

REPRESENTATION AND BIAS IN MULTILINGUAL NLP: INSIGHTS FROM CONTROLLED EXPERIMENTS ON CONDITIONAL LANGUAGE MODELING

Anonymous authors

Paper under double-blind review

ABSTRACT

Inspired by the phenomenon of performance disparity between languages in machine translation, we investigate whether and to what extent languages are equally hard to “conditional-language-model”. Our goal is to improve our understanding and expectation of the relationship between language, data representation, size, and performance. We study one-to-one, bilingual conditional language modeling through a series of systematically controlled experiments with the Transformer and the 6 languages from the United Nations Parallel Corpus. We examine character, byte, and word models in 30 language directions and 5 data sizes, and observe indications suggesting a script bias on the character level, a length bias on the byte level, and a word bias that gives rise to a hierarchy in performance across languages. We also identify two types of sample-wise non-monotonicity — while word-based representations are prone to exhibit Double Descent, length can induce unstable performance across the size range studied in a novel meta phenomenon which we term *erraticity*. By eliminating statistically significant performance disparity on the character and byte levels by normalizing length and vocabulary in the data, we show that, in the context of computing with the Transformer, there is no complexity intrinsic to languages other than that related to their statistical attributes and that performance disparity is not a necessary condition but a byproduct of word segmentation. Our application of statistical comparisons as a fairness measure also serves as a novel rigorous method for the intrinsic evaluation of languages, resolving a decades-long debate on language complexity. While all these quantitative biases leading to disparity are mitigable through a shallower network, we find room for a human bias to be reflected upon. We hope our work helps open up new directions in the area of language and computing that would be fairer and more flexible and foster a new transdisciplinary perspective for DL-inspired scientific progress.

1 INTRODUCTION

With a transdisciplinary approach to explore a space at the intersection of Deep Learning (DL) / Neural Networks (NNs), language sciences, and language engineering, we report our undertaking in **use-inspired basic research** — with an application-related phenomenon as inspiration, we seek **fundamental scientific understanding** through empirical experimentation. This is *not* an application or machine translation (MT) paper, but one that strives to evaluate and seek new insights on language in the context of DL with a consideration to contribute to our evaluation, segmentation, and model interpretation practice in multilingual Natural Language Processing (NLP).

Our inspiration: performance disparity in MT The use case that inspired our investigation is the disparity of MT results reported in Junczys-Dowmunt et al. (2016). Of the 6 official languages of the United Nations (UN) — Arabic (AR), English (EN), Spanish (ES), French (FR), Russian (RU), and Chinese (ZH), results with target languages AR, RU, and ZH seem to be worse than those with EN/ES/FR, regardless of the algorithm, may it be from phrased-based Statistical MT (SMT/Moses

(Koehn et al., 2007)) or Neural MT (NMT).¹ The languages have the same amount of line-aligned, high-quality parallel data available for training, evaluation, and testing. This prompts the question: are some languages indeed harder to translate from or to?

Problem statement: are all languages equally hard to Conditional-Language-Model (CLM)?

A similar question concerning (monolingual) language modeling (LMing) was posed in Cotterell et al. (2018) and Mielke et al. (2019) along with the introduction of a method to evaluate LMs with multiway parallel corpora (multitexts) in information-theoretic terms. To explicitly focus on modeling the complexities that may or may not be *intrinsic* to the languages, we study the more fundamental process of CLMing without performing any translation. This allows us to eliminate confounds associated with generation and other evaluation metrics. One could think of our effort as estimating conditional probabilities with the Transformer, with a bilingual setup where perplexity of one target language (l_{trg}) is estimated given the parallel data in one source language (l_{src}), where $l_{\text{src}} \neq l_{\text{trg}}$. We focus on the very basics and examine the first step in our pipeline — input representation, holding everything else constant. Instead of measuring absolute cross-entropy scores, we evaluate the relative differences between languages from across 5 magnitudes of data sizes in 3 different representation types/levels. **We consider *bias* to be present when performance disparity in our Transformer models is statistically significant.**

1.1 SUMMARY OF FINDINGS AND CONTRIBUTIONS

In investigating performance disparity as a function of size and data with respect to language and representation on the Transformer in the context of CLMing, we find:

1. in a bilingual (one-to-one) CLMing setup, there is **neutralization of source language instances**, i.e. there are no statistically significant differences between source language pairs. Only pairs of target languages differ significantly (see Table 1).
2. We identify 2 types of **sample-wise non-monotonicity** on each of the primary representation levels we studied:
 - (a) **Double Descent** (Belkin et al., 2019; Nakkiran et al., 2020): on the word level, for all languages, performance at 10^2 lines is typically better than at 10^3 before it improves again at 10^4 and beyond. This phenomenon can also be observed in character models with ZH as a target language as well as on the word level with non-neural n-gram LMs;
 - (b) **erraticity**: performance is irregular and exhibits great variance across runs. We find sequence length to be predictive of this phenomenon. We show that this can be rectified by data transformation or hyperparameter tuning. In our study, erraticity affects AR and RU on the byte level where the sequences are too long with UTF-8 encoding and ZH when decomposed into strokes on the character level.
3. In eliminating performance disparity through lossless data transformation on the character and byte levels, we resolve language complexity (§ 4 and App. J). We show that, in the context of computing with the Transformer, unless word-based methods are used, there is no linguistic/morphological complexity applicable or necessary. There is **no complexity that is intrinsic to a language aside from its statistical properties**. Hardness in modeling is relative to and bounded by its representation level (**representation relativity**). On the character and byte levels, hardness is correlated with statistical properties concerning sequence length and vocabulary of a language, irrespective of its linguistic typological, phylogenetic, historical, or geographical profile, and can be eliminated. On the word level, hardness is correlated with vocabulary, and a complexity hierarchy arises through the manual preprocessing step of word tokenization. This complexity/disparity effected by word segmentation cannot be eliminated due to the fundamental qualitative differences in the definition of a “word” being one that neither holds universally nor is suitable/consistent for fair crosslinguistic comparisons. We find clarification of this expectation of disparity necessary because more diligent error analyses need to be afforded instead of simply accepting massively disparate results or inappropriately attributing under-performance to linguistic reasons.
4. Representational units of **finer granularity** can help close the gap in performance disparity.
5. Bigger/overparameterized models can **magnify/exacerbate the effects of differences in data statistics**. Quantitative biases that lead to disparity are mitigable through numerical methods.

¹We provide a re-visualization of these grouped in 6 facets by target language in Figure 4 in Appendix A.

Outline of the paper In § 2, we define our method and experimental setup. We present our results and analyses on the primary representations in § 3 and those from secondary set of controls in § 4 in a progressive manner to ease understanding. Meta analyses on fairness evaluation, non-monotonic behavior, and discussion on biases are in § 5. Additional related work is in § 6. We refer our readers to the Appendices for more detailed descriptions/discussions and reports on supplementary experiments.

2 METHOD AND DEFINITIONS

Controlled experiments as basic research for scientific understanding Using the United Nations Parallel Corpus (Ziems et al., 2016), the data from which the MT results in Junczys-Dowmunt et al. (2016) stem, we perform a series of controlled experiments on the Transformer, holding the hyperparameter settings for all 30 one-to-one language directions from the 6 languages constant. We control for size (from 10^2 to 10^6 lines) and language with respect to representational granularity. We examine 3 primary representation types — character, byte (UTF-8), and word, and upon encountering some unusual phenomena, we perform a secondary set of controls with 5 alternate representations — on the character level: Pinyin and Wubi (ASCII representations for ZH phones and character strokes, respectively), on the byte level: code page 1256 (for AR) and code page 1251 (for RU), and on the word level: Byte Pair Encoding (BPE) (Sennrich et al., 2016), an adapted compression algorithm from Gage (1994). These symbolic variants allow us to manipulate the statistical properties of the representations, while staying as “faithful” to the language as possible. We adopt this symbolic data-centric approach because we would like to more directly interpret the confounds, if any, that make language data different from other data types. We operate on a smaller data size range as this is more common in traditional domain sciences and one of our higher goals is to bridge an understanding between language sciences and engineering (the latter being the dominant focus in NLP). We run statistical tests to identify the strongest correlates of performance and to assess whether the differences between the mean performance of different groups are indeed significant. **We are concerned *not* with the absolute scores, but with the *relations* between scores from different languages and the generalizations derived therefrom.**

Information-theoretic, fair evaluation with multitexts Most sequence-to-sequence models are optimized using a cross-entropy loss (see Appendix B for definition). Cotterell et al. (2018) propose to use “renormalized” perplexity (PP) to evaluate LMs fairly using the total number of bits divided by some constant. In our case, we choose instead a simpler method of using an “unnormalized” PP, directly using the total number of bits needed to encode the development (dev) set, which has a constant size of 3,077 lines per language.

Disparity/Inequality In the context of our CLMing experiments, we consider there to be “disparity” or “inequality” between languages l_1 and l_2 if there are significant differences between the performance distributions of these two languages with respect to each representation. Here, by performance we mean the number of bits required to encode the held-out data using a trained CLM. With 30 directions, there are 15 pairs of source languages (l_{src1}, l_{src2}) and 15 pairs of target languages (l_{trg1}, l_{trg2}) possible. To assess whether the differences are significant, we perform unpaired two-sided significance tests with the null hypothesis that the score distributions for the two languages are not different. Upon testing for normality with the Shapiro-Wilk test (Shapiro & Wilk, 1965; Royston, 1995), we use the parametric unpaired two-sample Welch’s t-test (Welch, 1947) (when normal) or the non-parametric unpaired Wilcoxon test (Wilcoxon, 1945) (when not normal) for the comparisons. We use the implementation in R (R Core Team, 2014) for these 3 tests. To account for the multiple comparisons we are performing, we correct all p-values using Bonferroni’s correction (Benjamini & Heller, 2008; Dror et al., 2017) and follow Holm’s procedure² (Holm, 1979; Dror et al., 2017) to identify the pairs of l_1 and l_2 with significant differences after correction. We report all 3 levels of significance ($\alpha \leq 0.05, 0.01, 0.001$) for a more comprehensive evaluation.

Experimental setup The systematic, identical treatment we give to our data is described as follows with further preprocessing and hyperparameter details in Appendices B and C, respectively. The distinctive point of our experiment is that the training regime is the same for all (intuition in App. O.1).

²using implementation from <https://github.com/rtmdrr/replicability-analysis-NLP>

After filtering length to 300 characters maximum per line in parallel for the 6 languages, we made 3 subsets of the data with 1 million lines each — one having lines in the order of the original corpus (dataset A) and two other randomly sampled (without replacement) from the full corpus (datasets B & C). Lines in all datasets are extracted in parallel and remain fully aligned for the 6 languages. For each run and each representation, there are 30 pairwise directions (i.e. one l_{src} to one l_{trg}) that result from the 6 languages. We trained all 150 (for 5 sizes) 6-layer Transformer models for each run using the SOCKEYE Toolkit (Hieber et al., 2018). We optimize using PP and use early stopping if no PP improvement occurs after 3 checkpoints up to 50 epochs maximum, taking the best checkpoint. Characters and bytes are supposed to mitigate the out-of-vocabulary (OOV) problem on the word level. In order to assess the effect of modeling with finer granularity more precisely, all vocabulary items appearing once in the train set are accounted for (i.e. full vocabulary on train, as in Gerz et al. (2018a;b)). But we allow our system to categorize all unknown items in the dev set to be unknown (UNK) so to measure OOVs (open vocabulary on dev (Jurafsky & Martin, 2009)). To identify correlates of performance, we perform Spearman’s correlation (Spearman, 1904) with some basic statistical properties of the data (e.g. length, vocabulary size ($|V|$), type-token-ratio, OOV rate) as metrics — a complete list thereof is provided in Appendix F. For each of the 3 primary representations — character, byte, and word, we performed 5 runs total in 5 sizes (10^2 - 10^6 lines) (runs A0, B0, C0, A1, & A2) and 7 more runs in 4 sizes (10^2 - 10^5 lines) (A3-7, B1, & C1), also controlling for seeds. For the alternate/secondary representations, we ran 3 runs each in 5 sizes (10^2 - 10^6 lines) (A0, B0, & C0).

3 EXPERIMENTAL RESULTS OF PRIMARY REPRESENTATIONS

Subfigures 1a, 1b, and 1c present the mean results across 12 runs of the 3 primary representations — character, byte, and word, respectively. The x-axis represents data size in number of lines and y-axis the total conditional cross-entropy, measured in bits (Eq. 1 in Appendix B). Each line connects 5 data points corresponding to the number of bits the CLMs (trained with training data of 10^2 , 10^3 , 10^4 , 10^5 , and 10^6 lines) need to encode the target language dev set given the corresponding text in the source language. These are the same data in the same 30 language directions and 5 sizes with the same training regime, just preprocessed/segmented differently. This confirms **representation relativity** — languages (or any objects being modeled) need to be evaluated relative to their representation. “One size does not fit all” (Durrani et al., 2019), our conventional way of referring to “language” (as a socio-cultural product or with traditional word-based approaches, or even for most multilingual tasks and competitions) is too coarse-grained (see also Fisch et al. (2019) and Ponti et al. (2020)).

Subfigures 1d, 1e, and 1f display the corresponding information sorted into facets by target language, source languages represented as line types. Through these we see more clearly that results can be grouped rather neatly by target language (cf. figures sorted by source language in Appendix H) — as implicit in the Transformer’s architecture, the decoder is unaware of the source language in the encoder. As shown in Table 1 in § 5 summarizing the number of source and target language pairs with significant differences, there are **no significant differences across any source language pairs**. The Transformer neutralizes source language instances. This could explain why transfer learning or multilingual/zero-shot translation (Johnson et al., 2017) is possible at all on a conceptual level.

In general, for character and byte models, most language directions do seem to converge at 10^4 lines to similar values across all target languages, with few notable exceptions. There are some fluctuations past 10^4 , indicating further tuning of hyperparameters would be beneficial due to our present setting possibly working most favorably at 10^4 . On the character level, target language ZH (ZH_{trg}) shows a different learning pattern throughout. And on the byte level, AR_{trg} and RU_{trg} display non-monotonic and unstable behavior, which we refer to as *erratic*. Word models exhibit Double Descent across the board (note the spike at 10^3), but overall, difficult/easy languages stay consistent, with AR and RU being the hardest, followed by ES and FR, then EN and ZH. A practical takeaway from this set of experiments: in order to obtain more robust training results, use bytes for ZH (as suggested in Li et al. (2019a)) and characters for AR and RU (e.g. Lee et al. (2017)) — also if one wanted to avoid any “class” problems in performance disparity with words. Performance disparity for these representations is reported in Table 1 under “CHAR”, “BYTE”, and “WORD”. Do note, however, that the intrinsic performance of ZH with word segmentation is not particularly subpar. But this often does not correlate with its poorer downstream tasks results (recall results from Junczys-Dowmunt et al. (2016)). Since the notion of word in ZH is highly contested and

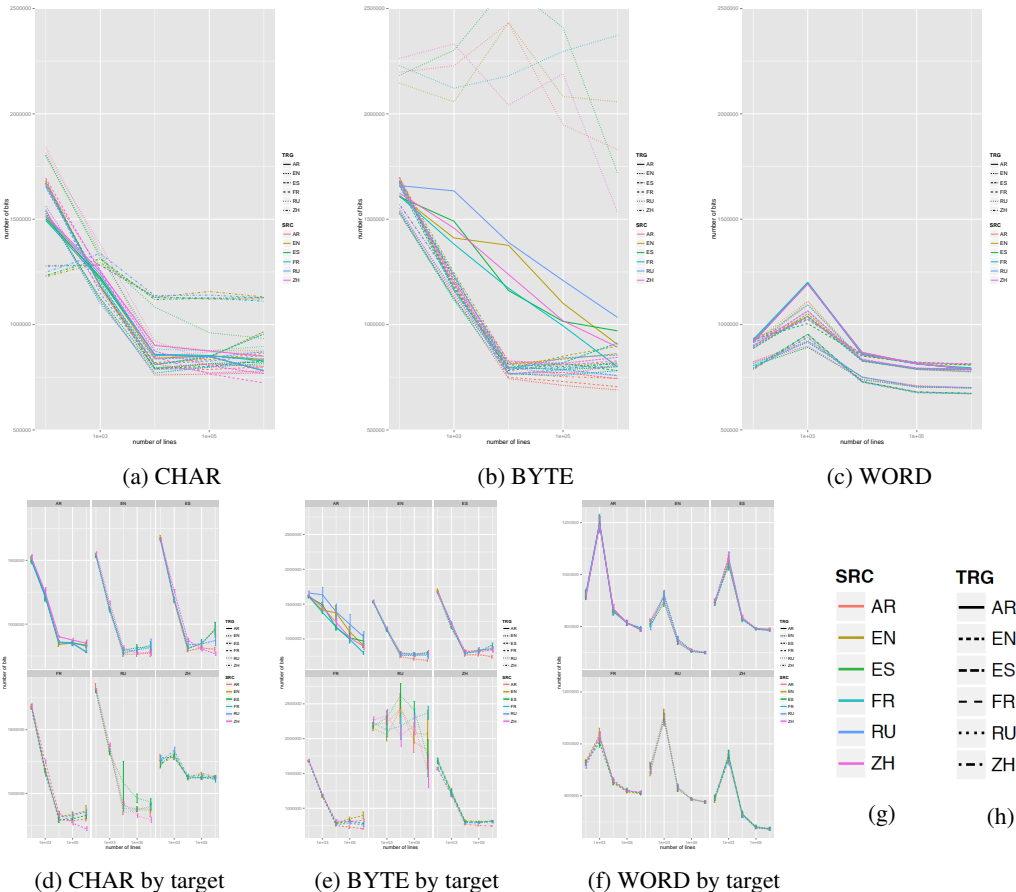


Figure 1: Number of bits (the lower the better) as a function of data size plotted for all 30 directions. Subfigures 1d, 1e, and 1f depict the corresponding information as in 1a, 1b, and 1c (showing mean across 12 runs), respectively, but sorted in 6 facets by target language and with error bars. Legend in Subfigure 1g shows the correspondence between colors and source languages, in Subfigure 1h between line types and target languages. (These figures are also shown enlarged in Appendix G.)

ambiguous — 1) it is often aimed to align with that in other languages so to accommodate manual feature engineering and academic theories, 2) there is great variation among different conventions, 3) native ZH speakers identify characters as words — there are reasons to rethink this procedure now that fairer and language-independent processing in finer granularity is possible (cf. Li et al. (2019b) as well as Duanmu (2017) for a summary on the contested nature of wordhood in ZH). A more native analysis of ZH, despite being considered a high-resource language, has not yet been recognized in NLP.

4 UNDERSTANDING THE PHENOMENA WITH ALTERNATE REPRESENTATIONS

To understand why some languages show different results than others, we carried out a secondary set of control experiments with representations targeting the problematic statistical properties of the corresponding target languages. (An extended version of this section is provided in Appendix P.)

Character level We reduced the high $|V|$ in ZH with representations in ASCII characters — Pinyin and Wubi. The former is a romanization of ZH characters based on their pronunciations and the latter an input algorithm that decomposes character-internal information into stroke shape and ordering and matches these to 5 classes of radicals (Lunde, 2008). We replaced the ZH data in these formats *only on the target side* and reran the experiments involving ZH_{trg} on the character level. Results in Figure 2 and Table 1 show that the elimination of disparity on character level is possible if ZH is represented

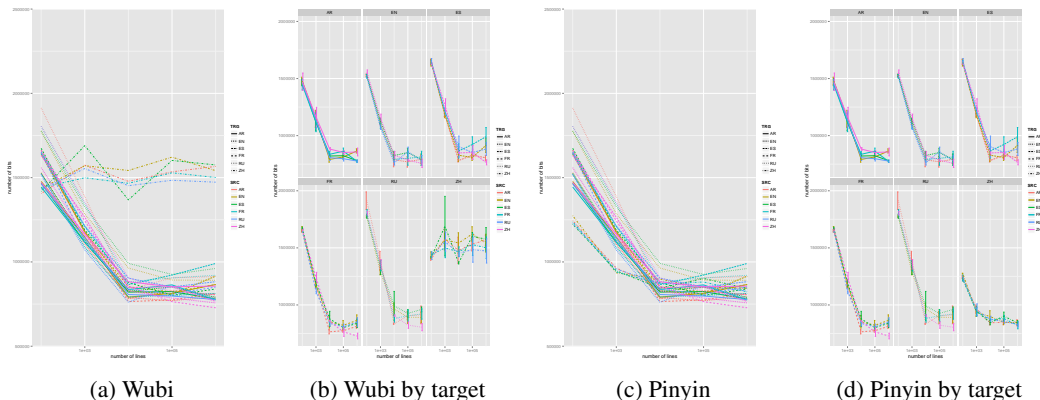


Figure 2: Character-level remedies for ZH: Wubi vs. Pinyin.

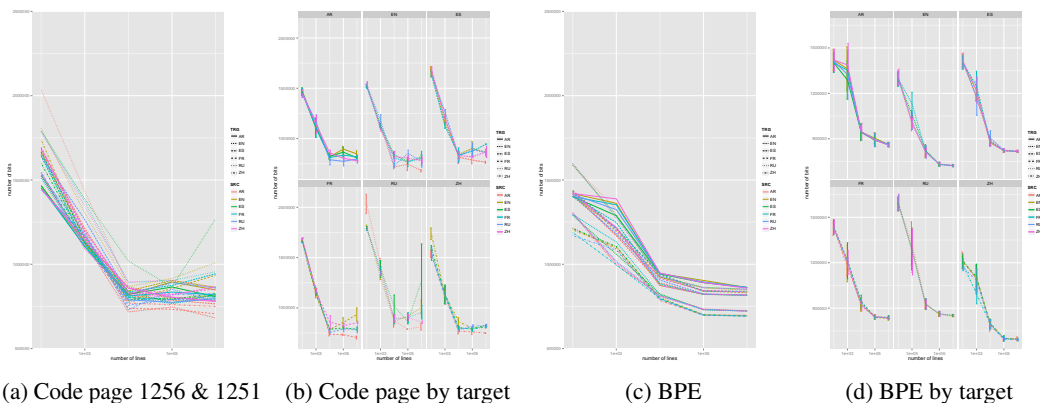


Figure 3: Byte-level (Subfigures 3a & 3b) remedies with code page 1256 for target AR and 1251 for target RU, and word-level (Subfigures 3c & 3d) remedy with BPE for all languages.


through Pinyin (transliteration), as in Subfigure 2c. But models with ZH logographic scripts display a behavioral tendency unlike those with other (phonetic) alphabetic scripts (Subfigure 2a). Work published thus far using Wubi with the Transformer seems to have needed some form of architectural modification (Gao et al., 2020) or a different architecture altogether (Nikolov et al., 2018; Zhang et al., 2019), suggesting a possible script bias (to be further discussed in § 5 under “Basis for biases”).

