LM task, which extends the findings of Assylbekov et al. (2017). One potential reason is that character-level modelling requires a larger model capacity due to the longer sequences, however, that is also an advantage towards syllables. Besides, other subword pieces with a closer sequence length to syllables (BPE, MORPH or POLY) were still outperformed.

Finally, in Appendix D, we further discuss the relationship between the syllable type/token ratio with the word vocabulary growth and perplexity.

4 Low-resource Machine Translation

After observing the positive impact on LM, we focus on syllables for MT, which adds complexity to the process, as there is at least one extra language involved. In contrast to prior work, we (i) study syllable-based MT for a distant and low-resource language-pair, Spanish–Shipibo-Konibo; (ii) compare syllables against the most widespread unsupervised segmentation method (BPE) with automatic metrics and human evaluation; and (iii) anal-

yse the applicability of syllables on multilingual translation systems. The last element is significant, as a multilingual setting is the state-of-the-art approach for leveraging low-resource language-pairs performance (Siddhant et al., 2022). Moreover, we decided to apply syllabification only on Shipibo-Konibo, a highly synthetic⁶ language with scarce textual data and with a transparent orthography⁷.

For this reason, we focus in three settings. First, MONO, a pairwise system where each source and target is segmented with a different method. Second, JOINT, another pairwise system where the BPE baseline is jointly trained with the source and target data. Third, O2M, a multilingual one-to-many⁸ system, where the BPE baseline is jointly trained with all the languages (we added Spanish–English in our experiments).

4.1 Experimental setup

Data For Spanish–Shipibo-Konibo (es–shp), we use the dataset provided by the AmericasNLP workshop (Mager et al. (2021); Galarreta et al. (2017); Gómez Montoya et al. (2019)), and perform the same split as Mager et al. (2022) for the dev and test subsets, to make the results comparable to their morphological segmentation experiments. For the multilingual case, we use the Spanish–English (es–en) train set from EuroParl (Koehn, 2005) and newscomentary-v8, and the NEWSTEST2013.ES-EN (Bojar et al., 2013) evaluation sets.

Segmentation (i) BPE (Sennrich et al., 2016) is our baseline segmentation method, and we use the implementation of SentencePiece (Kudo and Richardson, 2018). Similar to Mager et al. (2022), we fix the best vocabulary size at 5000 pieces for the MONO setting, after trying different values from 1k to 10k. JOINT and O2M use 5000 and 16000 pieces, respectively.

(ii) Syllabification (SYL) for Shipibo-Konibo is adapted from Alva and Oncevay (2017). The original method uses syllables to verify whether a word is composed by consistent syllables for spell-checking. In our experiments, when a word can not be syllabified, we split it into characters for the

⁶With a high ratio of number of morphemes per word.

⁷We attempted to use syllables on Spanish and English as well, but with negative results. With large data, unsupervised segmentation methods like BPE can obtain more significant and overlapping subwords from source and target.

⁸We do not consider the many-to-one direction due to resource constraints, and because we observed that the improvements by syllables are noted when decoding Shipibo-Konibo.

MONO setting, and we use the joint-BPE segmentation model for the JOINT and O2M settings.

Model and training We reproduce Mager et al. (2022)’s settings, by using the fairseq toolkit (Ott et al., 2019), and a Transformer model (Vaswani et al., 2017) with smaller dimensions (Guzmán et al., 2019). For the multilingual O2M setting, we use a sampling approach with 5 of temperature (Aharoni et al., 2019). See details in Appendix E.

Evaluation We use chrF (Popović, 2015) from SACREBLEU (Post, 2018)⁹ and also perform a human evaluation of 100 samples per system (BPE and SYL), following the annotation protocol used in the AmericasNLP shared task (Mager et al., 2021).