Byte level Length is the most salient statistical attribute that makes AR and RU outliers. To shorten their sequence lengths, we tested with alternate encodings on AR_{trg} and RU_{trg} — code page 1256 and 1251, which provide 1-byte encodings specific to AR and RU, respectively. Results are shown in Subfigures 3a and 3b. Not only is erraticity resolved, the number of 15 possible target language pairs with significant differences reduces from 8 with the UTF-8 byte representation to 0 (Table 1 under “ $ARRU_t$ ”), indicating that we eliminated disparity with this optimization heuristic. Since our heuristic is a lossless and reversible transform, it shows that **a complexity that is intrinsic and necessary in language³ does not exist** in computing, however diverse they may be, as our 6 are, from the conventional linguistic typological, phylogenetic, historical, or geographical perspectives. Please refer to Appendix J for our discussion on language complexity.

Word level The main difference between word and character/byte models is length not being a top contributing factor correlating with performance, but instead $|V|$ is. This is understandable as word segmentation neutralizes sequence lengths. To remedy the OOV problem, we use BPE, which learns a fixed vocabulary of variable-length character sequences (on word level, as it presupposes word

³aside from its statistical properties related to length and vocabulary. “Language” here refers to language represented through all representations.

Table 1: Number of language pairs out of 15 with significant differences, with respective p-values. $ARRU_t$ refers to AR & RU being optimized only on the target side; whereas $ARRU_{s,t}$ denotes optimization on both source and target sides (relevant for directions AR-RU and RU-AR).

p-value	CHAR		Pinyin		Wubi		BYTE		$ARRU_t$		$ARRU_{s,t}$		WORD		BPE	
	src	trg	src	trg	src	trg	src	trg	src	trg	src	trg	src	trg	src	trg
0.05	0	7	0	4	0	8	0	9	0	4	0	4	0	11	0	10
0.01	0	5	0	2	0	6	0	8	0	3	0	4	0	8	0	8
 0.001	0	3	0	0	0	5	0	8	0	0	0	2	0	8	0	7

segmentation) from the training data. It is more fine-grained than word segmentation and is known for its capability to model subword units for morphologically complex languages (e.g. AR and RU). We use the same vocabulary of 30,000 as specified in Junczys-Dowmunt et al. (2016). This reduced our averaged OOV token rate by 89-100% across the 5 sizes. The number of language pairs with significant differences reduced to 7 from 8 for word models, showing how **finer-grained modeling has a positive effect on closing the disparity gap**.

5 META-RESULTS, ANALYSIS, AND DISCUSSION

Performance disparity Table 1 lists the number of language pairs with significant differences under the representations studied. Considering how it is **possible** for our character and byte models to effect no performance disparity for the same languages on the same data, this indicates that disparity is not a necessary condition. In fact, the customary expectation that languages ought to perform differently stems from our word segmentation practice. Furthermore, the order of AR/RU > ES/FR > EN/ZH (Figure 1c) resembles the idea of morphological complexity. Considering there are character-internal meaningful units in languages with logographic script such as ZH (cf. Zhang & Komachi (2018)) that are rarely captured, studied, or referred to as “morphemes”, this goes to show that linguistic morphology, along with its complexity, as is practiced today⁴ and that which has occurred in the NLP discourse thus far, has only been relevant on and is bounded to the “word” level. The definition of word, however, has been recognized as problematic for a very long time in the language sciences (see Haspelmath (2011) and references therein from the past century). Since the conventional notion of word, which has been centered on English and languages with alphabetic scripts, has a negative impact on languages both morphologically rich (see Minkov et al. (2007), Seddah et al. (2010), inter alia), AR and RU in our case, as well as morphologically “frugal” (Koehn, 2005), as in ZH, finer-grained modeling with characters and bytes (or n-gram variants/pieces thereof) is indeed a more sensible option and enables a greater variety of languages to be handled with more simplicity, fairness, independence, and flexibility.

While the lack of significant differences between pairs of source languages would signify neutralization of source language instances, it does not mean that source languages have no effect on target. For our byte solutions with code pages, we experimented also with source side optimization in the directions that involve AR/RU as source. This affected the distribution of the disparity results for that representation — with 2 pairs being significantly different (see Table 1 under “ $ARRU_{s,t}$ ”). We defer further investigation on the nature of source language neutralization to future work.

Sample-wise Double Descent (DD) Sample-wise non-monotonicity/DD (Nakkiran et al., 2020) denotes a degradation followed by an improvement in performance with increasing data size. We notice word models and character models with ZH_{trg} , i.e. models with high target $|V|$, are prone to exhibit a spike at 10^3 . A common pattern for these is the **ratio of target training token count to number of parameters** falls into $O(10^{-4})$ for 10^2 lines, $O(10^{-3})$ at 10^3 , $O(10^{-2})$ at 10^4 , and $O(10^{-1})$ for 10^5 lines and so on. But for more atomic units such as alphabetic (not logographic) characters (may it be Latin, Cyrillic, or Abjad) and for bytes, this progression instead begins at $O(10^{-3})$ at 10^2 lines. Instead of thinking this spike of 10^3 as irregular, we may instead want to

⁴But there are no reasons why linguistics or linguistic typology cannot encompass a statistical science of language beyond/without “words”, or with continuous representations of characters and bytes. In fact, that could complement the needs of language engineering and the NNs/DL/ML communities better.

think of this learning curve as shifted by 1 order of magnitude to the right for characters and bytes and/or the performance at 10^2 lines for words and ZH-characters due to being overparameterized and hence abnormal. This would also fit in with the findings by Belkin et al. (2019) and Nakkiran et al. (2020) attributing DD to overparameterization. If we could use this ratio and logic of higher $|V|$ to automatically detect “non-atomic” units, ones that can be further decomposed, this observation could potentially be beneficial for advancing other sciences, e.g. biology. From a cognitive modeling perspective, the similarity in behavior of ZH characters and words of other languages can affirm the interpretation of wordhood for those ZH speakers who identify ZH characters as words (see also last paragraph in § 3 and Appendix J). While almost all work attribute DD to algorithmic reasons, concurrent work by Chen et al. (2020) corroborates our observation and confirms that DD arises due to “the interaction between the properties of the data and the inductive biases of learning algorithms”. Other related work on DD and its more recent development can also be found in their work.

We performed additional experiments testing our setting on the datasets used by the Nakkiran et al. (2020) and testing our data on a non-neural LM. Results support our findings and are provided in Appendix K. Number of model parameters can be found in Appendix L.

Erraticity We observe another type of sample-wise non-monotonicity, one that signals irregular and unstable performance across data sizes and runs. Within one run, erraticity can be observed directly as changes in direction on the y-axis. Across runs, large variance can be observed, even with the same dataset (see Figure 18 in Appendix M). Erraticity can also be observed indirectly through a negative correlation between data size and performance. Many work on length bias in NMT have focused on solutions related to search, e.g. Murray & Chiang (2018). Our experiments show that a kind of length bias can surface already with CLMing, without generation taking place. If the connection between erraticity and length bias can indeed be drawn, it could strengthen the case for global conditioning (Soutsov & Sarawagi, 2016). (See Appendix M for more discussion and results.)

Script bias, erraticity, word bias — are these necessary conditions? To assess whether the observed phenomena are particular to this one setting, we performed one run with dataset A in 4 sizes with the primary representations on 1-layer Transformers (see Appendix N). We observed no significant disparity across the board. It shows that **larger/overparameterized models can magnify/exacerbate the differences in the data statistics**. That hyperparameter tuning — in this case, through the reduction of the number of layers — can mitigate effects from data statistics is, to the best of our knowledge, a novel insight, suggesting also that a general expectation of monotonic development as data size increases can indeed be held. Our other findings remain consistent (representational relativity, source language neutralization, and DD on word level).

Bases for biases Recall in § 1, we “consider *bias* to be present when performance disparity in our Transformer models is statistically significant”. As shown in our data statistics and analysis (Appendices D and P respectively), script bias, length bias wrt erraticity in CLMing, and word bias are all evident in the vocabulary and length information in the data statistics. Hence these disparities in performance are really a result of the Transformer being able to model these **differences in data** at such a magnitude that the differences are statistically significant. The meta phenomenon of erraticity, however, warrants an additional consideration indicative of the **empirical limits of our compute** (cf. Xu et al. (2020)), even when the non-monotonicity is not observed during the training of each model.

In eliminating performance disparity in character and byte models by normalizing vocabulary and length statistics in the data, we demonstrated that performance disparity as expected from the morphological complexity hierarchy is due to word tokenization, not intrinsic or necessary in language. This is the word bias. Qualitative issues in the concept of word will persist and make crosslinguistic comparison involving “words” unfair even if one were to be able to find a quantitative solution to mitigate the OOV issue, the bottleneck in word-based processing. We humans have a choice in how we see/process languages. That some might still prefer to continue with a crosslinguistic comparison with “words” and exert the superiority of “word” tokenization speaks for a view that is centered on “privileged” languages — in that case, **word bias is a human bias**.

And, in eliminating performance disparity across the board with our one-layer models, we show that all quantitative differences in data statistics between languages can also be modeled in a “zoomed-

out”/“desensitized” mode, suggesting that while languages can be perceived as being fundamentally different in different ways in different granularities, they can also be viewed as fundamentally similar.

6 ADDITIONAL RELATED WORK

Similar to our work in testing for hardness are Cotterell et al. (2018), Mielke et al. (2019), and Bugliarello et al. (2020). The first two studied (monolingual) LMs — the former tested on the Europarl languages (Koehn, 2005) with n-gram and character models and concluded that morphological complexity was the culprit to hardness, the latter studied 62 languages of the Bible corpus (Mayer & Cysouw, 2014) in addition and refuted the relevance of linguistic features in hardness based on character and BPE models on both corpora in word-tokenized form. Bugliarello et al. (2020) compared translation results of the Europarl languages with BPEs at one data size and concluded that it is easier to translate out of EN than into it, statistical significance was, however, not assessed. In contrast, we ablated away the confound of generation and studied CLMing with controls with a broader range of languages with more diverse statistical profiles in 3 granularities and up to 5 orders of magnitude in data size. That basic data statistics are the driver of success in performance in multilingual modeling has so far only been explicitly argued for in Mielke et al. (2019). We go beyond their work in monolingual LMs to study CLMs and evaluate also in relation to data size, representational granularity, and quantitative and qualitative fairness.

Bender (2009) advocated the relevance of linguistic typology for the design of language-independent NLP systems based on crosslinguistic differences in word-based structural notions, such as parts of speech. Ponti et al. (2019) found typological information to be beneficial in the few-shot setting on the character level for 77 languages with Latin scripts. But no multilingual work has thus far explicitly examined the relation between linguistic typology and the statistical properties of the data, involving languages with diverse statistical profiles in different granularities.

As obtaining training data is often the most difficult part of an NLP or Machine Learning (ML) project, Johnson et al. (2018) introduced an extrapolation methodology to directly model the relation between data size and performance. Our work can be viewed as one preliminary step towards this goal. To the best of our knowledge, there has been no prior work on demonstrating the neutralization of source language instances through statistical comparisons, a numerical analysis on DD for sequence-to-sequence models, the meta phenomenon of a sample-wise non-monotonicity (erraticity) being related to length, or the connection between effects of data statistics and modification in architectural depth.

7 CONCLUSION

Summary We performed a novel, rigorous relational assessment of performance disparity across different languages, representations, and data sizes in CLMing with the Transformer. Different disparity patterns were observed on different representation types (character, byte, and word), which can be traced back to the data statistics. The disparity pattern reflected on the word level corresponds to the morphological complexity hierarchy, reminding us that the definition of morphology is predicated on the notion of word and indicating how morphological complexity can be modeled by the Transformer simply through word segmentation. As we were able to eliminate disparity on the same data on the character and byte levels by normalizing length and vocabulary, we showed that morphological complexity is not a necessary concept but one that results from word segmentation and is bounded to the word level, orthogonal to the performance of character or byte models. Representational units of finer granularity were shown to help eliminate performance disparity though at the cost of longer sequence length, which can have a negative impact on robustness. In addition, we found all word models and character models with ZH_{trg} to behave similarly in their being prone to exhibit a peak (as sample-wise DD) around 10^3 lines in our setting. While bigger/overparameterized models can magnify the effect of data statistics, exacerbating the disparity, we found a decrease in model depth can eliminate these quantitative biases, leaving only the qualitative aspect of “word” and the necessity of word segmentation in question.

Outlook Machine learning has enabled greater diversity in NLP (Joshi et al., 2020). Fairness, in the elimination of disparity, does not require big data. This paper made a pioneering attempt to bridge research in DL/NNs, language sciences, and language engineering through a data-centric perspective.

We believe a **statistical** science for NLP as a data science can well complement algorithmic analyses with an empirical view contributing to a more generalizable pool of knowledge for NNs/DL/ML. A more comprehensive study not only can lead us to new scientific frontiers, but also better design and evaluation, benefitting the development of a more general, diverse and inclusive Artificial Intelligence.

REFERENCES

- Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019. ISSN 0027-8424. doi: 10.1073/pnas.1903070116. URL <https://www.pnas.org/content/116/32/15849>.
- E. M. Bender. *Linguistic Fundamentals for Natural Language Processing: 100 Essentials from Morphology and Syntax*. 2013.
- Emily M. Bender. Linguistically naïve != language independent: Why NLP needs linguistic typology. In *Proceedings of the EACL 2009 Workshop on the Interaction between Linguistics and Computational Linguistics: Virtuous, Vicious or Vacuous?*, pp. 26–32, Athens, Greece, March 2009. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/W09-0106>.
- Yoav Benjamini and Ruth Heller. Screening for partial conjunction hypotheses. *Biometrics*, 64(4): 1215–1222, 2008. ISSN 0006341X, 15410420. URL <http://www.jstor.org/stable/25502204>.
- Christian Bentz, Tatyana Ruzsics, Alexander Koplenig, and Tanja Samardžić. A comparison between morphological complexity measures: Typological data vs. language corpora. In *Proceedings of the Workshop on Computational Linguistics for Linguistic Complexity (CLALC)*, pp. 142–153, Osaka, Japan, December 2016. The COLING 2016 Organizing Committee. URL <https://www.aclweb.org/anthology/W16-4117>.
- Emanuele Bugliarello, Sabrina J. Mielke, Antonios Anastasopoulos, Ryan Cotterell, and Naoaki Okazaki. It’s easier to translate out of English than into it: Measuring neural translation difficulty by cross-mutual information. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 1640–1649, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.149. URL <https://www.aclweb.org/anthology/2020.acl-main.149>.
- Mauro Cettolo, Christian Girardi, and Marcello Federico. Wit³: Web inventory of transcribed and translated talks. In *Proceedings of the 16th Conference of the European Association for Machine Translation (EAMT)*, pp. 261–268, Trento, Italy, May 2012.
- Lin Chen, Yifei Min, Mikhail Belkin, and Amin Karbasi. Multiple descent: Design your own generalization curve, 2020.
- Stanley F. Chen and Joshua Goodman. An empirical study of smoothing techniques for language modeling. *Computer Speech Language*, 13(4):359 – 394, 1999. ISSN 0885-2308. doi: <https://doi.org/10.1006/csla.1999.0128>. URL <http://www.sciencedirect.com/science/article/pii/S0885230899901286>.
- Tianqi Chen, Mu Li, Yutian Li, Min Lin, Naiyan Wang, Minjie Wang, Tianjun Xiao, Bing Xu, Chiyuan Zhang, and Zheng Zhang. Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems. *arXiv preprint arXiv:1512.01274*, 2015.
- Colin Cherry, George Foster, Ankur Bapna, Orhan Firat, and Wolfgang Macherey. Revisiting character-based neural machine translation with capacity and compression. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 4295–4305. Association for Computational Linguistics, 2018. URL <http://aclweb.org/anthology/D18-1461>.

- Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of neural machine translation: Encoder–decoder approaches. In *Proceedings of SSST-8, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation*, pp. 103–111, Doha, Qatar, October 2014. Association for Computational Linguistics. doi: 10.3115/v1/W14-4012. URL <https://www.aclweb.org/anthology/W14-4012>.
- Ryan Cotterell, Sabrina J. Mielke, Jason Eisner, and Brian Roark. Are all languages equally hard to language-model? In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pp. 536–541, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2085. URL <https://www.aclweb.org/anthology/N18-2085>.
- Rotem Dror, Gili Baumer, Marina Bogomolov, and Roi Reichart. Replicability analysis for natural language processing: Testing significance with multiple datasets. *Transactions of the Association for Computational Linguistics*, 5:471–486, 2017. URL <http://aclweb.org/anthology/Q17-1033>.
- San Duanmu. Word and wordhood, modern. In Rint Sybesma (ed.), *Encyclopedia of Chinese Language and Linguistics*, pp. 543–549. Brill, 2017.
- Nadir Durrani, Fahim Dalvi, Hassan Sajjad, Yonatan Belinkov, and Preslav Nakov. One size does not fit all: Comparing NMT representations of different granularities. 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)*, pp. 1504–1516, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1154. URL <https://www.aclweb.org/anthology/N19-1154>.
- Adam Fisch, Jiang Guo, and Regina Barzilay. Working hard or hardly working: Challenges of integrating typology into neural dependency parsers. 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)*, pp. 5713–5719, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1574. URL <https://www.aclweb.org/anthology/D19-1574>.
- Philip Gage. A new algorithm for data compression. *C Users J.*, 12(2):23–38, February 1994. ISSN 0898-9788. URL <http://dl.acm.org/citation.cfm?id=177910.177914>.
- Yingqiang Gao, Nikola I. Nikolov, Yuhuang Hu, and Richard H.R. Hahnloser. Character-level translation with self-attention. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 1591–1604, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.145. URL <https://www.aclweb.org/anthology/2020.acl-main.145>.
- Martin Gellerstam. Translationese in Swedish novels translated from English. In Lars Wollin and Hans Lindquist (eds.), *Translation Studies in Scandinavia*, pp. 88–95. CWK Gleerup, 1986.
- Daniela Gerz, Ivan Vulić, Edoardo Ponti, Jason Naradowsky, Roi Reichart, and Anna Korhonen. Language modeling for morphologically rich languages: Character-aware modeling for word-level prediction. *Transactions of the Association for Computational Linguistics*, 6:451–465, 2018a. doi: 10.1162/tacl_a_00032. URL <https://www.aclweb.org/anthology/Q18-1032>.
- Daniela Gerz, Ivan Vulić, Edoardo Maria Ponti, Roi Reichart, and Anna Korhonen. On the relation between linguistic typology and (limitations of) multilingual language modeling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 316–327, Brussels, Belgium, October–November 2018b. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/D18-1029>.
- Martin Haspelmath. The indeterminacy of word segmentation and the nature of morphology and syntax. *Folia Linguistica*, 2011.

- Kenneth Heafield. KenLM: Faster and smaller language model queries. In *Proceedings of the Sixth Workshop on Statistical Machine Translation*, pp. 187–197, Edinburgh, Scotland, July 2011. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/W11-2123>.
- Kenneth Heafield, Ivan Pouzyrevsky, Jonathan H. Clark, and Philipp Koehn. Scalable modified Kneser-Ney language model estimation. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 690–696, Sofia, Bulgaria, August 2013. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/P13-2121>.
- Felix Hieber, Tobias Domhan, Michael Denkowski, David Vilar, Artem Sokolov, Ann Clifton, and Matt Post. The Sockeye neural machine translation toolkit at AMTA 2018. In *Proceedings of the 13th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Papers)*, pp. 200–207. Association for Machine Translation in the Americas, 2018. URL <http://aclweb.org/anthology/W18-1820>.
- Sture Holm. A simple sequentially rejective multiple test procedure. *Scandinavian Journal of Statistics*, 6(2):65–70, 1979. ISSN 03036898, 14679469. URL <http://www.jstor.org/stable/4615733>.
- Mark Johnson, Peter Anderson, Mark Dras, and Mark Steedman. Predicting accuracy on large datasets from smaller pilot data. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 450–455. Association for Computational Linguistics, 2018. URL <http://aclweb.org/anthology/P18-2072>.
- Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. Google’s multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics*, 5:339–351, 2017. doi: 10.1162/tacl_a_00065. URL <https://www.aclweb.org/anthology/Q17-1024>.
- Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 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*, pp. 6282–6293, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.560. URL <https://www.aclweb.org/anthology/2020.acl-main.560>.
- Marcin Junczys-Dowmunt, Tomasz Dwojak, and Hieu Hoang. Is neural machine translation ready for deployment? A case study on 30 translation directions. In *IWSLT 2016, Seattle*, October 2016. URL <https://www.microsoft.com/en-us/research/publication/neural-machine-translation-ready-deployment-case-study-30-translation-directions/>.
- Daniel Jurafsky and James H. Martin. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. Pearson Prentice Hall, second edition, 2009.
- Reinhard Kneser and Hermann Ney. Improved backing-off for m-gram language modeling. In *1995 International Conference on Acoustics, Speech, and Signal Processing*, volume 1, pp. 181–184. IEEE, 1995.
- Philipp Koehn. Europarl: A Parallel Corpus for Statistical Machine Translation. In *Conference Proceedings: the tenth Machine Translation Summit*, pp. 79–86, Phuket, Thailand, 2005. AAMT, AAMT. URL <http://mt-archive.info/MTS-2005-Koehn.pdf>.
- Philipp Koehn and Rebecca Knowles. Six challenges for neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation*, pp. 28–39. Association for Computational Linguistics, 2017. doi: 10.18653/v1/W17-3204. URL <http://aclweb.org/anthology/W17-3204>.

- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, and Evan Herbst. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pp. 177–180. Association for Computational Linguistics, 2007. URL <http://aclweb.org/anthology/P07-2045>.
- Jason Lee, Kyunghyun Cho, and Thomas Hofmann. Fully character-level neural machine translation without explicit segmentation. *Transactions of the Association for Computational Linguistics*, 5: 365–378, December 2017. doi: 10.1162/tacl_a_00067. URL <https://www.aclweb.org/anthology/Q17-1026>.
- Bo Li, Yu Zhang, Tara Sainath, Yonghui Wu, and William Chan. Bytes are all you need: End-to-end multilingual speech recognition and synthesis with bytes. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5621–5625. IEEE, 2019a.
- Xiaoya Li, Yuxian Meng, Xiaofei Sun, Qinghong Han, Arianna Yuan, and Jiwei Li. Is word segmentation necessary for deep learning of Chinese representations? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 3242–3252, Florence, Italy, July 2019b. Association for Computational Linguistics. doi: 10.18653/v1/P19-1314. URL <https://www.aclweb.org/anthology/P19-1314>.
- Ken Lunde. *CJKV Information Processing*. O’Reilly Media, Inc., 2nd edition, 2008. ISBN 0596514476, 9780596514471.
- Thomas Mayer and Michael Cysouw. Creating a massively parallel Bible corpus. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC-2014)*, pp. 3158–3163, Reykjavik, Iceland, May 2014. European Languages Resources Association (ELRA).
- Sabrina J. Mielke, Ryan Cotterell, Kyle Gorman, Brian Roark, and Jason Eisner. What kind of language is hard to language-model? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4975–4989, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1491. URL <https://www.aclweb.org/anthology/P19-1491>.
- Einat Minkov, Kristina Toutanova, and Hisami Suzuki. Generating complex morphology for machine translation. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pp. 128–135, Prague, Czech Republic, June 2007. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/P07-1017>.
- Kenton Murray and David Chiang. Correcting length bias in neural machine translation. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pp. 212–223, Belgium, Brussels, October 2018. Association for Computational Linguistics. doi: 10.18653/v1/W18-6322. URL <https://www.aclweb.org/anthology/W18-6322>.
- Preetum Nakkiran. More data can hurt for linear regression: Sample-wise double descent, 2019.
- Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep double descent: Where bigger models and more data hurt. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=Blg5sA4twr>.
- Nikola Nikolov, Yuhuang Hu, Mi Xue Tan, and Richard H.R. Hahnloser. Character-level Chinese-English translation through ASCII encoding. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pp. 10–16, Belgium, Brussels, October 2018. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/W18-6302>.
- Franz Josef Och and Hermann Ney. A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29(1):19–51, 2003. doi: 10.1162/089120103321337421. URL <https://www.aclweb.org/anthology/J03-1002>.

- M Opper, W Kinzel, J Kleinz, and R Nehl. On the ability of the optimal perceptron to generalise. *Journal of Physics A: Mathematical and General*, 23(11):L581–L586, jun 1990. doi: 10.1088/0305-4470/23/11/012. URL <https://doi.org/10.1088%2F0305-4470%2F23%2F11%2F012>.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of NAACL-HLT 2019: Demonstrations*, 2019.
- Edoardo Maria Ponti, Ivan Vulić, Ryan Cotterell, Roi Reichart, and Anna Korhonen. Towards zero-shot language modeling. 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)*, pp. 2900–2910, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1288. URL <https://www.aclweb.org/anthology/D19-1288>.
- Edoardo Maria Ponti, Helen O’Horan, Yevgeni Berzak, Ivan Vulić, Roi Reichart, Thierry Poibeau, Ekaterina Shutova, and Anna Korhonen. Modeling language variation and universals: A survey on typological linguistics for natural language processing, 2020.
- R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2014. URL <http://www.R-project.org/>.
- Patrick Royston. Remark as r94: A remark on algorithm as 181: The w-test for normality. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 44(4):547–551, 1995. ISSN 00359254, 14679876. URL <http://www.jstor.org/stable/2986146>.
- Djame Seddah, Sandra Koebler, and Reut Tsarfaty (eds.). *Proceedings of the NAACL HLT 2010 First Workshop on Statistical Parsing of Morphologically-Rich Languages*, Los Angeles, CA, USA, June 2010. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/W10-1400>.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 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)*, pp. 1715–1725, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/P16-1162. URL <https://www.aclweb.org/anthology/P16-1162>.
- S. S. Shapiro and M. B. Wilk. An analysis of variance test for normality (complete samples)†. *Biometrika*, 52(3-4):591–611, 12 1965. ISSN 0006-3444. doi: 10.1093/biomet/52.3-4.591. URL <https://doi.org/10.1093/biomet/52.3-4.591>.
- Pavel Sountsov and Sunita Sarawagi. Length bias in encoder decoder models and a case for global conditioning. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 1516–1525, Austin, Texas, November 2016. Association for Computational Linguistics. doi: 10.18653/v1/D16-1158.
- C. Spearman. The proof and measurement of association between two things. *The American Journal of Psychology*, 15(1):72–101, 1904. ISSN 00029556. URL <http://www.jstor.org/stable/1412159>.
- Felix Stahlberg and Bill Byrne. On NMT search errors and model errors: Cat got your tongue? 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)*, pp. 3356–3362, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1331. URL <https://www.aclweb.org/anthology/D19-1331>.
- Reut Tsarfaty, Djame Seddah, Yoav Goldberg, Sandra Kuebler, Yannick Versley, Marie Candito, Jennifer Foster, Ines Rehbein, and Lamia Tounsi. Statistical parsing of morphologically rich languages (SPMRL) what, how and whither. In *Proceedings of the NAACL HLT 2010 First Workshop on Statistical Parsing of Morphologically-Rich Languages*, pp. 1–12, Los Angeles, CA, USA, June 2010. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/W10-1401>.

- Reut Tsarfaty, Djamé Seddah, Sandra Kübler, and Joakim Nivre. Parsing morphologically rich languages: Introduction to the special issue. *Computational Linguistics*, 39(1):15–22, 2013. doi: 10.1162/COLI_a_00133. URL <https://www.aclweb.org/anthology/J13-1003>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 30*, pp. 5998–6008. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf>.
- B. L. Welch. The Generalization of ‘Student’s’ Problem when Several Different Population Variances are Involved. *Biometrika*, 34(1-2):28–35, 01 1947. ISSN 0006-3444. doi: 10.1093/biomet/34.1-2.28. URL <https://doi.org/10.1093/biomet/34.1-2.28>.
- Frank Wilcoxon. Individual comparisons by ranking methods. *Biometrics Bulletin*, 1(6):80–83, 1945. ISSN 00994987. URL <http://www.jstor.org/stable/3001968>.
- Yilun Xu, Shengjia Zhao, Jiaming Song, Russell Stewart, and Stefano Ermon. A theory of usable information under computational constraints. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=r1eBeyHFDH>.
- Longtu Zhang and Mamoru Komachi. Neural machine translation of logographic language using sub-character level information. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pp. 17–25, Belgium, Brussels, October 2018. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/W18-6303>.
- Wei Zhang, Feifei Lin, Xiaodong Wang, Zhenshuang Liang, and Zhen Huang. Subcharacter Chinese-English neural machine translation with Wubi encoding, 2019.
- Michał Ziemiński, Marcin Junczys-Dowmunt, and Bruno Poulliquen. The United Nations Parallel Corpus v1.0. In Nicoletta Calzolari (Conference Chair), Khalid Choukri, Thierry Declerck, Sara Goggi, Marko Grobelnik, Bente Maegaard, Joseph Mariani, Helene Mazo, Asuncion Moreno, Jan Odijk, and Stelios Piperidis (eds.), *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*, Paris, France, may 2016. European Language Resources Association (ELRA). ISBN 978-2-9517408-9-1.

APPENDICES

A	RE-VISUALIZATION OF FIGURE 1 IN JUNCZYS-DOWMUNT ET AL. (2016) IN 6 FACETS BY TARGET LANGUAGE	17
B	DATA SELECTION AND PREPROCESSING DETAILS	17
C	HYPERPARAMETER SETTING	19
D	DATA STATISTICS	20
E	SCORE TABLES	24
F	CORRELATION STATISTICS	25
G	ENLARGED FIGURES FOR ALL 30 LANGUAGE DIRECTIONS (AGGREGATE RESULTS FROM ALL RUNS)	26
H	SAMPLE FIGURES FROM RUN A0, ALSO SORTED BY SOURCE LANGUAGE FOR CONTRAST	42
I	LANGUAGE PAIRS WITH SIGNIFICANT DIFFERENCES	43
J	LANGUAGE COMPLEXITY	44
K	SAMPLE-WISE DOUBLE DESCENT (DD)	46
	K.1 OUR EXPERIMENTAL FRAMEWORK ON DD DATASETS FROM (NAKKIRAN ET AL., 2020)	46
	K.2 TOKEN-TO-PARAMETER RATIO FOR NON-NEURAL MONOLINGUAL LMS	47
L	NUMBER OF MODEL PARAMETERS	50
M	ERRATICITY	53
	M.1 ERRATICITY AS LARGE VARIANCE: EVIDENCE FROM DIFFERENT RUNS OF THE SAME DATA	53
	M.2 ADDITIONAL EXPERIMENT WITH LENGTH FILTERING TO 300 BYTES	53
N	EXPERIMENTS WITH ONE-LAYER TRANSFORMER	56
O	PAQS (PREVIOUSLY ASKED QUESTIONS)	57
	O.1 ONE SETTING FOR ALL	57
	O.2 TRANSLATIONESE / WORD ORDER	57
P	UNDERSTANDING THE PHENOMENA WITH ALTERNATE REPRESENTATIONS (EXTENDED VERSION)	59

A RE-VISUALIZATION OF FIGURE 1 IN JUNCZYS-DOWMUNT ET AL. (2016) IN 6 FACETS BY TARGET LANGUAGE

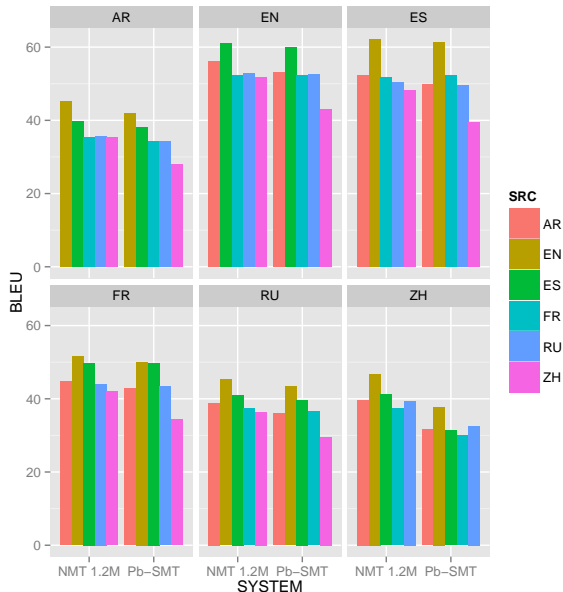


Figure 4: Results of the Moses baseline systems (right group in each facet) and neural models (left) with 1.2 million iterations (1 iteration corresponds to 1 mini-batch) for the 30 directions of the 6-way UN corpus, tokenized (ZH segmented), lowercased, and length filtered to 100 BPE tokens.

B DATA SELECTION AND PREPROCESSING DETAILS

The UN Parallel Corpus v1.0 (Ziemski et al., 2016) consists of manually translated UN documents from 1990 to 2014 in the 6 official UN languages. Therein is a subcorpus that is fully aligned by line, comprising the 6-way parallel corpus we use. We tried to have as little preprocessing or filtering as necessary to eliminate possible confounds. But as the initial runs of our experiment failed due to insufficient memory on a single GPU with 12 GB VRAM⁵, we filtered out lines with more than 300 characters in any language in lockstep with one another for all the 6 languages such that the subcorpora would remain parallel, thereby keeping the material of each language semantically equivalent to one another. 8,944,859 lines for each language were retained as our training data which cover up to the 75th percentile in line length for all 6 languages. In order to monitor the effect of data size, we made subcorpora of each language in 5 sizes by heading the first 10^2 , 10^3 , 10^4 , 10^5 , 10^6 lines⁶. We refer to this as dataset A. In addition, to better understand and verify the consistency of the phenomena observed, we made 2 supplemental datasets by shuffling the 8,944,859 lines two different times randomly and heading the number of lines in our 5 sizes for each language, again in lockstep with one another (datasets B and C).

⁵GPUs used for experiments in this paper range from a NVIDIA TITAN RTX (24 GB), NVIDIA GeForce RTX 2080 Ti (11 GB), a GTX Titan X (12 GB), to a GTX 1080 (8 GB). All jobs were run on a single GPU setting. Some word-level experiments involving AR_{trg} or RU_{trg} at 10^6 had to be run on a CPU as 24 GB VRAM were not sufficient. Models with higher maximum sequence lengths (e.g. byte models) were trained with 24 GB VRAM. Difference in equipment does not necessarily lead to degradation/improvement in scores.

⁶The terms “line” and “sentence” have been used interchangeably in the NLP literature. We use “line” to denote a sequence that ends with a newline character and “sentence” as one with an ending punctuation. Most parallel corpora, such as ours, are aligned by line, as a line may be part of a sentence or without an ending punctuation (e.g. a header/title). Using a standardized unit such as “line” would also be a fairer measure to *linguae/scriptiones continuae* (languages/scripts with no explicit punctuation).

For character modeling, we used a dummy symbol to denote each whitespace. For byte, we turned each UTF-8-encoded character into a byte string in decimal value, such that each token is a number between 0 and 255, inclusive. For word, we followed (Junczys-Dowmunt et al., 2016) and used the Moses tokenizer (Koehn et al., 2007) as is standard in NMT practice when word tokenization is applied and Jieba⁷ for segmentation in ZH.

For Pinyin, we used the implementation from <https://github.com/lxyu/pinyin> in the numerical format such that each character/syllable is followed by a single digit indicating its lexical tone in Mandarin. For Wubi, we used the dictionary from the implementation from <https://github.com/arcsecw/wubi>.

We have implemented all representations such that they would be reversible even when the sequence contains code-mixing.

We used the official dev set as provided in (Ziemski et al., 2016), 3,077 lines per language remained from 4,000 after filtering line length to 300 characters. Data statistics is provided in Appendix D for reference.

The systematic training regime that we give to our language directions are identical for all. For each primary representation type (character, byte, and word), we performed:

- 5 runs in 5 sizes ($10^2 - 10^6$): A0 (seed=13), B0 (13), C0 (9948), A1 (9948), A2 (265), and
- 7 more runs in 4 sizes ($10^2 - 10^5$): A3 (777), A4 (42), A5 (340589), A6 (1000), A7 (83146), B1 (9948), & C1 (13).

For each run and each size, there are 30 pairwise directions (i.e. 1 source language to 1 target language, e.g. AR-EN for Arabic to English) that result from the 6 languages. We trained all 150 jobs for each run and representation using the Transformer model (Vaswani et al., 2017) as supported by the SOCKEYE Toolkit (Hieber et al., 2018) (version 1.18.85), based on MXNet (Chen et al., 2015). A detailed description of the architecture of the Transformer can be found in (Vaswani et al., 2017). The same set of hyperparameters applies to all and its values are listed in Appendix C.

Notes on training time Each run of 30 directions in 5 sizes took approximately 8-12 days for character and byte models. Byte models generally took longer — hence training time is positively correlated with length (concurring with observations by Cherry et al. (2018) as they compared character with BPE models). A maximum length of 300 characters entails a maximum length of *at least* 300 bytes in UTF-8. Each run of word models (30 directions, 5 sizes) took about 6 days (excluding the training of some 7-9 directions out of 30 per run involving AR_{trg} or RU_{trg} at 10^6 on word level which took about 12-18 hours *each direction* to train on a CPU as these required more space and would run out of memory (OOM) on our GPUs otherwise). These figures do not include the additional probing experiments described in § 4.

Evaluation metric Most sequence-to-sequence models are optimized using a cross-entropy loss, defined as:

$$H(\mathbf{t}, \mathbf{s}) = - \sum_{i=1}^N \log_2 p(t_i | \mathbf{t}_{<i}, \mathbf{s}) \quad (1)$$

where \mathbf{t} is the sequence of tokens to be predicted, t_i refers to the i^{th} token in that sequence, \mathbf{s} is the sequence of tokens conditioned on, and $N = |\mathbf{t}|$. It is customary to report scores as PP, which is $2^{\frac{1}{N}H(\mathbf{t}, \mathbf{s})}$, i.e. 2 to the power of the cross-entropy averaged by the number of tokens (based on whichever granularity of unit is used for training) in the data. Cotterell et al. (2018) propose to use “renormalized” PP to evaluate LMs fairly through the division of an arbitrary constant. In our case, we choose instead a simpler method of using an “unnormalized” PP, i.e. the total number of bits needed to encode the development (dev) set, which has a constant size of 3,077 lines per language (after length filtering of the same dev set used in Junczys-Dowmunt et al. (2016)) for all various training sizes. As the implementation we used (SOCKEYE (Hieber et al., 2018)) only reports PP, we transform it back to entropy as defined above by noting that $H(\mathbf{t}, \mathbf{s}) = \log_2 PP(\mathbf{t}|\mathbf{s}) \times N$.

⁷<https://github.com/fxsjy/jieba>

C HYPERPARAMETER SETTING

- encoder transformer;
- decoder transformer;
- num-layers 6:6;
- num-embed 512:512;
- transformer-model-size 512;
- transformer-attention-heads 8;
- transformer-feed-forward-num-hidden 2048;
- transformer-activation-type relu;
- transformer-positional-embedding-type fixed;
- transformer-preprocess d; transformer-postprocess drn;
- transformer-dropout-attention 0.1;
- transformer-dropout-act 0.1;
- transformer-dropout-prepost 0.1;
- batch-size 15;
- batch-type sentence;
- max-num-checkpoint-not-improved 3;
- max-num-epochs 50;
- optimizer adam;
- optimized-metric perplexity;
- optimizer-params epsilon: 0.000000001, beta1: 0.9, beta2: 0.98;
- label-smoothing 0.0;
- learning-rate-reduce-num-not-improved 4;
- learning-rate-reduce-factor 0.001;
- loss-normalization-type valid;
- max-seq-len 300 for character, word, and BPE, 672 for all bytes, 688 for Wubi, 680 for Pinyin;
- checkpoint-frequency/interval 4000.

(For smaller datasets, the end of 50 epochs is often reached before the first checkpoint. Since SOCKEYE only outputs scores at checkpoints, we adjusted the checkpoint frequency as follows to get a score outputted by the end of 50 epochs: 1000 for 100 lines for all character & byte instances, 400 for 100 lines for word and 500 for 100 lines BPE, 3450 for 1000 lines for word & BPE. For the very few cases that this default does not suffice due to bucketing of similar length sequences, we manually set the checkpoint frequency to the last batch.)

D DATA STATISTICS

- Number of types, i.e. vocabulary size (|V|). Note that Sockeye adds for its calculation 4 additional types: <pad>, <eos>, </s>, <unk>.
- Number of tokens. This excludes the 1 EOS/BOS (end-/beginning-of-sentence) marker added by Sockeye to each line.
- Out-of-vocabulary (OOV) type rate (in %), i.e. the fraction of the types in the dev data that is not covered by the types in the training data.
- OOV token rate (in %), i.e. the fraction of tokens in the dev data that is treated as UNKNOWNS.
- Type-token-ratio (in %) i.e. the ratio between the number of types and tokens in the data. This is a rough proxy for lexical diversity in that a value of 1 would indicate that no type is ever seen twice, and a value very close to 0 would indicate that very few distinct types account for almost all of the data.
- Line length (excl. EOS/BOS marker): mean \pm standard deviation, and the 0/25/50/75/100-th percentile.

Statistics for dataset A

	100	1,000	100,000	1,000,000	100,000,000	1,000,000,000	100,000,000,000	1,000,000,000,000	100,000,000,000,000	1,000,000,000,000,000
Number of TYPES										
AR	82	107	140	176	220	260	300	340	380	420
EN	81	122	151	205	250	280	300	310	330	340
FR	81	131	167	211	260	290	310	320	330	340
RU	95	106	126	172	200	220	230	240	250	260
ZH	90	108	110	118	125	130	135	140	145	150
ZH_vocab	70	86	90	100	105	110	115	120	125	130
ZH_vocab_vocab	100	106	113	120	125	130	135	140	145	150
Number of TOKENS										
AR	5,079	128,382	1,883,817	10,625,917	602,661,280	3,210,496,880	18,287,060,160	105,687,520,000	602,661,280,000	3,210,496,880,000
EN	5,079	128,382	1,883,817	10,625,917	602,661,280	3,210,496,880	18,287,060,160	105,687,520,000	602,661,280,000	3,210,496,880,000
FR	5,079	128,382	1,883,817	10,625,917	602,661,280	3,210,496,880	18,287,060,160	105,687,520,000	602,661,280,000	3,210,496,880,000
RU	5,079	128,382	1,883,817	10,625,917	602,661,280	3,210,496,880	18,287,060,160	105,687,520,000	602,661,280,000	3,210,496,880,000
ZH	5,079	128,382	1,883,817	10,625,917	602,661,280	3,210,496,880	18,287,060,160	105,687,520,000	602,661,280,000	3,210,496,880,000
ZH_vocab	5,079	128,382	1,883,817	10,625,917	602,661,280	3,210,496,880	18,287,060,160	105,687,520,000	602,661,280,000	3,210,496,880,000
ZH_vocab_vocab	5,079	128,382	1,883,817	10,625,917	602,661,280	3,210,496,880	18,287,060,160	105,687,520,000	602,661,280,000	3,210,496,880,000
OOV (per type) (%)										
AR	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
EN	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
FR	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
RU	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
ZH	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
ZH_vocab	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
ZH_vocab_vocab	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
OOV (per token) (%)										
AR	0.75	0.35	0.19	0.10	0.05	0.03	0.02	0.01	0.01	0.01
EN	0.75	0.35	0.19	0.10	0.05	0.03	0.02	0.01	0.01	0.01
FR	0.75	0.35	0.19	0.10	0.05	0.03	0.02	0.01	0.01	0.01
RU	0.75	0.35	0.19	0.10	0.05	0.03	0.02	0.01	0.01	0.01
ZH	0.75	0.35	0.19	0.10	0.05	0.03	0.02	0.01	0.01	0.01
ZH_vocab	0.75	0.35	0.19	0.10	0.05	0.03	0.02	0.01	0.01	0.01
ZH_vocab_vocab	0.75	0.35	0.19	0.10	0.05	0.03	0.02	0.01	0.01	0.01
TYPE (%)										
AR	0.09	0.09	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10
EN	0.09	0.09	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10
FR	0.09	0.09	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10
RU	0.09	0.09	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10
ZH	0.09	0.09	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10
ZH_vocab	0.09	0.09	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10
ZH_vocab_vocab	0.09	0.09	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10
Mean line length (excl. EOS/BOS marker)										
AR	90.79 \pm 5.45	123.88 \pm 47.89	188.88 \pm 53.79	306.85 \pm 54.88	520.66 \pm 57.25	835.91 \pm 58.23	1,212.98 \pm 60.33	1,711.10 \pm 62.44	2,344.58 \pm 64.56	3,144.72 \pm 66.67
EN	90.79 \pm 5.45	123.88 \pm 47.89	188.88 \pm 53.79	306.85 \pm 54.88	520.66 \pm 57.25	835.91 \pm 58.23	1,212.98 \pm 60.33	1,711.10 \pm 62.44	2,344.58 \pm 64.56	3,144.72 \pm 66.67
FR	90.79 \pm 5.45	123.88 \pm 47.89	188.88 \pm 53.79	306.85 \pm 54.88	520.66 \pm 57.25	835.91 \pm 58.23	1,212.98 \pm 60.33	1,711.10 \pm 62.44	2,344.58 \pm 64.56	3,144.72 \pm 66.67
RU	90.79 \pm 5.45	123.88 \pm 47.89	188.88 \pm 53.79	306.85 \pm 54.88	520.66 \pm 57.25	835.91 \pm 58.23	1,212.98 \pm 60.33	1,711.10 \pm 62.44	2,344.58 \pm 64.56	3,144.72 \pm 66.67
ZH	90.79 \pm 5.45	123.88 \pm 47.89	188.88 \pm 53.79	306.85 \pm 54.88	520.66 \pm 57.25	835.91 \pm 58.23	1,212.98 \pm 60.33	1,711.10 \pm 62.44	2,344.58 \pm 64.56	3,144.72 \pm 66.67
ZH_vocab	90.79 \pm 5.45	123.88 \pm 47.89	188.88 \pm 53.79	306.85 \pm 54.88	520.66 \pm 57.25	835.91 \pm 58.23	1,212.98 \pm 60.33	1,711.10 \pm 62.44	2,344.58 \pm 64.56	3,144.72 \pm 66.67
ZH_vocab_vocab	90.79 \pm 5.45	123.88 \pm 47.89	188.88 \pm 53.79	306.85 \pm 54.88	520.66 \pm 57.25	835.91 \pm 58.23	1,212.98 \pm 60.33	1,711.10 \pm 62.44	2,344.58 \pm 64.56	3,144.72 \pm 66.67