4.2 Results and discussion

Table 2 shows the translation performance in all settings, and we observe that syllables are statistically better than the BPE baseline when translating from Spanish into Shipibo-Konibo, but not in the other direction. This fact indicates that syllables support the decoding more than the encoding step of a language with a transparent orthography. Also, the JOINT setting reduces the gap between BPE and SYL, probably due to the shared roots between the two languages (i.e., loanwords from Spanish into Shipibo-Konibo). Furthermore, we note that the impact of syllables is not minimised in a multilingual system (O2M), where the performance for es→shp has drastically improved, and the other language-pair (es→en) retains a comparable result.

Moreover, our MONO experiments are comparable with the study of Mager et al. (2022), where they tested several unsupervised and supervised morphological segmentation methods against BPE for MT in four polysynthetic languages (including Shipibo-Konibo). Our result with syllables in es→shp outperforms all other approaches, such as LMVR (Ataman et al., 2017), with a 38.99 chrF score. This indicates that syllables are a robust alternative to morphologically-aware methods when we are dealing with limited data and translating into a polysynthetic language.

4.3 Human evaluation

We also conducted a small human evaluation of system outputs using a 5-points scale for the adequacy and fluency of the Spanish→Shipibo-Konibo translation, which is the translation direction that

⁹chrF2+numchars.6+space.false+v.1.5.0.

	BPE	SYL	BPE SYL	
	es→shp		es→en	
MONO	37.62±1.87	41.27* ±0.54	53.99	53.85
JOINT	40.41±0.82	41.74* ±0.95		
O2M	48.30	51.25*		
	shp→es			
MONO	33.37 ±0.79	32.85*±1.22		
JOINT	34.55 ±0.56	33.13*±0.75		

Table 2: chrF scores in the test subsets. MONO: single BPE model (5k pieces) for each source and target. JOINT: joint BPE model (5k) for both source and target. O2M: joint BPE model (16k) for ES, EN and SHP. For the first two settings, we run three experiments and present the mean and standard deviation. The latter only has one run due to resource constraints, and we report es→en scores as a reference. Syllabification (SYL) is only applied on the Shipibo-Konibo side, and (*) indicates a p-value ≤ 0.05 against the BPE baseline.

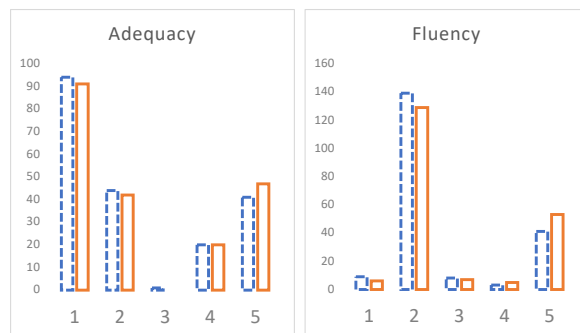


Figure 1: Adequacy and fluency scores (1-5) for 200 outputs of two approaches: BPE (dashed blue) and SYL (solid orange), from the best es→shp O2M system.

benefited from the syllable segmentation. The annotation protocol and annotator’s information is provided in Appendix F.

Figure 1 shows the scores annotated for adequacy and fluency, where we compare BPE and SYL in the O2M setting, which obtained the best performance for both segmentation methods. We observe that the adequacy is very poor for both systems (1-2), but there is an advantage for SYL in the smaller batch of highest adequacy (5), with 3% more of the total samples. Regarding fluency, both systems mostly obtain a low score (2), but there is a consistent advantage for SYL over BPE in the highest value (5), with 6.5% more of the total samples. The differences are very small to determine whether a segmentation works better than the other from human judgement, but they are consistent with the automatic evaluation provided previously. A larger sample, an extra annotator, or more robust

systems could aid to clarify other potential benefits.

5 Limitations and opportunities

Syllables only cannot offer a universal solution to the subword segmentation problem for all languages, as the syllabification tools are language-dependent. Besides, the analysis should be extended to more scripts and morphological types. Furthermore, we do not encode any semantics in the syllable-vector space, with a few exceptions like in Korean (Choi et al., 2017). Nevertheless, our results confirm that syllables are reliable for LM and MT, and building a syllable splitter might require less effort than annotating morphemes to train a robust supervised tool¹⁰.

Specifically for MT, syllables could be useful when: (i) we are dealing with extremely low-resource data, which affects unsupervised word segmentation, (ii) we are translating into a language with a high synthesis, which has been observed as a factor that impacts on MT performance (Oncevay et al., 2022), and (iii) we are working with a language with a transparent orthography. This is the scenario for several languages from the Americas, where their writing systems have been recently standardised for documentation and revitalisation purposes (Mager et al., 2018), and some resources for MT have been compiled (Mager et al., 2021).

6 Conclusion

We have proved that syllables are valuable for generation tasks such as: (i) Open-vocabulary LM, where they behave positively even for languages with deep orthography, and overcome character and subword baselines. (ii) Low-resource and multilingual MT, outperforming BPE pieces when we translate into a language with a transparent orthography and complex morphology (high synthesis), even when the language-pair is not related.

Acknowledgements

The first author acknowledges the support of NVIDIA Corporation with the donation of a Titan Xp GPU used for the study. The last author acknowledges the Max Planck Institute for Evolutionary Anthropology, Department of Linguistic and

¹⁰For instance, the syllabification tool that we used for English is based on five general rules from: <https://www.howmanysyllables.com/divideintosyllables>. Their implementation should take less effort than annotating a UD treebank or building a Finite-State-Transducer for morphological analysis.

Cultural Evolution, for its support to the development of the Chana Field Station in the Amazonian region of Peru, and the support of CONCYTEC-ProCiencia, Peru, under the contract 183-2018-FONDECYT-BM-IADT-MU from the funding call E041-2018-01-BM.

References

- Roe Aharoni, Melvin Johnson, and Orhan Firat. 2019. [Massively multilingual neural machine translation](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3874–3884, Minneapolis, Minnesota. Association for Computational Linguistics.
- Carlo Alva and Arturo Oncevay. 2017. [Spell-checking based on syllabification and character-level graphs for a Peruvian agglutinative language](#). In *Proceedings of the First Workshop on Subword and Character Level Models in NLP*, pages 109–116, Copenhagen, Denmark. Association for Computational Linguistics.
- Zhenisbek Assylbekov, Rustem Takhanov, Bagdat Myrzakhmetov, and Jonathan N. Washington. 2017. [Syllable-aware neural language models: A failure to beat character-aware ones](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1866–1872, Copenhagen, Denmark. Association for Computational Linguistics.
- Duygu Ataman, Matteo Negri, Marco Turchi, and Marcello Federico. 2017. Linguistically motivated vocabulary reduction for neural machine translation from turkish to english. *arXiv preprint arXiv:1707.09879*.
- Terra Blevins and Luke Zettlemoyer. 2019. [Better character language modeling through morphology](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1606–1613, Florence, Italy. Association for Computational Linguistics.
- Ondřej Bojar, Christian Buck, Chris Callison-Burch, Christian Federmann, Barry Haddow, Philipp Koehn, Christof Monz, Matt Post, Radu Soricut, and Lucia Specia. 2013. [Findings of the 2013 Workshop on Statistical Machine Translation](#). In *Proceedings of the Eighth Workshop on Statistical Machine Translation*, pages 1–44, Sofia, Bulgaria. Association for Computational Linguistics.
- Susanne R Borgwaldt, Frauke M Hellwig, and Annette MB De Groot. 2005. Onset entropy matters—letter-to-phoneme mappings in seven languages. *Reading and Writing*, 18(3):211–229.
- Elisabeth Borleffs, Ben AM Maassen, Heikki Lyytinen, and Frans Zwarts. 2017. Measuring orthographic