Statistics for dataset B

Representation	100	1,000	10,000	100,000	1,000,000	10,000,000	100,000,000	1,000,000,000	10,000,000,000	100,000,000,000	1,000,000,000,000	10,000,000,000,000
Number of Files	87	124	124	124	124	124	124	124	124	124	124	124
Number of TYPES	87	124	124	124	124	124	124	124	124	124	124	124
AK	56	88	124	124	124	124	124	124	124	124	124	124
EN	50	79	131	131	131	131	131	131	131	131	131	131
FR	80	122	122	122	122	122	122	122	122	122	122	122
RU	100	142	130	130	130	130	130	130	130	130	130	130
ZH	71	124	124	124	124	124	124	124	124	124	124	124
ZH_main	71	124	124	124	124	124	124	124	124	124	124	124
ZH_work	87	124	124	124	124	124	124	124	124	124	124	124
AK_spl256	9,798	100,399	1,019,696	10,238,976	102,481,816	1,024,818,166	10,248,181,666	102,481,816,666	1,024,818,166,666	10,248,181,666,666	102,481,816,666,666	1,024,818,166,666,666
EN_spl256	13,544	135,396	1,353,960	13,539,600	135,396,000	1,353,960,000	13,539,600,000	135,396,000,000	1,353,960,000,000	13,539,600,000,000	135,396,000,000,000	1,353,960,000,000,000
FR_spl256	13,173	131,732	1,317,320	13,173,200	131,732,000	1,317,320,000	13,173,200,000	131,732,000,000	1,317,320,000,000	13,173,200,000,000	131,732,000,000,000	1,317,320,000,000,000
RU_spl256	3,829	38,290	382,900	3,829,000	38,290,000	382,900,000	3,829,000,000	38,290,000,000	382,900,000,000	3,829,000,000,000	38,290,000,000,000	382,900,000,000,000
ZH_spl256	11,801	118,010	1,180,100	11,801,000	118,010,000	1,180,100,000	11,801,000,000	118,010,000,000	1,180,100,000,000	11,801,000,000,000	118,010,000,000,000	1,180,100,000,000,000
ZH_main_spl256	10,155	101,550	1,015,500	10,155,000	101,550,000	1,015,500,000	10,155,000,000	101,550,000,000	1,015,500,000,000	10,155,000,000,000	101,550,000,000,000	1,015,500,000,000,000
ZH_work_spl256	10,155	101,550	1,015,500	10,155,000	101,550,000	1,015,500,000	10,155,000,000	101,550,000,000	1,015,500,000,000	10,155,000,000,000	101,550,000,000,000	1,015,500,000,000,000
AK_spl128	4,783	47,830	478,300	4,783,000	47,830,000	478,300,000	4,783,000,000	47,830,000,000	478,300,000,000	4,783,000,000,000	47,830,000,000,000	478,300,000,000,000
EN_spl128	6,727	67,270	672,700	6,727,000	67,270,000	672,700,000	6,727,000,000	67,270,000,000	672,700,000,000	6,727,000,000,000	67,270,000,000,000	672,700,000,000,000
FR_spl128	20,313	203,130	2,031,300	20,313,000	203,130,000	2,031,300,000	20,313,000,000	203,130,000,000	2,031,300,000,000	20,313,000,000,000	203,130,000,000,000	2,031,300,000,000,000
RU_spl128	2,671	26,710	267,100	2,671,000	26,710,000	267,100,000	2,671,000,000	26,710,000,000	267,100,000,000	2,671,000,000,000	26,710,000,000,000	267,100,000,000,000
ZH_spl128	6,184	61,840	618,400	6,184,000	61,840,000	618,400,000	6,184,000,000	61,840,000,000	618,400,000,000	6,184,000,000,000	61,840,000,000,000	618,400,000,000,000
ZH_main_spl128	5,325	53,250	532,500	5,325,000	53,250,000	532,500,000	5,325,000,000	53,250,000,000	532,500,000,000	5,325,000,000,000	53,250,000,000,000	532,500,000,000,000
ZH_work_spl128	24,313	243,130	2,431,300	24,313,000	243,130,000	2,431,300,000	24,313,000,000	243,130,000,000	2,431,300,000,000	24,313,000,000,000	243,130,000,000,000	2,431,300,000,000,000
OOV type rate (%)	0.31	0.01	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
AK	0.31	0.01	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
EN	0.05	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
FR	0.07	0.01	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
RU	0.22	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
ZH	0.03	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
ZH_main	0.09	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
ZH_work	0.09	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
AK_spl256	0.31	0.01	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
EN_spl256	0.05	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
FR_spl256	0.07	0.01	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
RU_spl256	0.22	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
ZH_spl256	0.03	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
ZH_main_spl256	0.09	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
ZH_work_spl256	0.09	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Mean line length (tokens)	97.98 ± 55.47	100.60 ± 50.00	101.97 ± 45.00	102.40 ± 58.36	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39
AK	97.98 ± 55.47	100.60 ± 50.00	101.97 ± 45.00	102.40 ± 58.36	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39
EN	97.98 ± 55.47	100.60 ± 50.00	101.97 ± 45.00	102.40 ± 58.36	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39
FR	97.98 ± 55.47	100.60 ± 50.00	101.97 ± 45.00	102.40 ± 58.36	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39
RU	97.98 ± 55.47	100.60 ± 50.00	101.97 ± 45.00	102.40 ± 58.36	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39
ZH	97.98 ± 55.47	100.60 ± 50.00	101.97 ± 45.00	102.40 ± 58.36	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39
ZH_main	97.98 ± 55.47	100.60 ± 50.00	101.97 ± 45.00	102.40 ± 58.36	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39
ZH_work	97.98 ± 55.47	100.60 ± 50.00	101.97 ± 45.00	102.40 ± 58.36	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39
AK_spl256	97.98 ± 55.47	100.60 ± 50.00	101.97 ± 45.00	102.40 ± 58.36	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39
EN_spl256	97.98 ± 55.47	100.60 ± 50.00	101.97 ± 45.00	102.40 ± 58.36	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39
FR_spl256	97.98 ± 55.47	100.60 ± 50.00	101.97 ± 45.00	102.40 ± 58.36	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39
RU_spl256	97.98 ± 55.47	100.60 ± 50.00	101.97 ± 45.00	102.40 ± 58.36	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39
ZH_spl256	97.98 ± 55.47	100.60 ± 50.00	101.97 ± 45.00	102.40 ± 58.36	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39
ZH_main_spl256	97.98 ± 55.47	100.60 ± 50.00	101.97 ± 45.00	102.40 ± 58.36	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39
ZH_work_spl256	97.98 ± 55.47	100.60 ± 50.00	101.97 ± 45.00	102.40 ± 58.36	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39	102.48 ± 58.39

Statistics for development (dev) set

As a different set of vocabulary is learned from each training dataset and data size, BPE has a distinct dev set for each.

Representation	CHAR		BYTE		WORD		BPE_A		BPE_B		BPE_C		BPE_C	
	100	1,000	100,000	1,000,000	100	1,000	100,000	1,000,000	100	1,000	100,000	1,000,000	100	1,000
Number of lines in train set														
Number of TYPES														
AR	130	123	13,859	708	8,232	8,994	11,473	12,430	805	3,991	10,755	13,000	908	4,119
EN	67	102	7,209	660	2,927	5,088	7,461	7,536	801	3,189	6,538	7,557	823	3,335
FR	113	118	8,332	714	2,821	6,055	8,596	8,666	827	3,407	7,859	8,665	883	3,863
RU	150	150	13,609	824	3,853	6,905	10,000	10,000	875	4,000	10,000	10,000	900	4,000
ZH	1,959	133	7,433	3,224	4,215	6,288	7,541	7,634	3,201	4,483	6,235	7,718	1,300	4,521
ZH_pinyin														
ZH_vocab	8													
ZH_pinyin_vocab	130													
ZH_vocab_vocab	130													
Number of TOKENS														
AR	331,328	603,316	61,371	167,574	115,603	83,001	79,284	79,327	110,689	97,843	73,311	68,328	110,623	97,231
EN	127,141	177,100	106,614	170,433	113,983	88,634	82,281	81,341	155,657	103,748	84,014	80,599	154,746	103,172
FR	138	118	8,332	714	2,821	6,055	8,596	8,666	827	3,407	7,859	8,665	883	3,863
RU	433,333	433,333	133,628	177,839	113,628	79,726	72,840	71,081	103,319	101,294	75,103	70,000	103,306	99,075
ZH	307,900	311,955	60,033	96,745	80,231	68,129	63,614	62,830	93,775	75,036	65,232	61,882	94,127	75,718
ZH_pinyin	258,979													
ZH_vocab	438,154													
AR_spl256	334,338													
RU_spl251	431,338													
TR (%)														
EN	0.01	0.02	2.79	0.42	1.04	1.04	15.04	17.62	0.58	4.08	14.61	18.90	0.51	4.21
EN	0.02	0.03	10.64	0.42	2.52	7.78	10.47	10.74	0.57	3.47	8.86	10.83	0.59	3.50
ES	0.03	0.03	10.65	0.45	2.46	7.80	10.54	10.93	0.59	3.50	8.86	10.83	0.59	3.50
RU	0.03	0.02	19.97	0.49	3.32	12.64	17.25	17.90	0.63	4.42	14.73	18.35	0.63	4.51
ZH	0.03	0.01	12.53	0.53	3.23	9.37	11.55	12.19	0.48	3.52	12.46	12.46	0.41	3.57
ZH_pinyin														
ZH_vocab														
AR_spl256	0.01													
RU_spl251	0.01													
Mean length of tokens														
AR	159.65 ± 59.61	159.33 ± 105.35	97.61 ± 69.51	157.61 ± 113.57	157.61 ± 113.57	157.61 ± 113.57	157.61 ± 113.57	157.61 ± 113.57	157.61 ± 113.57	157.61 ± 113.57	157.61 ± 113.57	157.61 ± 113.57	157.61 ± 113.57	157.61 ± 113.57
EN	300.110115277	61.071092777	203.1222777	202.777	202.777	202.777	202.777	202.777	202.777	202.777	202.777	202.777	202.777	202.777
EN	127.14 ± 68.64	127.14 ± 68.65	21.86 ± 11.81	50.64 ± 27.44	33.68 ± 18.67	25.05 ± 13.56	23.16 ± 12.42	22.86 ± 12.27	45.58 ± 24.86	29.53 ± 16.39	23.93 ± 12.95	22.54 ± 12.06	45.62 ± 24.67	29.57 ± 16.13
ES	4.771143267	3.781143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267
FR	142.37 ± 71.81	147.08 ± 60.39	25.59 ± 13.86	51.04 ± 26.55	37.27 ± 20.65	28.84 ± 15.76	26.79 ± 14.47	26.51 ± 14.35	50.78 ± 27.79	33.82 ± 18.74	27.66 ± 15.11	26.27 ± 14.30	49.97 ± 27.25	33.83 ± 18.80
RU	4.71143267	4.71143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267	1.1321143267
ZH	35.10 ± 13.48	35.10 ± 13.48	97.85 ± 52.10	31.44 ± 16.96	26.67 ± 14.19	21.14 ± 12.02	20.67 ± 11.16	20.41 ± 11.05	30.48 ± 16.47	24.58 ± 13.42	21.19 ± 11.40	20.11 ± 10.79	30.59 ± 16.45	24.61 ± 13.49
ZH_pinyin	6.671231573													
ZH_vocab	107.49 ± 56.73	108.66 ± 18.01												
AR_spl256	4.0010811294													
RU_spl251	4.0010811294													

F CORRELATION STATISTICS

Best correlating metrics, i.e. the union of top 3 metrics for all representations.

For each representation, the **top 3 metrics** are boldfaced.

All correlations are **highly significant** ($p < 10^{-30}$), except for min source length for WORD ($p \approx 0.0001$) and min target length for WORD ($p \approx 0.3861$).

Metric	CHAR	Pinyin	Wubi	BYTE	ARRU _t	ARRU _{s,t}	WORD	BPE
minimum length (target)	0.84	0.85	0.86	0.60	0.84	0.84	-0.02	0.65
minimum length (source)	0.82	0.84	0.85	0.57	0.84	0.84	0.10	0.64
number of tokens (source)	-0.78	-0.81	-0.82	-0.60	-0.81	-0.81	-0.59	-0.83
TTR (target)	0.83	0.83	0.84	0.48	0.81	0.81	0.61	0.83
$ V $ (source)	-0.54	-0.51	-0.51	-0.50	-0.67	-0.68	-0.63	-0.86
data size in lines	-0.80	-0.83	-0.83	-0.59	-0.81	-0.81	-0.62	-0.86
OOV token rate (target)	0.69	0.66	0.66	0.47	0.67	0.68	0.66	0.62
OOV type rate (target)	0.70	0.71	0.72	0.47	0.69	0.70	0.65	0.62
TTR (source)	0.67	0.71	0.71	0.60	0.81	0.81	0.56	0.82

The full list of metrics used for the correlation analysis is:

1. minimum length (source),
2. minimum length (target),
3. maximum length (source),
4. maximum length (target),
5. median length (source),
6. median length (target),
7. mean length (source),
8. mean length (target),
9. length std (source),
10. length std (target),
11. data size in lines,
12. number of parameters,
13. number of types ($|V|$) (source),
14. number of types ($|V|$) (target),
15. number of tokens (source),
16. number of tokens (target),
17. type-token-ratio (TTR) (source),
18. type-token-ratio (TTR) (target),
19. OOV type rate (source),
20. OOV type rate (target),
21. OOV token rate (source),
22. OOV token rate (target),
23. token ratio,
24. target type-to-parameter ratio,
25. target token-to-parameter ratio,
26. distance between the TTRs of source and target = $(1 - TTR_{src}/TTR_{trg})^2$,
27. token-to-parameter ratio (i) = $(\text{median length source} * \text{median length target} * \text{num_lines}) / \text{num_parameters}$,
28. token-to-parameter ratio (ii) = $(\text{num_source_tokens} * \text{num_target_tokens}) / \text{num_parameters}$.

G ENLARGED FIGURES FOR ALL 30 LANGUAGE DIRECTIONS (AGGREGATE RESULTS FROM ALL RUNS)

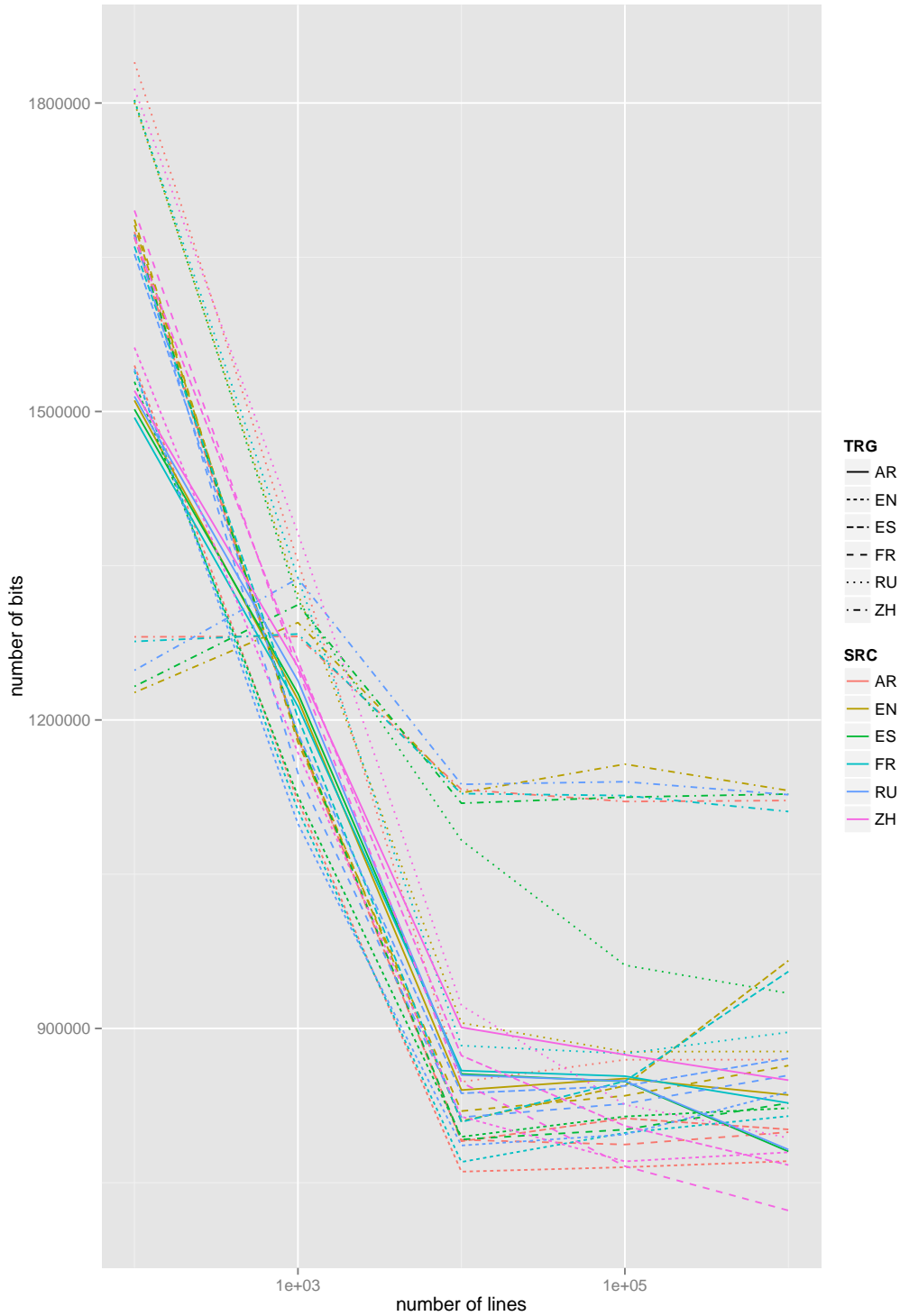


Figure 5: CHAR: character models

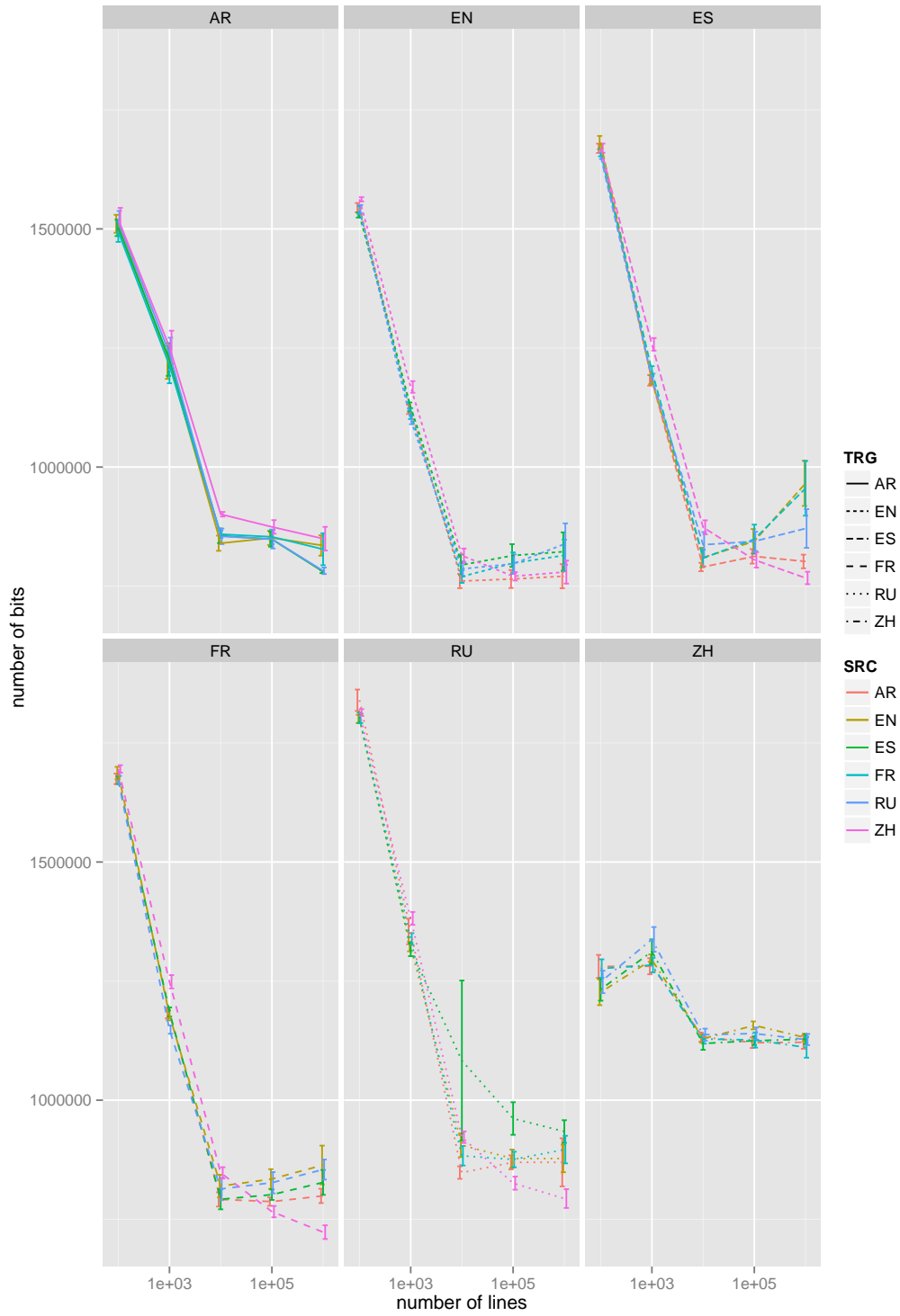


Figure 5: CHAR: character models (target language as facet)

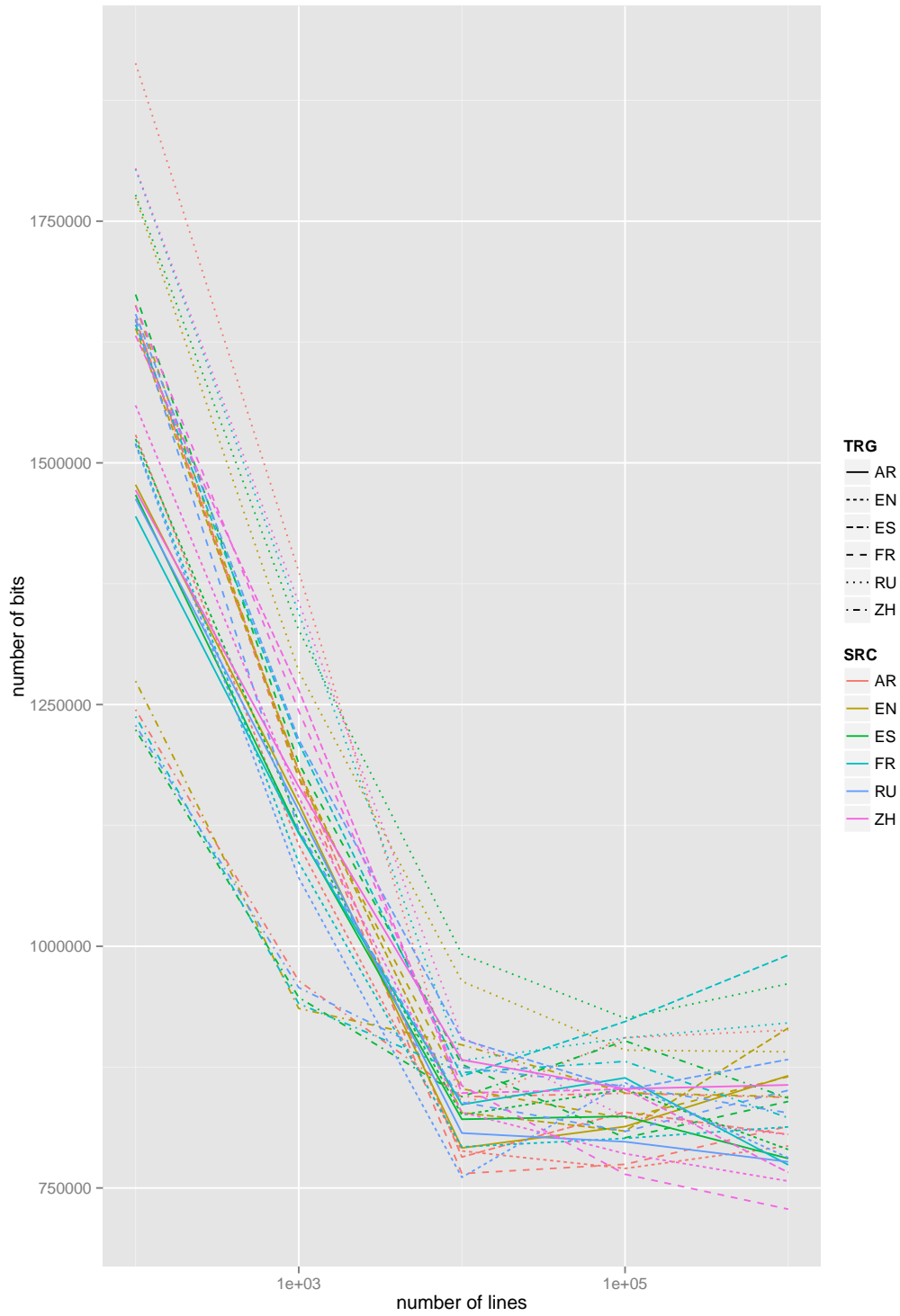


Figure 6: CHAR with Pinyin for ZH_{trg}

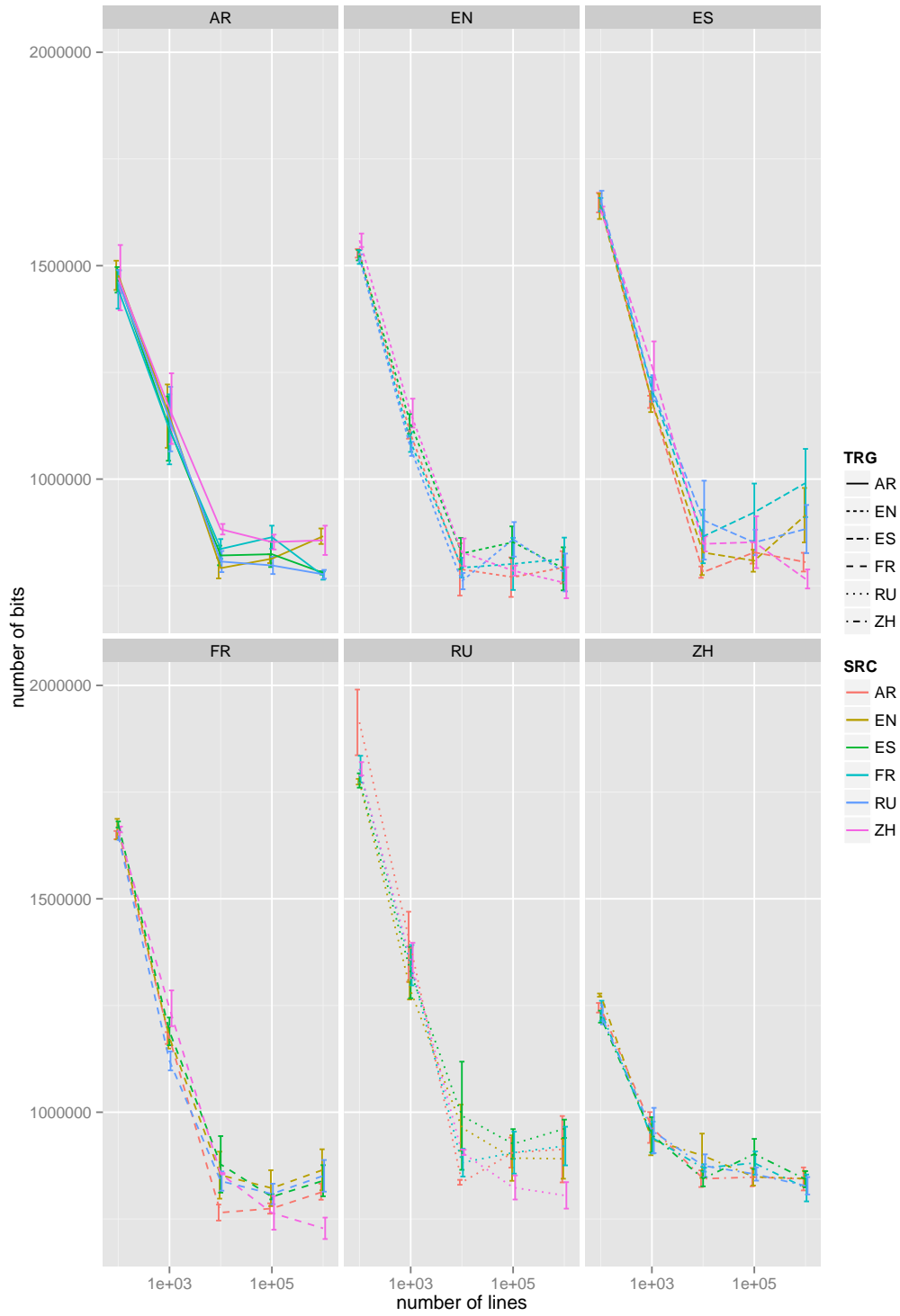


Figure 6: CHAR with Pinyin for ZH_{trg} (target language as facet)

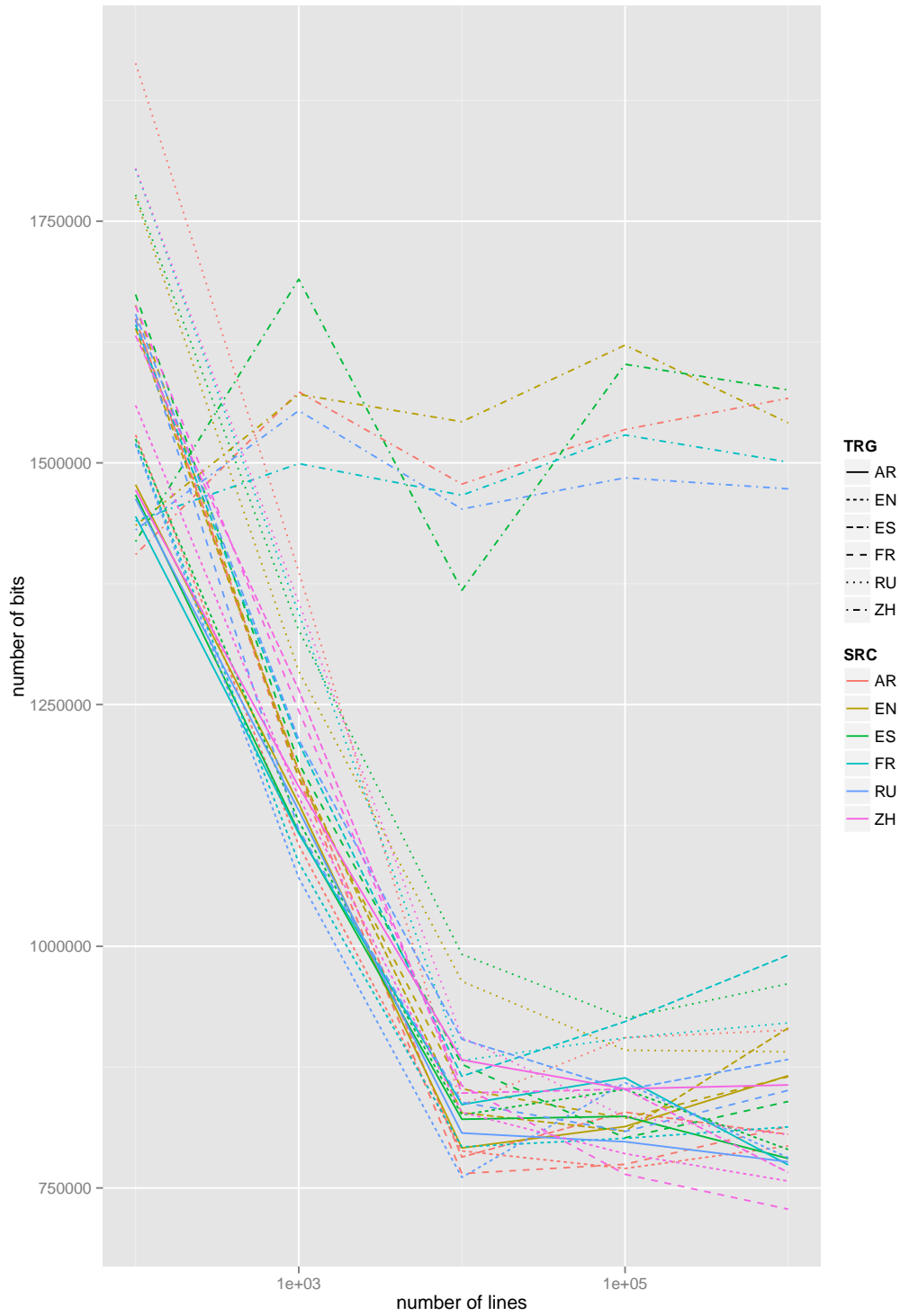


Figure 7: CHAR with Wubi for ZH_{trg}

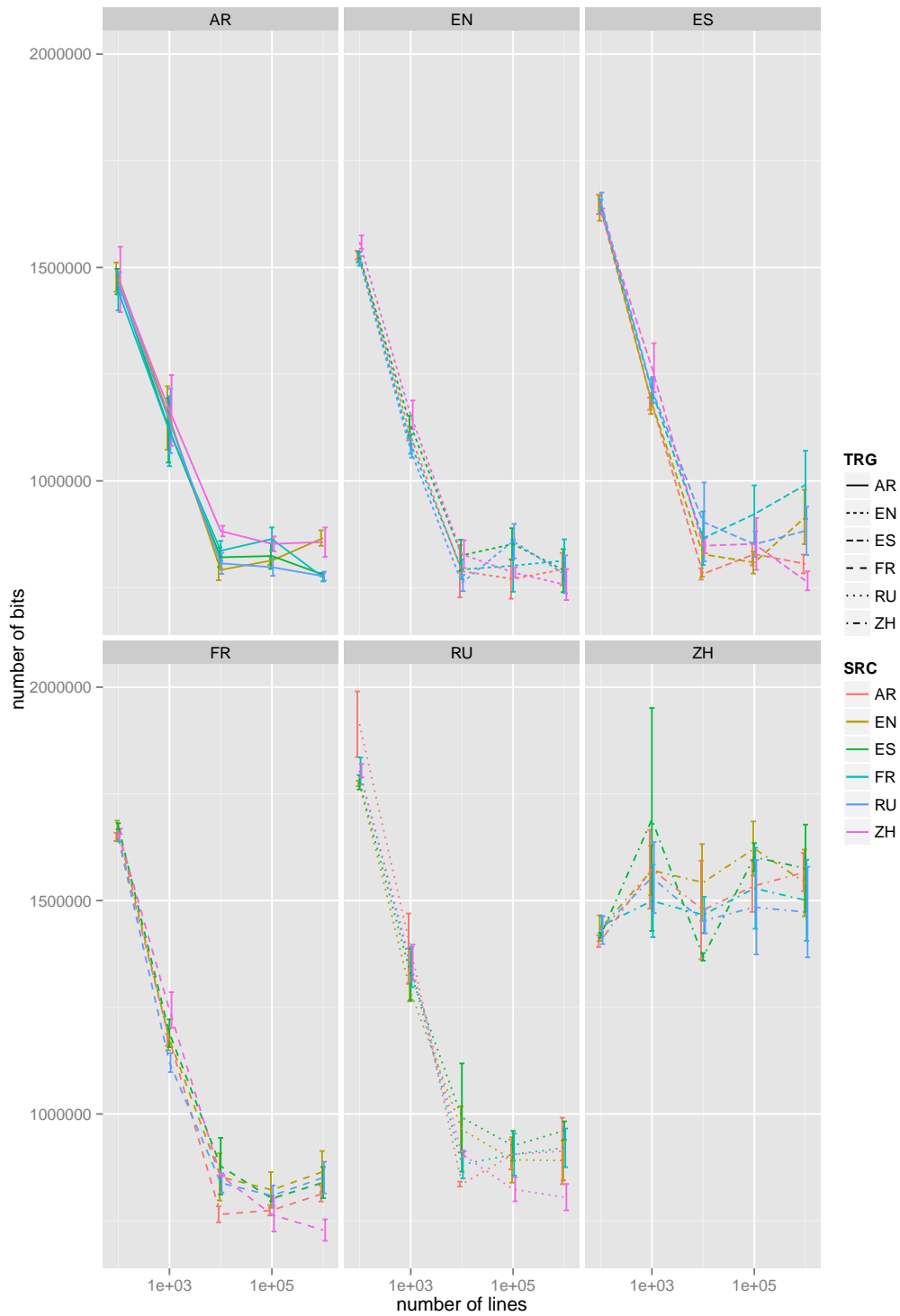


Figure 7: CHAR with Wubi for ZH_{trg} (target language as facet)

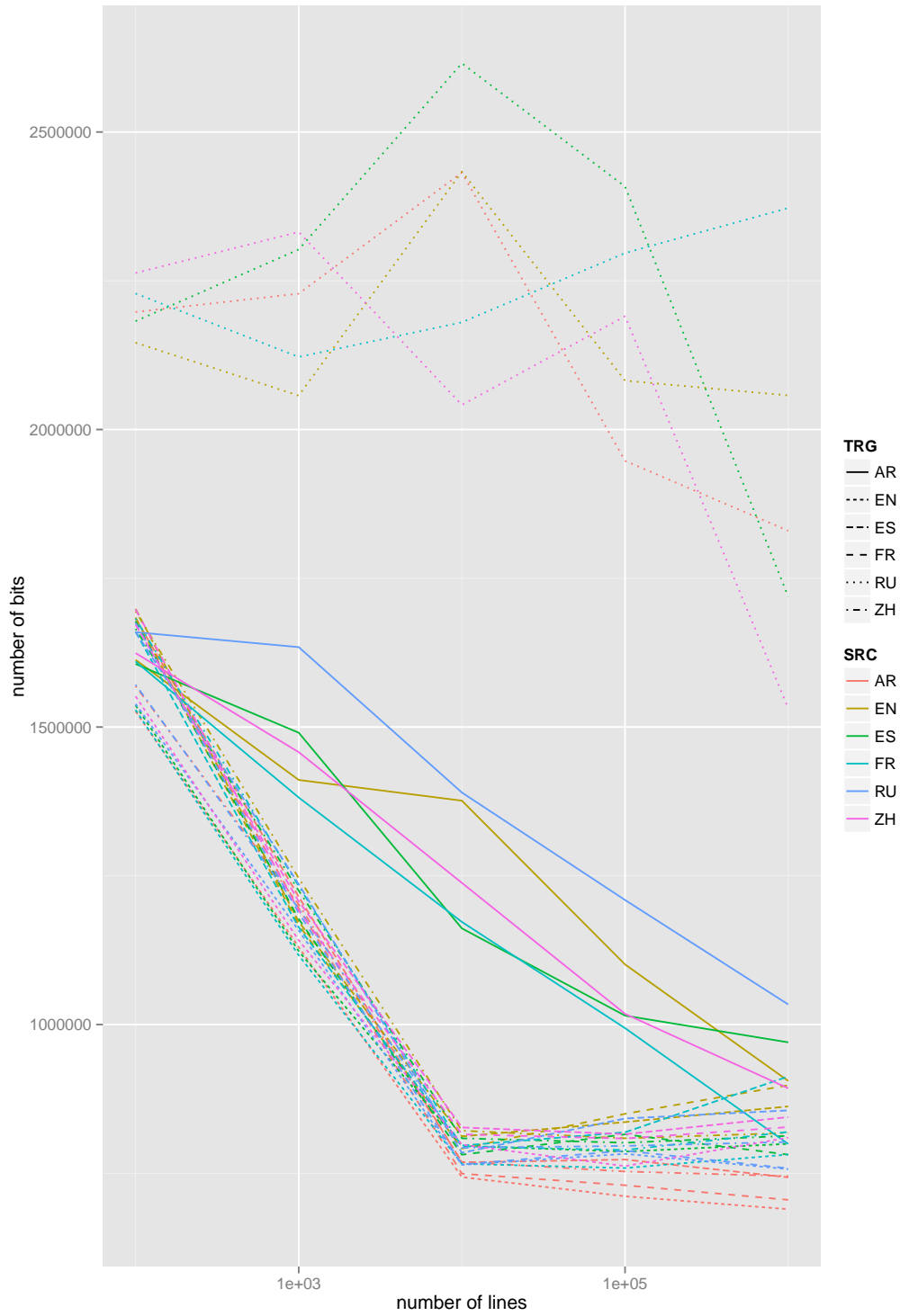


Figure 8: BYTE models with UTF-8 encoding

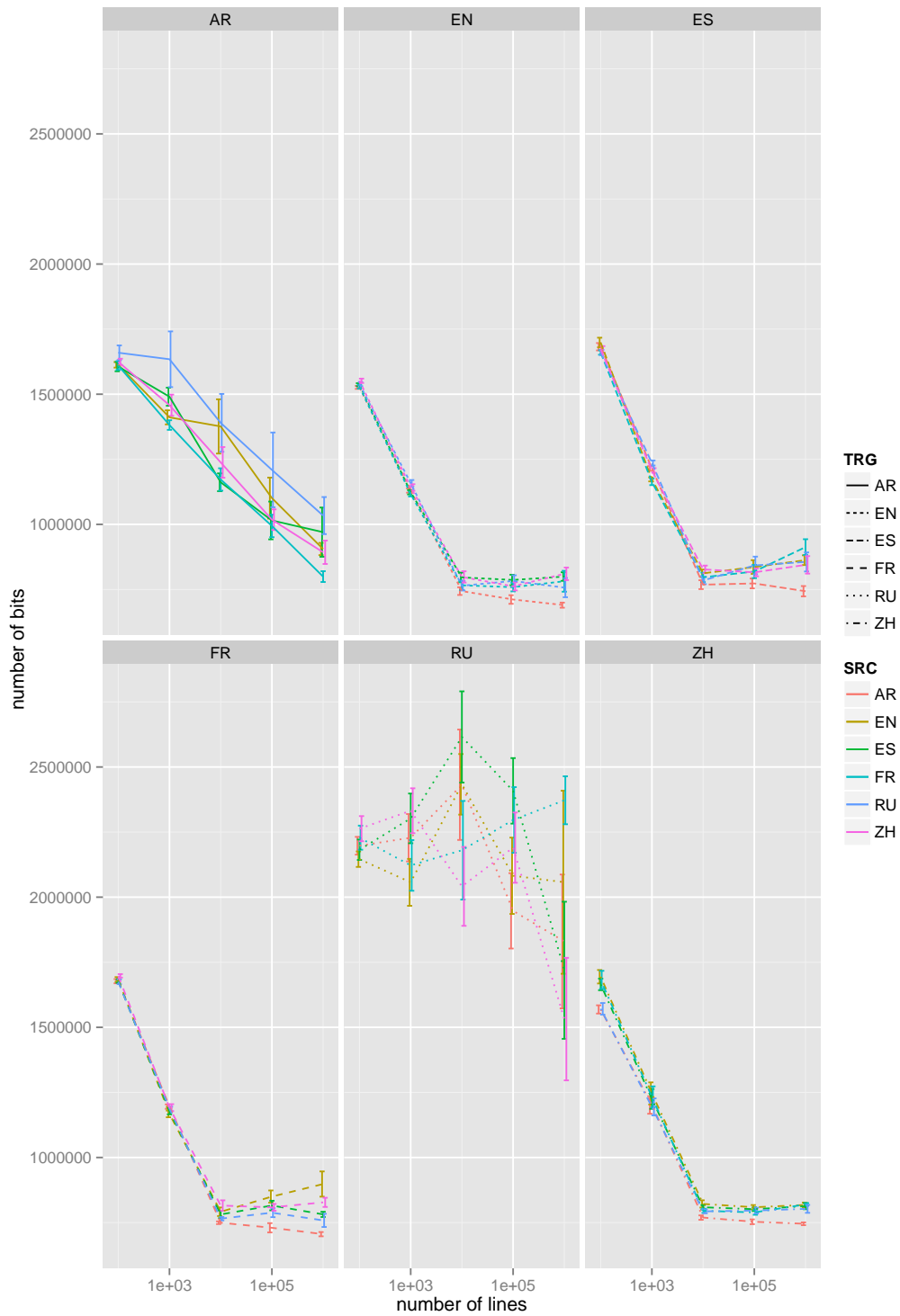


Figure 8: BYTE models with UTF-8 encoding (target language as facet)