- transparency and morphological-syllabic complexity in alphabetic orthographies: a narrative review. *Reading and writing*, 30(8):1617–1638.
- Sanghyuk Choi, Taeuk Kim, Jinseok Seol, and Sang-goo Lee. 2017. [A syllable-based technique for word embeddings of Korean words](#). In *Proceedings of the First Workshop on Subword and Character Level Models in NLP*, pages 36–40, Copenhagen, Denmark. Association for Computational Linguistics.
- Ana-Paula Galarreta, Andrés Melgar, and Arturo Oncevay. 2017. [Corpus creation and initial SMT experiments between Spanish and Shipibo-konibo](#). In *Proceedings of the International Conference Recent Advances in Natural Language Processing, RANLP 2017*, pages 238–244, Varna, Bulgaria. INCOMA Ltd.
- Héctor Erasmo Gómez Montoya, Kervy Dante Rivas Rojas, and Arturo Oncevay. 2019. [A continuous improvement framework of machine translation for Shipibo-konibo](#). In *Proceedings of the 2nd Workshop on Technologies for MT of Low Resource Languages*, pages 17–23, Dublin, Ireland. European Association for Machine Translation.
- Francisco Guzmán, Peng-Jen Chen, Myle Ott, Juan Pino, Guillaume Lample, Philipp Koehn, Vishrav Chaudhary, and Marc’Aurelio Ranzato. 2019. [The FLORES evaluation datasets for low-resource machine translation: Nepali–English and Sinhala–English](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6098–6111, Hong Kong, China. Association for Computational Linguistics.
- Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. [The state and fate of linguistic diversity and inclusion in the NLP world](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6282–6293, Online. Association for Computational Linguistics.
- Kazuya Kawakami, Chris Dyer, and Phil Blunsom. 2017. [Learning to create and reuse words in open-vocabulary neural language modeling](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1492–1502, Vancouver, Canada. Association for Computational Linguistics.
- Yoon Kim, Yacine Jernite, David Sontag, and Alexander M. Rush. 2016. [Character-aware neural language models](#). In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, AAAI’16*, pages 2741–2749. AAAI Press.
- Diederik P. Kingma and Jimmy Ba. 2015. [Adam: A method for stochastic optimization](#). In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Philipp Koehn. 2005. [Europarl: A parallel corpus for statistical machine translation](#). In *Proceedings of Machine Translation Summit X: Papers*, pages 79–86, Phuket, Thailand.
- Taku Kudo and John Richardson. 2018. [SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.
- Anoop Kunchukuttan and Pushpak Bhattacharyya. 2016. [Orthographic syllable as basic unit for SMT between related languages](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1912–1917, Austin, Texas. Association for Computational Linguistics.
- Wen Lai, Xiaobing Zhao, and Wei Bao. 2018. [Tibetan-Chinese neural machine translation based on syllable segmentation](#). In *Proceedings of the AMTA 2018 Workshop on Technologies for MT of Low Resource Languages (LoResMT 2018)*, pages 21–29, Boston, MA. Association for Machine Translation in the Americas.
- Manuel Mager, Ximena Gutierrez-Vasques, Gerardo Sierra, and Ivan Meza-Ruiz. 2018. [Challenges of language technologies for the indigenous languages of the Americas](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 55–69, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Manuel Mager, Arturo Oncevay, Abteen Ebrahimi, John Ortega, Annette Rios, Angela Fan, Ximena Gutierrez-Vasques, Luis Chiruzzo, Gustavo Giménez-Lugo, Ricardo Ramos, Ivan Vladimir Meza Ruiz, Rolando Coto-Solano, Alexis Palmer, Elisabeth Mager-Hois, Vishrav Chaudhary, Graham Neubig, Ngoc Thang Vu, and Katharina Kann. 2021. [Findings of the AmericasNLP 2021 shared task on open machine translation for indigenous languages of the Americas](#). In *Proceedings of the First Workshop on Natural Language Processing for Indigenous Languages of the Americas*, pages 202–217, Online. Association for Computational Linguistics.
- Manuel Mager, Arturo Oncevay, Elisabeth Mager, Katharina Kann, and Thang Vu. 2022. [BPE vs. morphological segmentation: A case study on machine translation of four polysynthetic languages](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 961–971, Dublin, Ireland. Association for Computational Linguistics.
- Xavier Marjou. 2021. [OTEANN: Estimating the transparency of orthographies with an artificial neural network](#). In *Proceedings of the Third Workshop on Computational Typology and Multilingual NLP*, pages 1–9, Online. Association for Computational Linguistics.

- Stephen Merity, Nitish Shirish Keskar, and Richard Socher. 2018. [Regularizing and optimizing LSTM language models](#). In *International Conference on Learning Representations*.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. [Pointer sentinel mixture models](#). *arXiv preprint arXiv:1609.07843*.
- Sabrina J. Mielke. 2019. Can you compare perplexity across different segmentations? Available in: <http://sjmielke.com/comparing-perplexities.htm>.
- Sabrina J. Mielke, Ryan Cotterell, Kyle Gorman, Brian Roark, and Jason Eisner. 2019. [What kind of language is hard to language-model?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4975–4989, Florence, Italy. Association for Computational Linguistics.
- Sabrina J. Mielke and Jason Eisner. 2019. [Spell once, summon anywhere: A two-level open-vocabulary language model](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6843–6850.
- Tomáš Mikolov, Ilya Sutskever, Anoop Deoras, Hainan Le, Stefan Kombrink, and Jan Cernocky. 2012. Subword language modeling with neural networks. *preprint (http://www.fit.vutbr.cz/imikolov/rnnlm/char.pdf)*, 8:67.
- Thazin Myint Oo, Ye Kyaw Thu, and Khin Mar Soe. 2019. [Neural machine translation between Myanmar \(Burmese\) and Rakhine \(Arakanese\)](#). In *Proceedings of the Sixth Workshop on NLP for Similar Languages, Varieties and Dialects*, pages 80–88, Ann Arbor, Michigan. Association for Computational Linguistics.
- Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Jan Hajič, Christopher D. Manning, Sampo Pyysalo, Sebastian Schuster, Francis Tyers, and Daniel Zeman. 2020. [Universal Dependencies v2: An evergrowing multilingual treebank collection](#). In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 4034–4043, Marseille, France. European Language Resources Association.
- Arturo Oincevay. 2021. [Peru is multilingual, its machine translation should be too?](#) In *Proceedings of the First Workshop on Natural Language Processing for Indigenous Languages of the Americas*, pages 194–201, Online. Association for Computational Linguistics.
- Arturo Oincevay, Duygu Ataman, Niels Van Berkel, Barry Haddow, Alexandra Birch, and Johannes Bjerva. 2022. [Quantifying synthesis and fusion and their impact on machine translation](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1308–1321, Seattle, United States. Association for Computational Linguistics.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. [fairseq: A fast, extensible toolkit for sequence modeling](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota. Association for Computational Linguistics.
- Maja Popović. 2015. [chrF: character n-gram F-score for automatic MT evaluation](#). In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.
- Matt Post. 2018. [A call for clarity in reporting BLEU scores](#). In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. [Neural machine translation of rare words with subword units](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- Yimong ShweSin, Win Pa Pa, and KhinMar Soe. 2019. [UCSYNLP-lab machine translation systems for WAT 2019](#). In *Proceedings of the 6th Workshop on Asian Translation*, pages 195–199, Hong Kong, China. Association for Computational Linguistics.
- Aditya Siddhant, Ankur Bapna, Orhan Firat, Yuan Cao, Mia Xu Chen, Isaac Caswell, and Xavier Garcia. 2022. [Towards the next 1000 languages in multilingual machine translation: Exploring the synergy between supervised and self-supervised learning](#).
- Ilya Sutskever, James Martens, and Geoffrey Hinton. 2011. [Generating text with recurrent neural networks](#). In *Proceedings of the 28th International Conference on Machine Learning*, ICML’11, pages 1017–1024, USA. Omnipress.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc.
- Sami Virpioja, Peter Smit, Stig-Arne Grönroos, and Mikko Kurimo. 2013. Morfessor 2.0: Python implementation and extensions for morfessor baseline. In *Aalto University publication series*. Department of Signal Processing and Acoustics, Aalto University.
- Seunghak Yu, Nilesh Kulkarni, Haejun Lee, and Jihie Kim. 2017. [Syllable-level neural language model for agglutinative language](#). In *Proceedings of the*