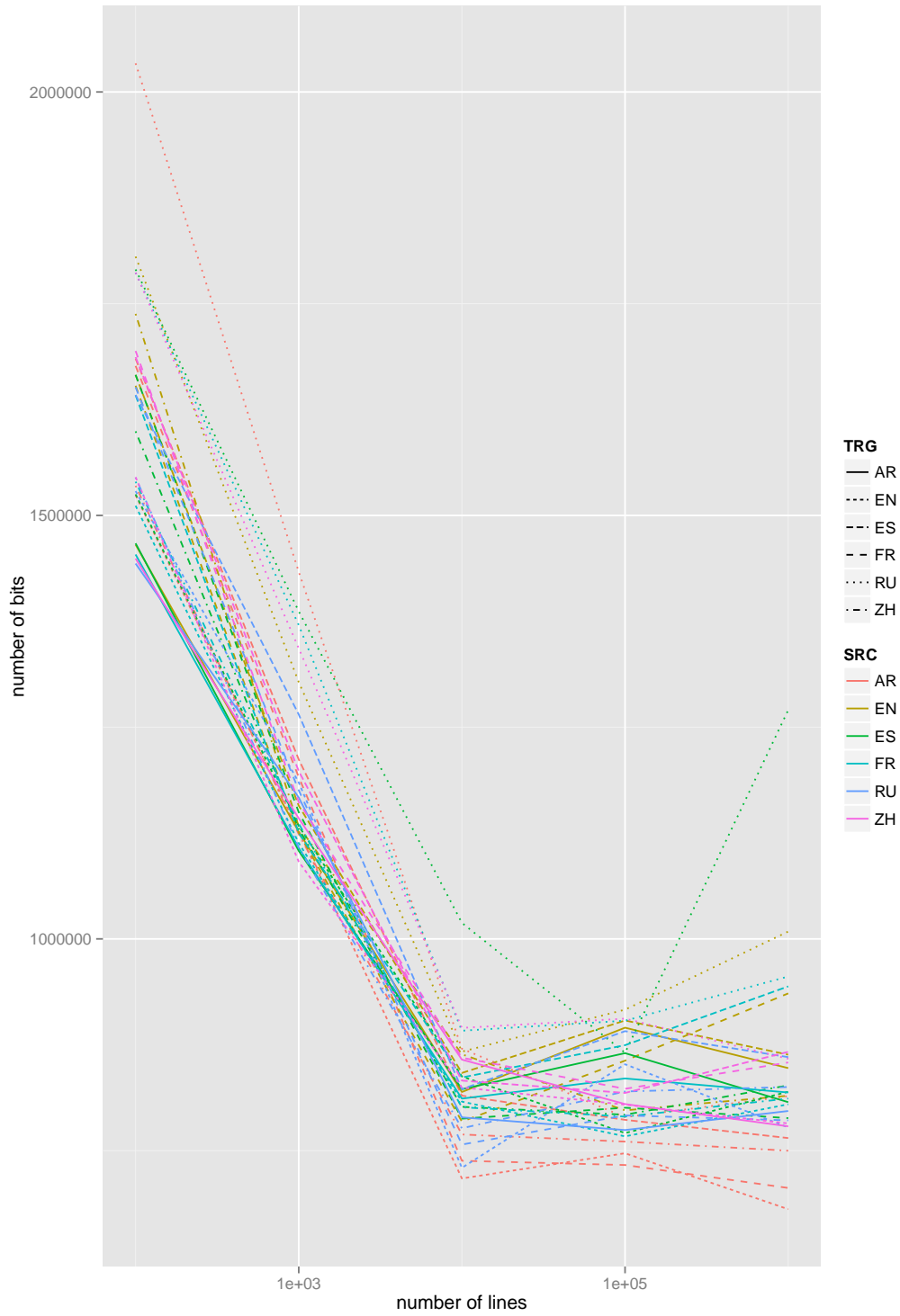


Figure 9: BYTE with AR_{trg} & RU_{trg} optimized with code pages 1256 & 1251 ($ARRU_{trg}$)

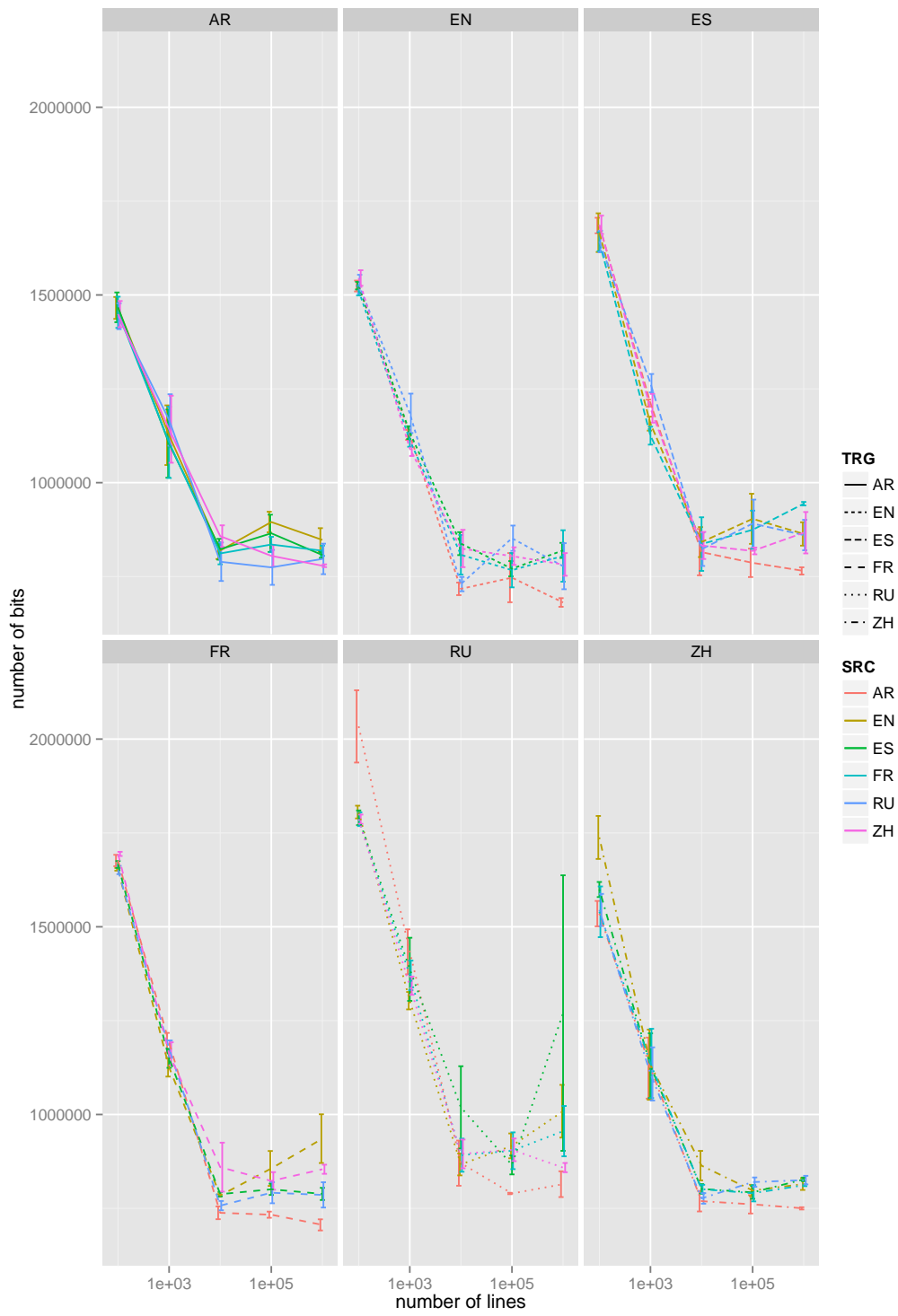


Figure 9: BYTE with AR_{trg} & RU_{trg} optimized with code pages 1256 & 1251 (target language as facet)

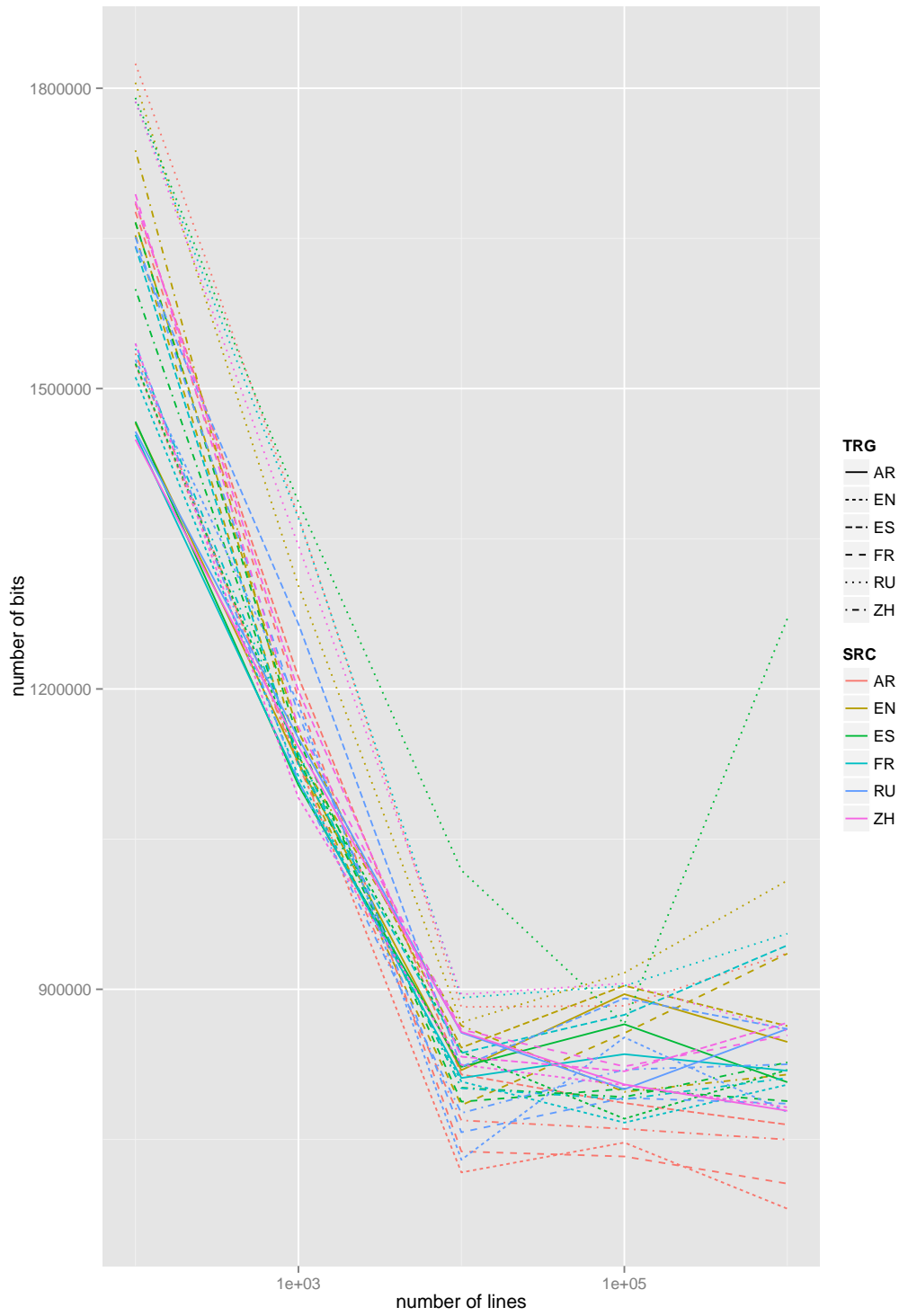


Figure 10: BYTE with directions AR-RU & RU-AR optimized on both source and target sides ($ARRU_{src,trg}$)

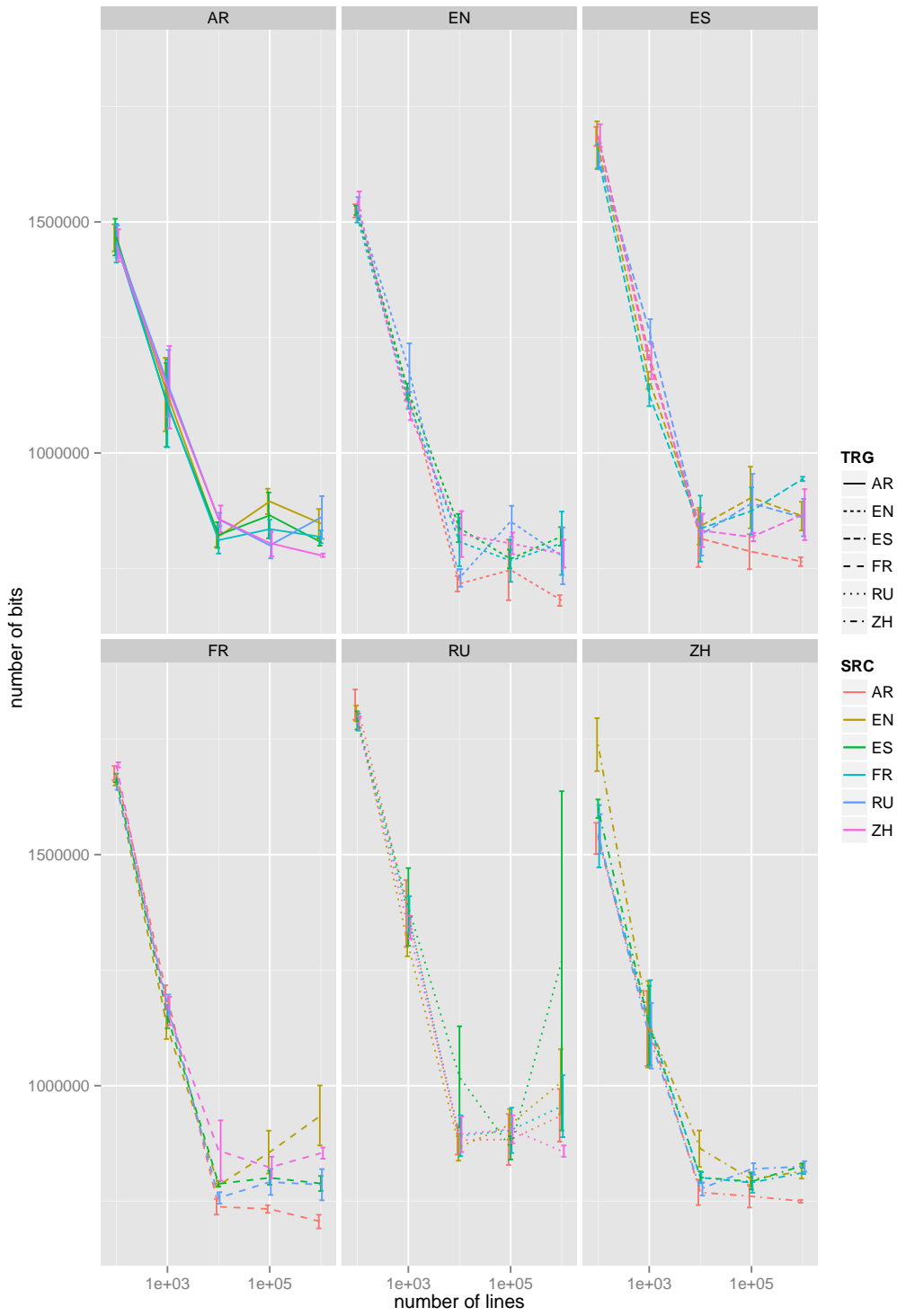


Figure 10: BYTE with directions AR-RU & RU-AR optimized on both source and target sides (target language as facet)

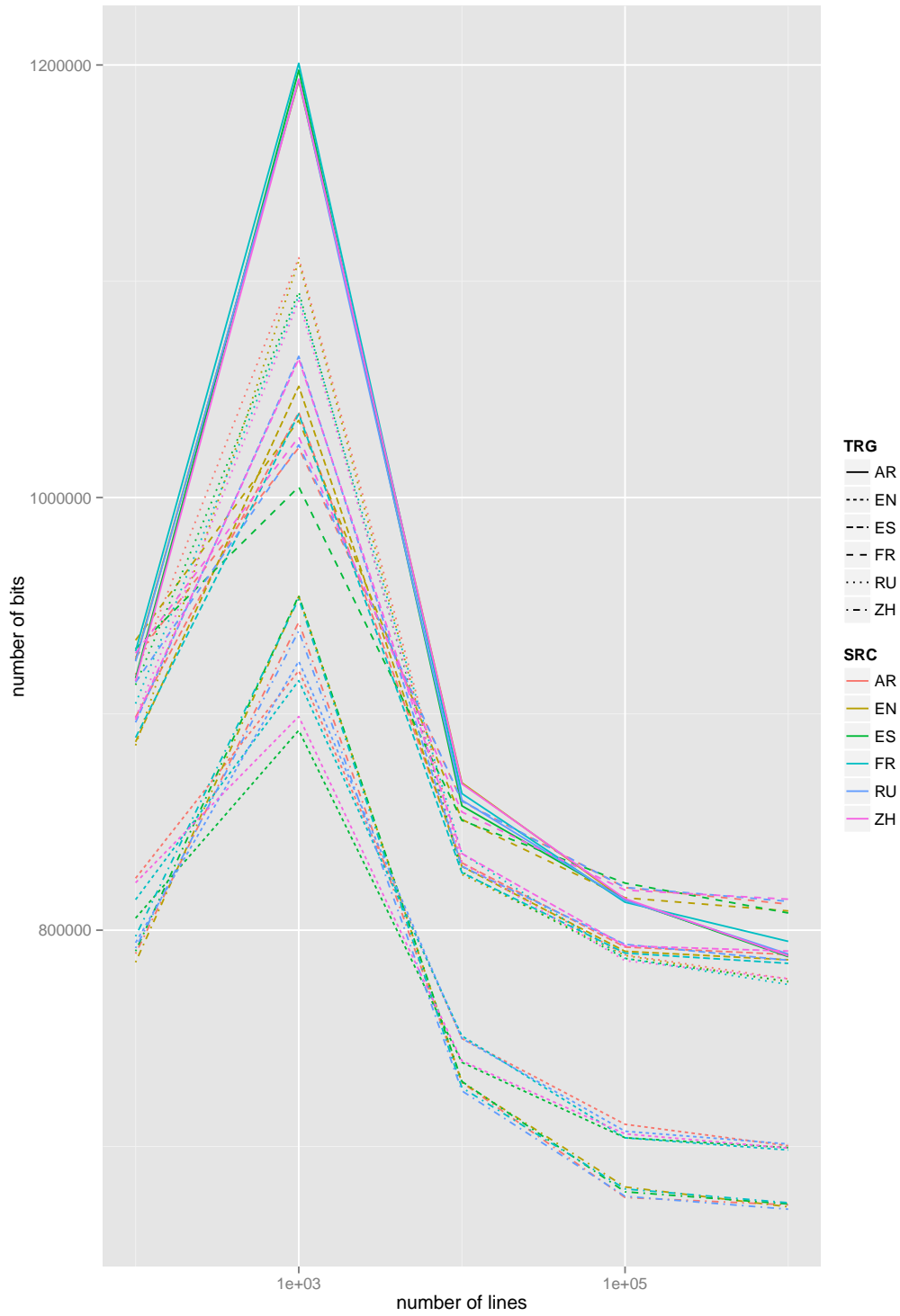


Figure 11: WORD models

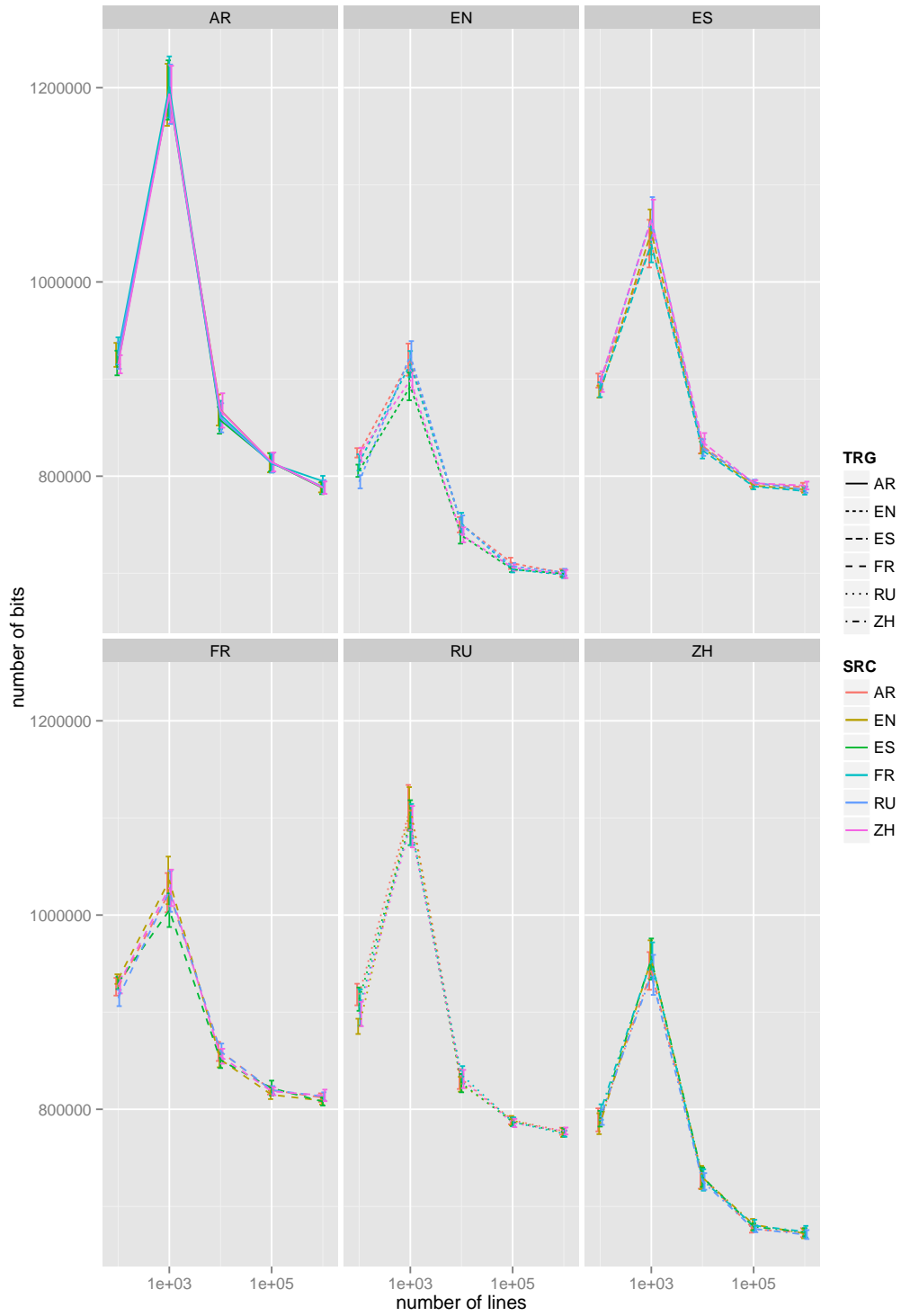


Figure 11: WORD models (target language as facet)

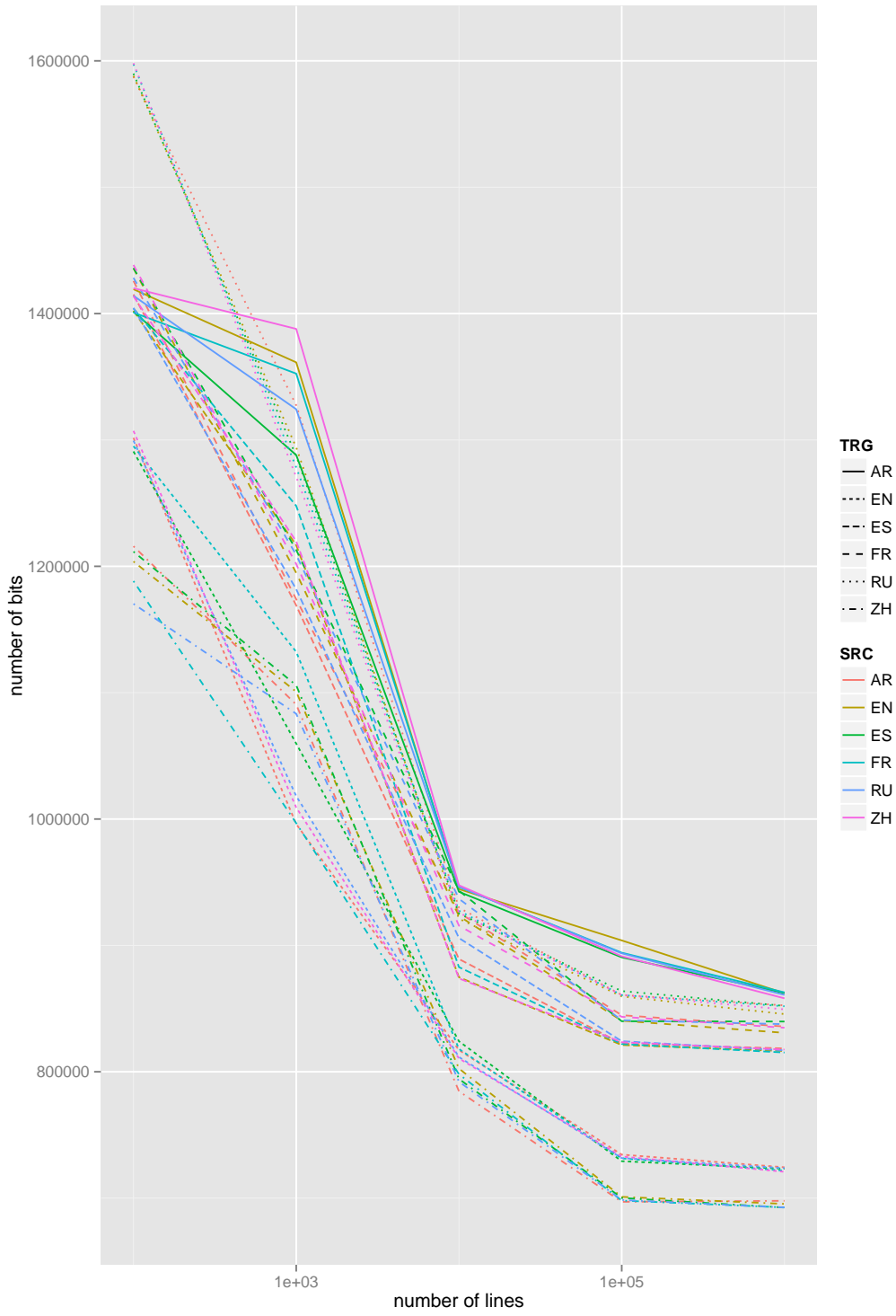


Figure 12: BPE models

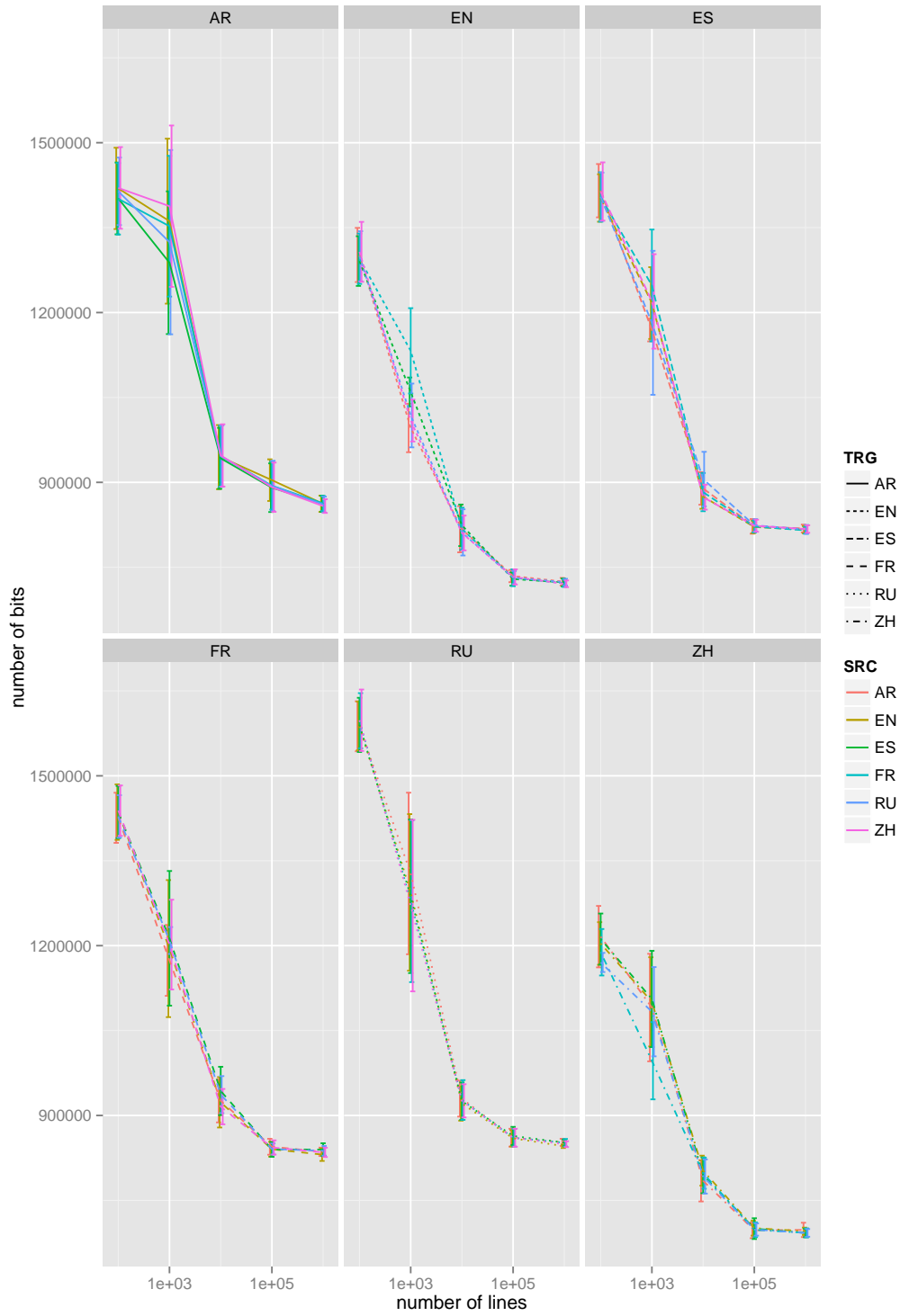


Figure 12: BPE models (target language as facet)

H SAMPLE FIGURES FROM RUN A0, ALSO SORTED BY SOURCE LANGUAGE FOR CONTRAST

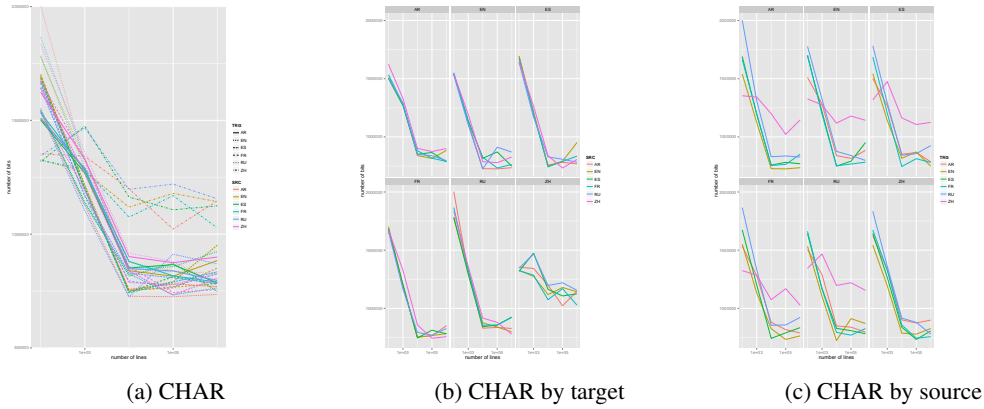


Figure 13: CHAR: character models from run A0

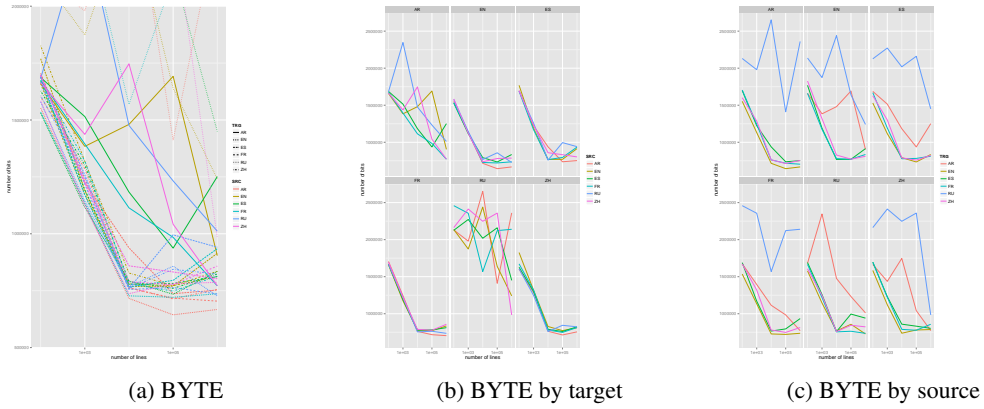


Figure 14: BYTE: byte models from run A0

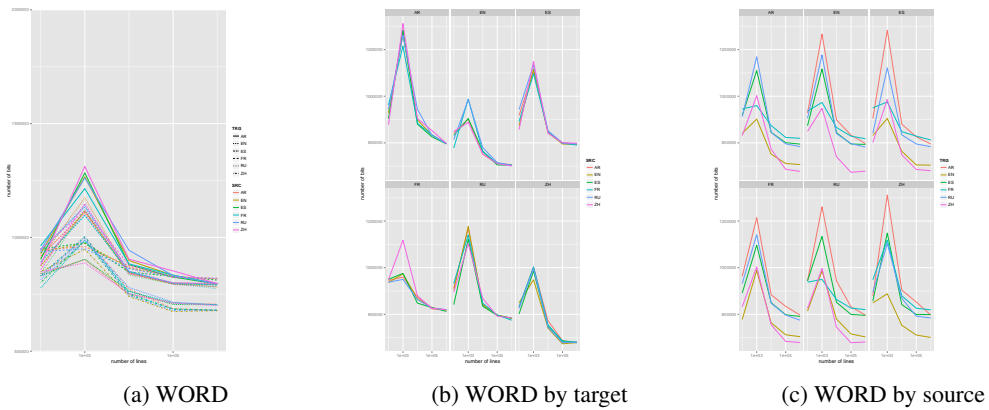


Figure 15: WORD: word models from run A0

I LANGUAGE PAIRS WITH SIGNIFICANT DIFFERENCES

15 (non-directional) language pairs total possible from 30 language directions, $p=0.001$.

LANG PAIR	CHAR	Pinyin	Wubi	BYTE	ARRU _t	ARRU _{s,t}	WORD	BPE
AR-EN				X			X	X
AR-ES								
EN-ES							X	
AR-FR				X				
EN-FR							X	X
ES-FR								
AR-RU				X				
EN-RU				X		X	X	X
ES-RU				X				
FR-RU				X				
AR-ZH	X		X	X			X	X
EN-ZH	X		X					
ES-ZH			X				X	X
FR-ZH	X		X				X	X
RU-ZH			X	X		X	X	X

Language pairs with significant differences indicate that the 2 languages are *not* equally/similarly good or equally/similarly bad.

- Character models with ZH behave differently but the disparity can be eliminated with Pinyin.
- Byte models with AR and RU exhibit unstable performance due to length but this can be rectified with compression on the target side only (ARRU_t).
- Word-based models, including BPE, however, consistently favor EN and ZH (though it is more of a “mis-segmentation” for the latter, see § 3 and Appendix J) and disfavor AR and RU (as morphologically complex languages with higher OOV rates).

J LANGUAGE COMPLEXITY

In the words of Bentz et al. (2016):

Languages are often compared with regard to their complexity from a computational, theoretical and learning perspective. In computational linguistics, it is generally known that methods mainly developed for the English language do not necessarily transfer well to other languages. The cross-linguistic variation in the amount of information encoded at the level of a word is, for instance, recognized as one of the main challenges for multilingual syntactic parsing (formulated as The Architectural Challenge (Tsarfaty et al., 2013)). Complexity of this kind is also found to influence machine translation: translating from morphologically rich languages into English is easier than the other way around (Koehn, 2005).

Morphology is “the study of the formation and internal structure of words”. Morphemes are “the smallest meaningful units of language”. (Bender, 2013)

AR and RU are traditionally considered morphologically complex (see e.g. Minkov et al. (2007), Seddah et al. (2010) and proceedings of related workshops in subsequent years), and ZH lacking morphological richness (Koehn, 2005). But this definition of morphology is predicated on the notion of word, defined primarily from an alphabetic perspective. As pointed out by Zhang & Komachi (2018), “the important differences between logographic and alphabetic writing systems have long been overlooked”. In logographic languages (i.e. languages with logographic scripts), there can be units within a character that carry semantic and phonetic information that have never been accounted for in the traditional practice of morphology or in the computation of morphological complexity. For example, in the comparison of different morphological complexity measures by Bentz et al. (2016), all measures studied are defined with the notion of word.⁸ Yet, there is **no universally valid definition of a “word”** — the form/idea (as in, the philosophical concept) of a “word” may be there for most languages/cultures (though that is certainly also debatable), but its instantiations are different in different languages/cultures, as well as in different genres/settings within one language. The variability in the definition of word is evident in the variation in language-specific word tokenization algorithms, along with the “indeterminacy of word segmentation” or a work-in-progress status for the definition of “word” advocated by Haspelmath (2011), as well as the contested nature of wordhood, esp. for logographic languages such as ZH (see Duanmu (2017) and Li et al. (2019b) for how some ZH speakers do indeed consider a ZH character to be a word or how “word”, as conventionally used in NLP, is not a native term or does not correspond with speakers’ judgement).

Our results with the Transformer indicate that a notion of morphological complexity can be modeled given our word tokenization scheme, confirming that morphological complexity is only predicated on the notion of word and bounded within the word level, and orthogonal to the performance of character or byte models. That is, unless word-based segmentation has been applied, there is no reason to attribute crosslinguistic performance disparity to differences in morphological complexity. In fact, on the character and byte level, we were able to achieve performance without disparity. Hence **disparity is not a necessary condition but an expectation that has been in mutual reinforcement with our practice of word segmentation, while the definitions of “morphological complexity” and “word” are in a circular dependency with each other.**

In this paper, we *resolve* language complexity, more specifically that of morphological complexity, in the context of computing through CLMing with the Transformer, in that we *explain away* the representation granularities and criteria relevant for such calculation.

TLDR: Up to the point of our taking up the subject of language complexity in this paper, there has been not a rigorous definition of “language complexity”. Conventionally, “language complexity” is synonymous to “linguistic complexity” (with the tradition of “linguistics” being primarily word-

⁸An exception could be that of the type/token ratio (TTR). One could imagine applying TTR on the character level for ZH, and that would be indicative of its morphological richness on the character level. However, that has thus far never been practiced or recognized in NLP.

based), and people just assume linguistic complexity, e.g. morphological/syntactic complexity, to be intrinsic and necessary in languages (across representation levels). Our findings show that linguistic complexity is relative to the representation granularity, i.e. since morphology is based on words, it is bounded to the word level.

An alternative perspective, with finer prints:

We have also developed a more rigorous interpretation. We take on the definition of “language complexity” as one that is related to the statistical attributes of languages. We assume and define *solving* as the elimination of statistically significant performance disparity.

In larger (6-layer) models, and according to the conventional definition of “language” — i.e. language as a whole, we solved language complexity with compression of AR and RU in byte representations. In smaller (1-layer) models, one can think of the situation as: i) no complexity has been modeled by the Transformer hence there is nothing to solve, or ii) there is no complexity between these languages to begin with, or iii) the Transformer solved the complexity.

With respect to each representation level/granularity in the larger models:

- **BYTE:** one can think of us as having solved complexity with byte representations or with 1-layer models — for these 6 languages empirically. Theoretically, there could be languages with longer sequence lengths than RU and AR, in those cases, we don’t claim to have solved the matter empirically but only resolved it conceptually. But this is the most that anyone could do at the moment, as there is no relevant parallel data available.
- **CHARACTER:** one can think of us as having solved it via bytes or 1-layer models. Whether we can be considered to have solved it via Pinyin for ZH depends on whether the evaluator accepts decomposition into a *phonetic representation only* qualifies as a solution for the ZH language.
- **WORD:** one can think of us as having solved it via bytes or 1-layer models. It is not possible to solve it strictly within the word level without creating word segmentation criteria that would be unrelatable to native speakers. And since “word” is exclusively a human concept, we must either claim that a universal solution is undefined or undefinable for computing, or retreat to a unit that is the greatest common factor crosslinguistically. Since some ZH speakers consider ZH characters as words, we return to the character-level solution.

It is beyond the scope of our paper to solve the qualitative disparity on the word level. However, we do advocate a more inclusive evaluation and critical reflection on the possibility of discontinuing the usage of “word” as such a non-technical term biases against both “morphologically complex” and “morphologically simple” languages. The world of languages in written form can be divided into those with logographic scripts and those with (phonetic) alphabetic ones, with the unit of character being the greatest common factor of them all, from the human perspective. For technical processing, esp. for fair multilingual sequence-to-sequence modeling with the Transformer, we recommend measures that are more standardized, such as those based on bytes or characters. There is room for improvement in the design of character encoding that complements the statistical profiles, e.g. with relative rank in sequence length, of different languages. We believe there is crosslinguistic systematicity on the character level to be leveraged.

One’s readiness to accept this as a solution to language complexity can be a subjective matter. One may insist that language complexity be solved exclusively with monolingual LMing (which lies outside the scope of the present work), instead of being confounded with the logic of one language being conditional on another. One may also object to the idea of (re-)solving morphological complexity being equivalent to or leading to solving language complexity as a whole, for there could also be e.g. syntactic complexity (although as substantial “information concerning syntactic units and relations is expressed at word level” in morphologically rich languages (Tsarfaty et al., 2010), the boundary between morphology and syntax is less distinct for some languages than others (Haspelmath, 2011)). If, however, our results could be extended, we wonder if syntactic complexity could be due to our sentence segmentation or a combination of word and sentence segmentation practice. That we leave for future work for those who are interested in the topic.

K SAMPLE-WISE DOUBLE DESCENT (DD)

K.1 OUR EXPERIMENTAL FRAMEWORK ON DD DATASETS FROM (NAKKIRAN ET AL., 2020)

Text experiments from previous work reporting sample-wise DD involved words (Belkin et al., 2019) and BPEs (Nakkiran et al., 2020).

We applied our experimental framework — by testing data points with 10^n lines — on the datasets reported in (Nakkiran et al., 2020) to exhibit DD. WMT’14⁹ EN-FR was reported to demonstrate model-wise DD and IWSLT’14 (Cettolo et al., 2012) DE-EN model-wise and sample-wise DD. We downloaded and prepared the data with scripts¹⁰ from the FAIRSEQ Toolkit (Ott et al., 2019). The WMT data was preprocessed with 40,000 BPE operations and IWSLT 10,000. Our focus is on sample-wise DD and hence our goal was to see if the spike at 10^3 we observed with the UN data would apply also to these datasets. We used the same training regime¹¹ with the Transformer and Adam on SOCKEYE as before and tested both language directions on the entirety of both datasets, with no subsampling. For the IWSLT dataset, we tested data sizes with $10^2 - 10^5$ lines, then at 160, 239 as that is the total number of lines available. For the WMT dataset, we tested from 10^2 to 10^7 , then at 35, 762, 532.

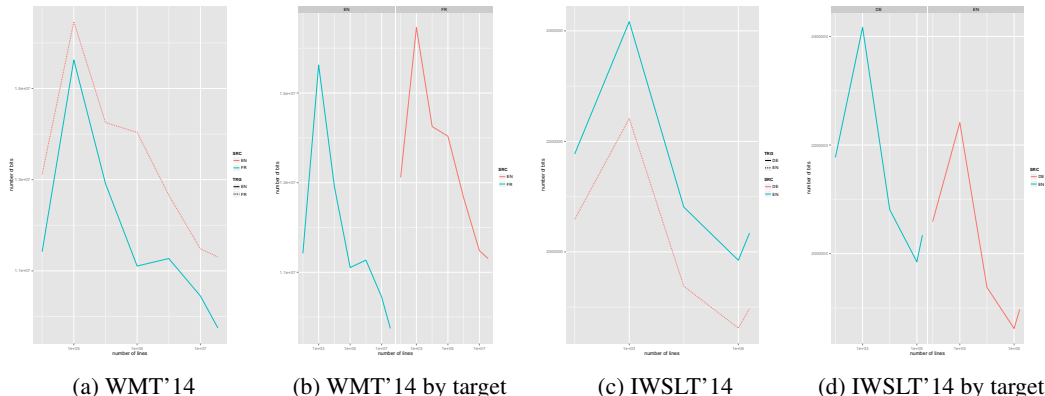


Figure 16: WMT’14 EN-FR and FR-EN and IWSLT’14 DE-EN and EN-DE: sample-wise DD shown at 10^3

Table 2: Target-Train-Token-to-Parameter ratio (TTT2P ratio) for WMT’14 EN-FR and FR-EN

	Number of lines						
	100	1,000	10,000	100,000	1,000,000	10,000,000	35,762,532
EN: num train tokens	3,248	33,768	313,154	3,123,129	30,852,455	308,640,462	1,174,344,513
FR: num train tokens	3,548	36,507	339,803	3,414,959	33,865,679	343,344,536	1,327,817,765
EN-FR num params	45,609,474	51,039,363	62,871,584	75,630,304	85,210,037	108,226,335	111,417,633
TTT2P ratio	0.000078	0.000715	0.005405	0.045153	0.397438	3.172468	11.917483
FR-EN num params	45,540,219	50,692,575	61,916,891	74,547,874	83,936,258	107,378,859	111,399,165
TTT2P ratio	0.000071	0.000666	0.005058	0.041894	0.367570	2.874313	10.541771

This shows that the effect we reported in § 5 also holds on these datasets: “the **ratio of target training token count to number of parameters** falls into $O(10^{-4})$ for 10^2 lines, $O(10^{-3})$ at 10^3 , $O(10^{-2})$ at 10^4 , and $O(10^{-1})$ for 10^5 lines and so on”.

⁹<http://www.statmt.org/wmt14/translation-task.html>

¹⁰<https://github.com/pytorch/fairseq/blob/master/examples/translation/prepare-wmt14en2fr.sh> and <https://github.com/pytorch/fairseq/blob/master/examples/translation/prepare-iwslt14.sh>

¹¹max-seq-len 300; checkpoint-frequency 4000 except for cases where 50 epochs would be reached before the first checkpoint: 400 for 10^2 lines and 3450 for 10^3 lines.