First Workshop on Subword and Character Level Models in NLP, pages 92–96, Copenhagen, Denmark. Association for Computational Linguistics.

Roberto Zariquiey, Harald Hammarström, Mónica Arakaki, Arturo Oncevay, John Miller, Aracelli García, and Adriano Ingunza. 2019. *Obsolescencia lingüística, descripción gramatical y documentación de lenguas en el Perú: hacia un estado de la cuestión*. *Lexis*, 43(2):271–337.

Johannes C Ziegler, Daisy Bertrand, Dénes Tóth, Valéria Csépe, Alexandra Reis, Luís Fáisca, Nina Saine, Heikki Lyytinen, Anniek Vaessen, and Leo Blomert. 2010. Orthographic depth and its impact on universal predictors of reading: A cross-language investigation. *Psychological science*, 21(4):551–559.

A The Shipibo-Konibo language

Shipibo-Konibo (shp) is the largest and most vital language within the Pano language family. With more than 30,000 speakers, the Shipibo-Konibo are among the largest indigenous groups in Peru. Shipibo-Konibo people mainly live in the Peruvian Amazonia (in the regions of Ucayali, Loreto, Huánuco and Madre de Dios), but there are also large groups of Shipibo-Konibo people living in the Peruvian coast (particularly in Lima and Ica).

Shipibo-Konibo is a well-documented language, although the publicly available data on this language is rather small (Zariquiey et al., 2019). It has a complex morphology due to its high synthesis (high ratio of morphemes per word, mostly by suffixation) and its agglutinative nature. Its orthography can be considered transparent, because its alphabet was recently standardised by the Peruvian Government (Alva and Oncevay, 2017), and the datasets we are using in all experiments are provided with the most recent writing standard (Mager et al., 2021). Machine translation research on Shipibo-Konibo has focused on the development of new parallel corpora (Galarreta et al., 2017; Gómez Montoya et al., 2019), the application of multilingual models (Oncevay, 2021), or the impact of morphological segmentation methods (Mager et al., 2022). However, neither of them has focused on syllables as a unit for segmentation. For this study, we adapt the syllabification function proposed by Alva and Oncevay (2017), which was used for spell-checking.

B Dataset details

Table 4 shows the size of the training, validation and test splits for all the datasets used in the LM task, while Table 3 shows details of the Spanish–Shipibo-Konibo and Spanish–English parallel corpora used in the MT task.

	train	dev	test
es–shp	13,102	587	1,030
es–en	2,140,175	5,003	3,000

Table 3: Total number of sentences in train, dev and test splits for the language-pairs used in the MT experiments.

C Segmentation details

Tools We list the tools for rule-based syllabification and dictionary-based hyphenation:

	Train			Valid			Test		
	Word	Syl	Char	Word	Syl	Char	Word	Syl	Char
en _w	2,089	4,894	10,902	218	505	1,157	246	568	1,304
bg	125	386	710	16	50	92	16	49	90
ca	436	1,123	2,341	59	152	317	61	157	327
cs	1,158	3,546	6,868	157	482	933	172	524	1,012
da	81	215	442	10	28	57	10	27	56
de	260	735	1,637	12	34	75	16	45	102
en	210	488	1,061	26	61	133	26	61	132
es	376	1,060	2,043	37	103	198	12	33	64
fi	165	595	1,224	19	67	137	21	76	155
fr	360	837	1,959	36	84	197	10	23	54
hr	154	484	930	20	62	119	23	75	145
it	263	762	1,504	11	32	64	10	28	57
lv	113	349	690	19	58	115	20	59	116
nl	187	488	1,074	12	30	66	11	31	68
pl	102	293	589	13	37	73	13	37	74
pt	192	551	1,040	10	29	54	9	27	51
ro	183	549	1,056	17	51	98	16	48	94
ru	867	2,707	5,411	118	364	722	117	360	717
sk	80	232	437	12	39	76	13	41	80
tk	38	126	242	10	33	63	10	33	64
uk	88	289	501	12	41	71	16	56	99
shp	43	141	398						

Table 4: Total number of tokens (in thousands) at word, syllable and character-level for all the splits.

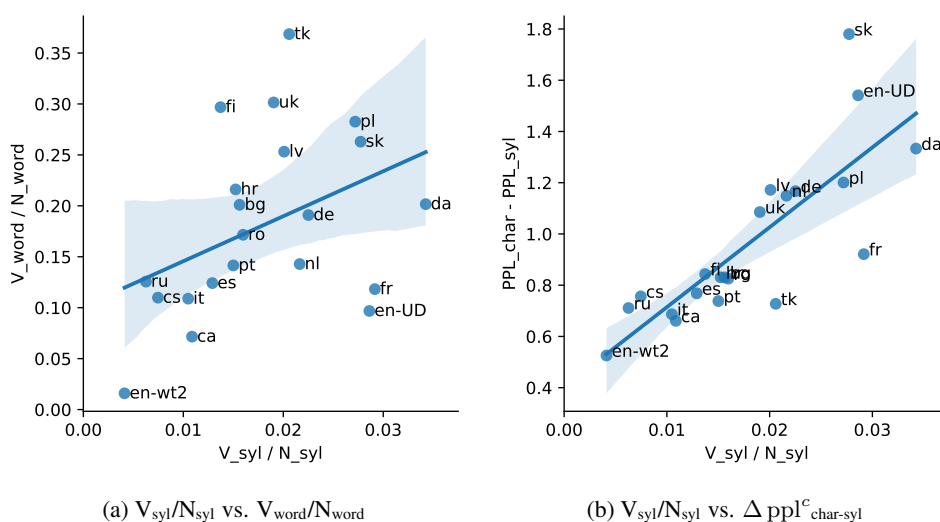


Figure 2: Left (a): Vocabulary growth rate of syllables (x-axis) versus words (y-axis). Right (b): Vocabulary growth rate of syllables (x-axis) versus the difference of ppl^c obtained by characters and syllables (y-axis).

- English syllabification: Extracted from <https://www.howmanysyllables.com/>
- Spanish syllabification: <https://pypi.org/project/pylabeador/>
- Russian syllabification: <https://github.com/Koziev/rusyllab>
- Finnish syllabification: <https://github.com/tsnaomi/finnsyll>
- Turkish syllabification: <https://github.com/MeteHanC/turkishnlp>
- Shipibo-Konibo syllabification: Alva and Onceva (2017)
- Hyphenation: PyPhen (<https://pyphen.org/>), which is based on Hunspell dictionaries.

Format For syllables in the LM task, we separate the subwords as: “A @ syl la ble @ con tains @ a @ sin gle @ vow el @ u nit”, where “@” is a special token that indicates the word boundary. We also evaluated syllables with a segmentation format like in Sennrich et al. (2016): “A syl@ la@ ble con@ tains a ...”, but we obtained lower performance in general. Whereas in the MT task, we adopt the segmentation format used by SentencePiece (Kudo and Richardson, 2018) for syllables: “_A _syl la

ble _contains _a _single _vowel _unit”.

D Type/token ratio of syllables in LM

In Figure 2a, we show a scatter plot of the token/type growth rate of syllables versus words for all languages and corpora. In other words, the ratio of syllable-types (syllabary or V_{syl}) per total number of syllable-tokens (N_{syl}) versus the type/token ratio of words ($V_{\text{word}}/N_{\text{word}}$) in the train set. The figure suggests at least a weak relationship, which agrees with the notion that a low word-vocabulary richness only requires a low syllabary richness for expressivity. Also, a richer vocabulary can use a richer syllabary or just longer words, so the distribution of the vocabulary richness could be larger.

We expected that the syllabary growth rate ($V_{\text{syl}}/N_{\text{syl}}$) for a low phonemic language like English would be relatively high, but wikitext-2 (en-wt2) is located in the bottom-left corner of the plot, probably caused by its large amount of word-tokens. However, we observed a large $V_{\text{syl}}/N_{\text{syl}}$ for the English (en-UD) and French (fr) treebanks, despite their low $V_{\text{word}}/N_{\text{word}}$ ratio, which is an expected pattern for deep orthographies.

We also observe that languages with a more transparent orthography, like Czech (cs) or Finnish (fi), are located in the left side of the figure, whereas Turkish (tr) is around the middle section. Nevertheless, our study does not aim to analyse the relationship between the level of phonemic orthography with the $V_{\text{syl}}/N_{\text{syl}}$ ratio. For that purpose, we might need an instrument to measure how deep or shallow a language orthography is (Marjou, 2021; Borgwaldt et al., 2005; Borleffs et al., 2017), and a multi-parallel corpus for a more fair comparison.

Finally, in Figure 2b we observe a stronger relationship of the syllable type/token ratio with the difference of CHAR’s ppl^c minus SYL’s ppl^c . In other words, if our dataset possesses a rich syllabary, we are fairly approximating the amount of word-level tokens, which reduces the ppl^c gain.

E Model and Training

LM In contrast with the default settings, we use a smaller embedding size of 500 units for faster training. Additionally, we have 3 layers of depth, 1152 of hidden layer size and a dropout of 0.15. We train for 25 epochs with a batch size of 64, a learning rate of 0.002 and Adam optimiser (Kingma and Ba, 2015) with default parameters. We fit the model using the one cycle policy and an early stopping of

4. We run our experiments in a NVIDIA Titan Xp.

MT Similar to Mager et al. (2022), we use a small Transformer model for our low-resource MT settings, following Guzmán et al. (2019): “with 5 encoder and 5 decoder layers, where the number of attention heads, embedding dimension and inner-layer dimension are 2, 512 and 2048, respectively”. For the pairwise systems, we train up to 100 epochs with an early stopping policy of 5 (validating every 5 epochs), whereas for the multilingual systems we train up to 30 epochs. For all the experiments, we use 4 NVIDIA GeForce GTX 1080 Ti GPUs.

F Human evaluation

F.1 Annotation protocol

Adapted and summarised from the AmericasNLP shared task (Mager et al., 2021): The expert received the source sentence in Spanish, the reference in Shipibo-Konibo, and an anonymized system output, which includes the baseline (BPE) and our syllable-based system (SYL). The expert received only 200 samples (per system, same entries) that were randomly selected and shuffled. They were asked to annotate **Adequacy** (Does the output sentence express the meaning of the reference?) from 1 to 5 (extremely bad, bad, neutral, sufficiently good, excellent), and **Fluency** (Is the output sentence easily readable and looks like a human-produced text?) from 1 to 5 as well.

F.2 About the annotator

The annotator is a native speaker of Shipibo-Konibo, a certified and professional translator, and a bilingual teacher in Peru. The annotator has experience in translating corpus for MT research, and performing human evaluation for Spanish–Shipibo-Konibo. This expertise is almost unique for Shipibo-Konibo, and we could not identify a second annotator with the same expertise to obtain inter-annotation agreement.