Table 3: Target-Train-Token-to-Parameter ratio (TTT2P ratio) for IWSLT’14 DE-EN and EN-DE

	Number of lines				
	100	1,000	10,000	100,000	160,239
DE: num train tokens	2,874	27,675	253,757	2,519,534	4,035,591
EN: num train tokens	2,739	26,416	245,659	2,461,879	3,949,114
DE-EN num params	45,297,348	49,410,683	53,639,825	55,189,376	55,428,584
TTT2P ratio	0.000060	0.000535	0.004580	0.044608	0.071247
EN-DE num params	45,405,078	49,809,797	54,300,056	56,245,643	56,564,366
TTT2P ratio	0.000063	0.000556	0.004673	0.044795	0.071345

K.2 TOKEN-TO-PARAMETER RATIO FOR NON-NEURAL MONOLINGUAL LMS

We experimented also on KenLM (Heafield, 2011; Heafield et al., 2013), a non-neural LM with modified Kneser-Ney smoothing (Kneser & Ney, 1995; Chen & Goodman, 1999), on our dataset A and found that on the word level, such a spike (or a hump) is common across all languages, see Figure 17. The target-token-to-parameter ratio is under 1 for most of these smaller data sizes. This seems related to the analytical findings in Oppen et al. (1990) where the pseudo-inverse solution to a simple learning problem was shown to exhibit non-monotonicity, with the peak exactly as the ratio of data to parameters (α) approaches 1.

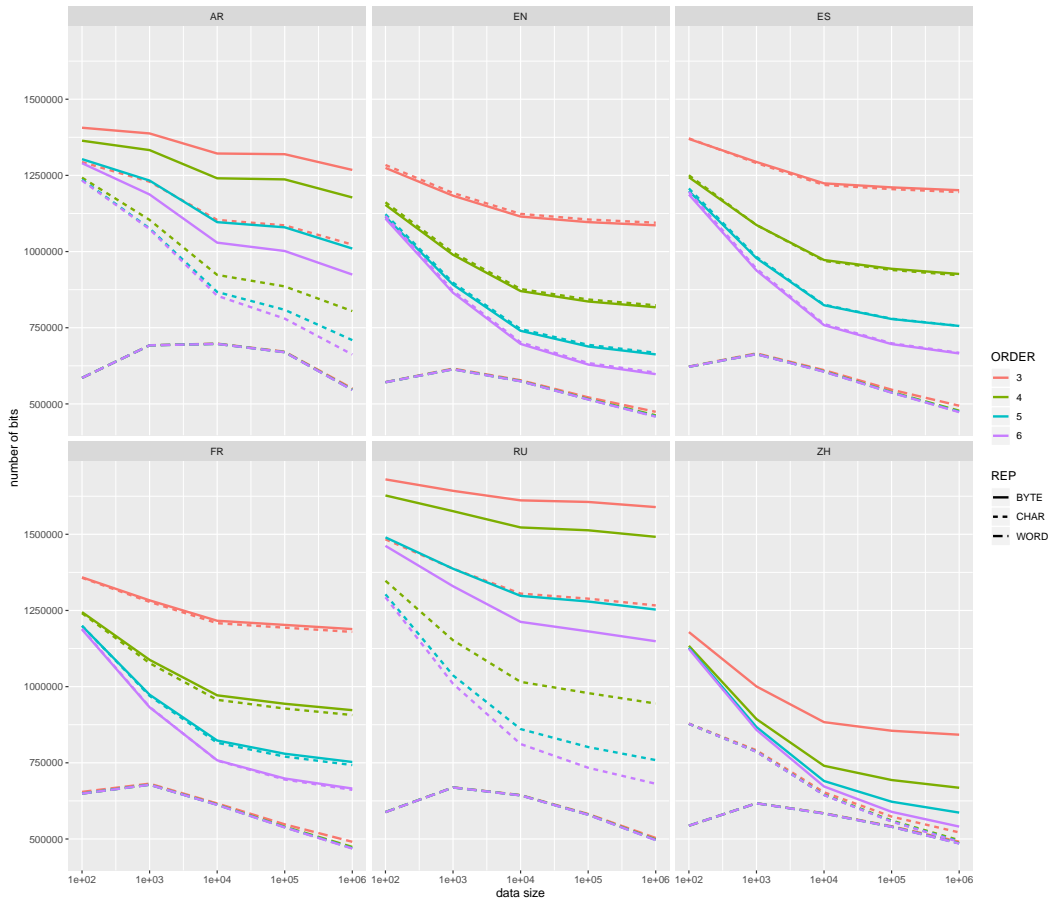


Figure 17: Kneser-Ney (monolingual) n-gram LMs on the same data (A) used for our neural CLMs

The number of parameters of a k -gram model is the number of unique n -grams, $1 \leq n \leq k$. Table 4 shows the ratios for our trigram model (all n -gram models of higher order exhibit the same effect).

On word level, where the function of number of bits to data size is not always monotonic, we observe less of a monotonic development whenever the token-to-parameter ratio is smaller than 1. This is more notably shown in the first 4 sizes in AR with a hump-like curve before the performance improves at 10^6 . This is different from the sharper descent for ES and FR, where only the first two data sizes have a non-monotonic relationship and a token-to-parameter ratio less than 1. Taking the token-to-parameter ratio as a rough proxy for over- (< 1) and under-parameterization (> 1), this can be seen as an instance of non-monotonicity with respect to data size in the “critical regime”, i.e. when the model transitions from being (heavily) over- to under-parameterized (Belkin et al., 2019; Nakkiran, 2019).

A remark on modeling with finer granularity Our KenLM results show the performance of bytes and characters is not on par with that of words with non-neural algorithms. NNs/DL has enabled much progress in this regard.

Table 4: Token-to-parameter ratios on non-neural monolingual trigram LMs

	lang_numlines	num_tokens	l unigrams l	l bigrams l	l trigrams l	num_params	tokens/params	
CHAR	AR_100	9079	85	925	2894	3904	2.325563525	
	AR_1000	123832	110	1577	8592	10279	12.04708629	
	AR_10000	1083517	152	3216	21479	24847	43.60755826	
	AR_100000	10625047	179	5114	44251	49544	214.4567859	
	AR_1000000	102064230	242	8517	90353	99112	1029.786807	
	EN_100	11730	78	806	2532	3416	3.433840749	
	EN_1000	159444	84	1215	5808	7107	22.43478261	
	EN_10000	1344001	125	2532	17181	19838	67.7488154	
	EN_100000	13132862	170	4231	36104	40505	324.2281694	
	EN_1000000	123491871	247	7126	70406	77779	1587.727677	
	ES_100	12374	87	781	2398	3266	3.788732394	
	ES_1000	171104	93	1210	5045	6348	26.95400126	
	ES_10000	1484804	117	2534	15462	18113	81.97449346	
	ES_100000	14549703	176	4261	33554	37991	382.9776263	
	ES_1000000	138596036	257	7217	67280	74754	1854.02836	
	FR_100	12456	89	836	2610	3535	3.523620934	
	FR_1000	179048	97	1259	5711	7067	25.33578605	
	FR_10000	1490983	133	2607	16282	19022	78.38203133	
	FR_100000	14528593	178	4390	35051	39619	366.707716	
	FR_1000000	138049189	259	7353	69522	77134	1789.732012	
	RU_100	11980	98	952	3051	4101	2.921238722	
	RU_1000	168156	111	1415	7106	8632	19.48053753	
	RU_10000	1436078	163	3506	20478	24147	59.4723154	
	RU_100000	14151728	190	5737	44071	49998	283.0458818	
	RU_1000000	134706120	263	10186	94975	105424	1277.755729	
	ZH_100	3318	605	2036	2634	5275	0.6290047393	
	ZH_1000	42572	1239	13266	24811	39316	1.082816156	
	ZH_10000	372003	2270	68178	175730	246178	1.511113909	
	ZH_100000	3659617	3403	241215	968852	1213470	3.015828162	
	ZH_1000000	34672612	4888	611213	3977112	4593213	7.548661906	
	BYTE	AR_100	16655	76	320	1163	1559	10.68313021
		AR_1000	227163	98	539	2070	2707	83.91688216
		AR_10000	1985014	133	1616	5974	7723	257.0262851
		AR_100000	19487689	148	2844	14274	17266	1128.674215
		AR_1000000	186171180	165	5219	40507	45891	4056.812447
		EN_100	11731	79	807	2533	3419	3.431120211
EN_1000		159449	85	1219	5812	7116	22.40711074	
EN_10000		1345771	130	2527	17139	19796	67.98196605	
EN_100000		13158948	154	3971	34985	39110	336.4599335	
EN_1000000		123705128	169	6422	66606	73197	1690.030029	
ES_100		12629	88	766	2414	3268	3.864443084	
ES_1000		175286	94	1146	4901	6141	28.54355968	
ES_10000		1513782	121	2409	14894	17424	86.87913223	
ES_100000		14821495	154	3925	31905	35984	411.8912572	
ES_1000000		141276766	169	6338	62199	68706	2056.250779	
FR_100		12875	90	830	2560	3480	3.699712644	
FR_1000		185227	99	1227	5497	6823	27.14744247	
FR_10000		1542105	133	2492	15615	18240	84.54523026	
FR_100000		15055657	156	4014	33105	37275	403.9076325	
FR_1000000		143495667	175	6423	64044	70642	2031.308103	
RU_100		21751	100	475	1365	1940	11.21185567	
RU_1000		309279	113	694	2732	3539	87.39163606	
RU_10000		2636591	151	1898	8430	10479	251.607119	
RU_100000		25990263	160	3364	18321	21845	1189.757977	
RU_1000000		247098758	169	6224	45935	52328	4722.113553	
ZH_100		8559	140	1524	3532	5196	1.647228637	
ZH_1000		116667	146	2706	12857	15709	7.426761729	
ZH_10000		1019969	156	5596	36176	41928	24.32667907	
ZH_100000		9990046	167	9228	81997	91392	109.3098521	
ZH_1000000		94268840	196	13407	160359	173962	541.893287	
WORD		AR_100	1776	869	1534	1669	4072	0.4361493124
		AR_1000	23460	5868	16064	20063	41995	0.5586379331
		AR_10000	206549	26108	116814	164062	306984	0.6728331118
		AR_100000	2035190	97997	776730	1383009	2257736	0.9014295737
		AR_1000000	19410502	304978	4297319	10005650	14607947	1.328763173
		EN_100	2071	682	1567	1869	4118	0.5029140359
	EN_1000	27398	3292	13148	19834	36274	0.7553068313	
	EN_10000	236659	12014	83397	155493	250904	0.9428665944	
	EN_100000	2339109	37264	428249	1117802	1583315	1.477349106	
	EN_1000000	21943139	122457	1818166	6505850	8446473	2.59790554	
	ES_100	2232	710	1605	1974	4289	0.5204010259	
	ES_1000	29461	3839	13199	20634	37672	0.7820397112	
	ES_10000	263024	15116	83900	160078	259094	1.01516824	
	ES_100000	2588791	49499	439584	1116177	1605260	1.612692648	
	ES_1000000	24654449	142809	1840029	6268684	8251522	2.987866844	
	FR_100	2298	745	1737	2072	4554	0.5046113307	
	FR_1000	32011	3881	14535	22608	41024	0.780299337	
	FR_10000	273195	13998	86815	170729	271542	1.006087456	
	FR_100000	2684982	42870	428339	1150965	1622174	1.655175092	
	FR_1000000	25595487	118204	1703399	6171437	7993040	3.202221808	
	RU_100	1854	886	1589	1734	4209	0.4404846757	
	RU_1000	24746	5433	15511	20035	40979	0.603870275	
	RU_10000	216638	23403	108516	162401	294320	0.7360627888	
	RU_100000	2150746	81342	670857	1306351	2058550	1.044786865	
	RU_1000000	20421965	236088	3295028	8617195	12148311	1.681053852	
	ZH_100	1751	630	1434	1614	3678	0.4760739532	
	ZH_1000	23568	3181	13998	19341	36520	0.6453450164	
	ZH_10000	207714	13137	96829	160642	270608	0.7675826287	
	ZH_100000	2038639	46941	554739	1278188	1879868	1.08445859	
	ZH_1000000	19361101	134492	2527710	8401311	11063513	1.749995774	

L NUMBER OF MODEL PARAMETERS

Number of model parameters for dataset A

Table with columns: Representation, Number of timesteps, CHAR, BYTE, WORD, BPE, and 1,000, 10,000, 100,000, 1,000,000, 10,000,000, 1,000,000,000. Rows list various representations like AR-EN, AR-ES, AR-FR, etc., with corresponding parameter counts for different datasets.

M ERRATICITY

Length has been an issue since the dawn of the encoder-decoder approach for NMT (Cho et al., 2014). Most work on length bias, except for that by e.g. Soutsov & Sarawagi (2016), seems to have focused on the evaluation of generated translation output and monitored performance degradation with respect to sequence length, often arguing that beam size plays a role (Koehn & Knowles, 2017; Murray & Chiang, 2018). (Related work in Stahlberg & Byrne (2019) provides a good summary on this issue.) While there could also be confounds in search, our experiments show that a kind of length bias can surface already with CLMing, without generation taking place. To our knowledge, length bias has not been expressed as a sample-wise non-monotonicity across a large data size range as ours. While the connection between erraticity in CLMs and length bias in NMT models remains to be verified on a case-by-case basis, the knowledge of length also contributing to robustness (not just consistently poor/poorer performance) could support further experimentation/replication of any study. Failed attempts to reproduce results may be explainable by erraticity.

One may argue that erraticity may not be relevant when each model is more optimally trained (as opposed to being treated with our one-setting-for-all regime). But we do want to stress that this very stark contrast between erratic and non-erratic behavior is possible, prompting a question on fairness: is there a one-for-all setting under which the languages with non-erratic behavior shown in our study would demonstrate erraticity and vice versa?

To the best of our knowledge, the meta phenomenon of erraticity, as a sample-wise non-monotonicity measured intrinsically with cross-entropy and contributing to large variance across runs, is a novel and original discovery and contribution to research in robustness. We hope our work would inspire further evaluation on other models/architectures, reflection and theories on our assumption of unbounded computation (e.g. Xu et al. (2020)), as well as new understanding and solutions that take data statistics and realistic computational aspects into account. We defer a more comprehensive analysis of erraticity with further experiments to future work.

M.1 ERRATICITY AS LARGE VARIANCE: EVIDENCE FROM DIFFERENT RUNS OF THE SAME DATA

To confirm that erraticity is not due to data-specific reasons, e.g. when certain data segments might be “easier” to model than others, we show figures from 2 runs (Figs. 18a and 18b) on the same dataset of wildly differing performance that only differ in seed. Note that changes in the y-direction can vary much, indicating large variance across runs.

By establishing that high variance holds across sample sizes, we showcased how it’d be possible to just test on 2 or 3 data points of smaller sizes to get a gauge on the robustness in higher order. It serves as a signal of when the system is being “stress-tested” and hyperparameters need re-tuning. Spot-testing on a couple of smaller data sizes can indeed save much time and energy. Take our run B0 byte models as an example: the training of the 10^2 -line model for EN-RU took 15 minutes, 10^3 40 minutes, 10^4 1 hour 50 minutes, and 10^5 3 hours 36 minutes. One can imagine how these would just be a fraction of training time for bigger models. (Likewise, for our ratio of target training token count to number of parameters — knowing when a representation might be prone to DD within a data size range could help prevent practitioners from prematurely declaring experimental results as negative or from unnecessarily rerunning an experiment because bigger data did not lead to better results.)

M.2 ADDITIONAL EXPERIMENT WITH LENGTH FILTERING TO 300 BYTES

Figure 19a and 19b show results of additional experiment with subset of data in byte (UTF-8) representation length-filtered to 300, including dev data:

Erraticity remains for AR and RU. Scores are lower, though they cannot be compared with the experiments in the main paper due to difference in dev data size (3,077 lines vs. 1,804 lines here). Number of total lines for train is 5,533,672 lines for each language, from which we took the initial 10^2 - 10^6 . As in our main experiments, we filtered out only whole lines, i.e. not by discarding the tails of longer lines. 300 bytes aren’t long sequences, but without data transform or hyperparameter tuning, things can look unfair. The EN translation of the longest RU line in this dataset is: “47. It is

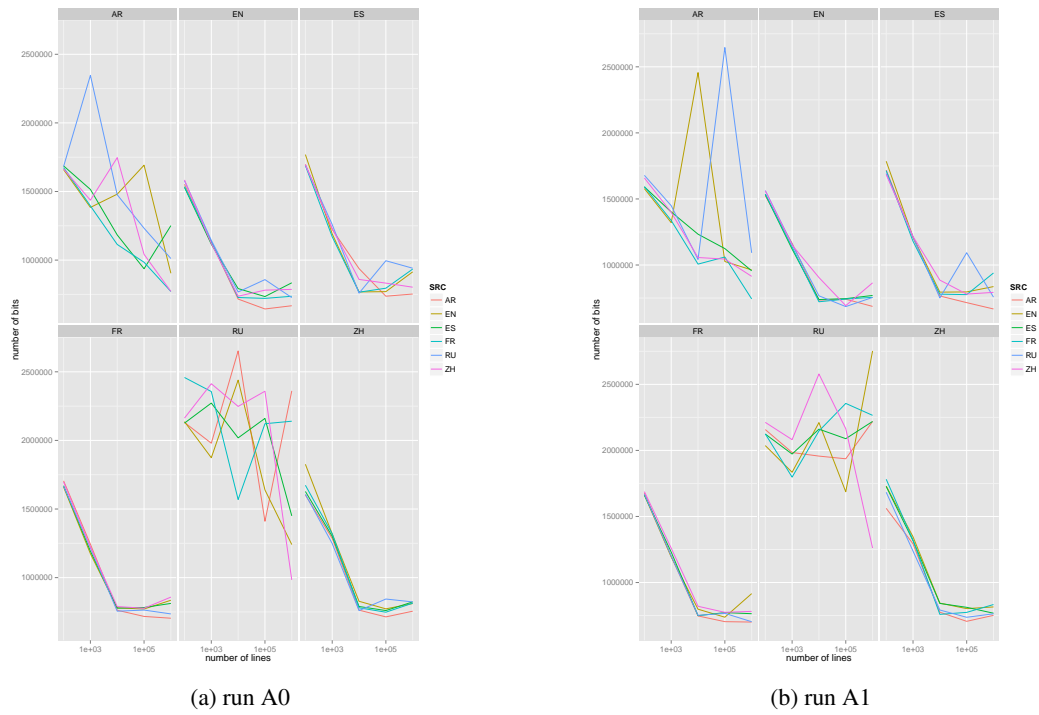


Figure 18: Same data with differing seeds

noted that there is a lack of information provided by the Government of Trinidad and Tobago with regard to the legal status of the Convention in the domestic legislation.”

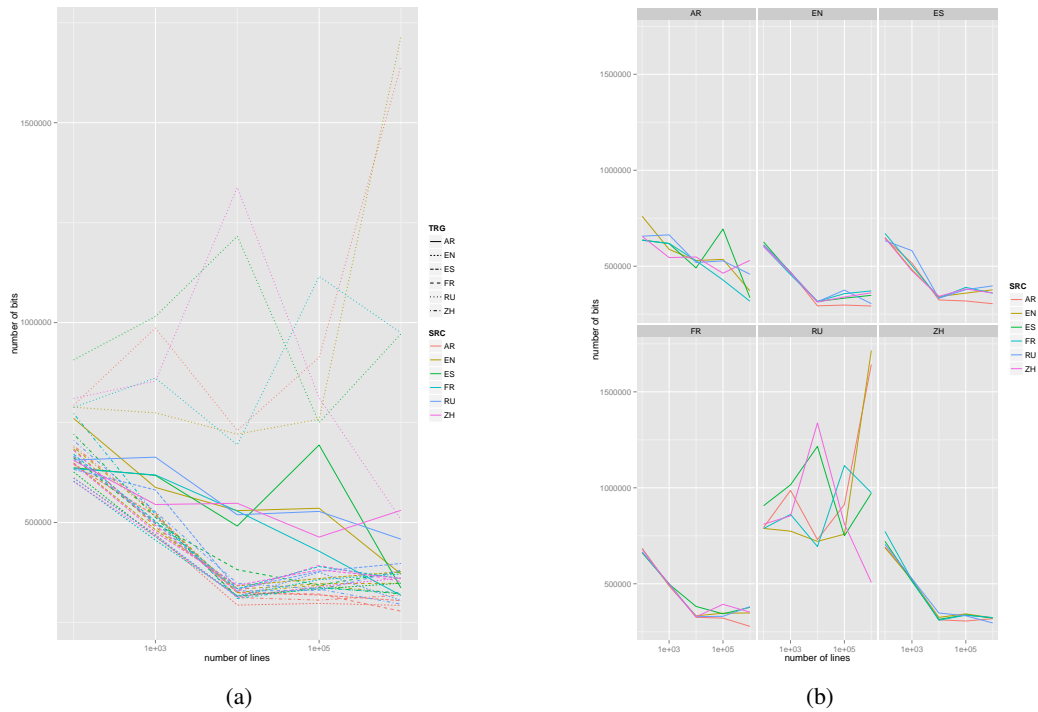


Figure 19: Additional experiment with maximum length of 300 bytes (with no hyperparameter tuning, in our blind one-setting-for-all evaluation). Considering there are languages with much higher character sequence length than RU, there is food for thought for the design of next-generation Multilingual Plane.

N EXPERIMENTS WITH ONE-LAYER TRANSFORMER

We performed 1 run with dataset A in 4 sizes (10^2 - 10^5 lines, seed=13) with the primary representations of characters, bytes, and words, on 1-layer Transformers (num-layers 1:1, all other hyperparameters remain the same as for our main experiments). We compared this against run A0 in 4 sizes with the same seed. (Based on how our null hypothesis is set up, the higher the number of runs, the more likely it is for there to be disparity. Important is that we evaluate based on an equal number of runs and on the same data for all candidates.) Results are shown in Table 5 with no statistically significant disparity observed on the models trained with 1 layer across the board.

Many are under the impression that big data is the cause to the neutralization of language instances in DL/NNs. But, as this set of experiments shows, it is possible for there to be no statistically significant differences between them, with as little as our smallest data size of 100 lines.

Table 5: Number of language pairs out of 15 with significant differences, with respective p-values. $\text{BYTE}_{6layers}$ is the representation with erratic AR_{trg} and RU_{trg} .

p-value	$\text{CHAR}_{6layers}$		$\text{BYTE}_{6layers}$		$\text{WORD}_{6layers}$		CHAR_{1layer}		BYTE_{1layer}		WORD_{1layer}	
	src	trg	src	trg	src	trg	src	trg	src	trg	src	trg
0.05	0	0	0	6	0	5	0	0	0	0	0	0
0.01	0	0	0	6	0	1	0	0	0	0	0	0
0.001	0	0	0	5	0	0	0	0	0	0	0	0

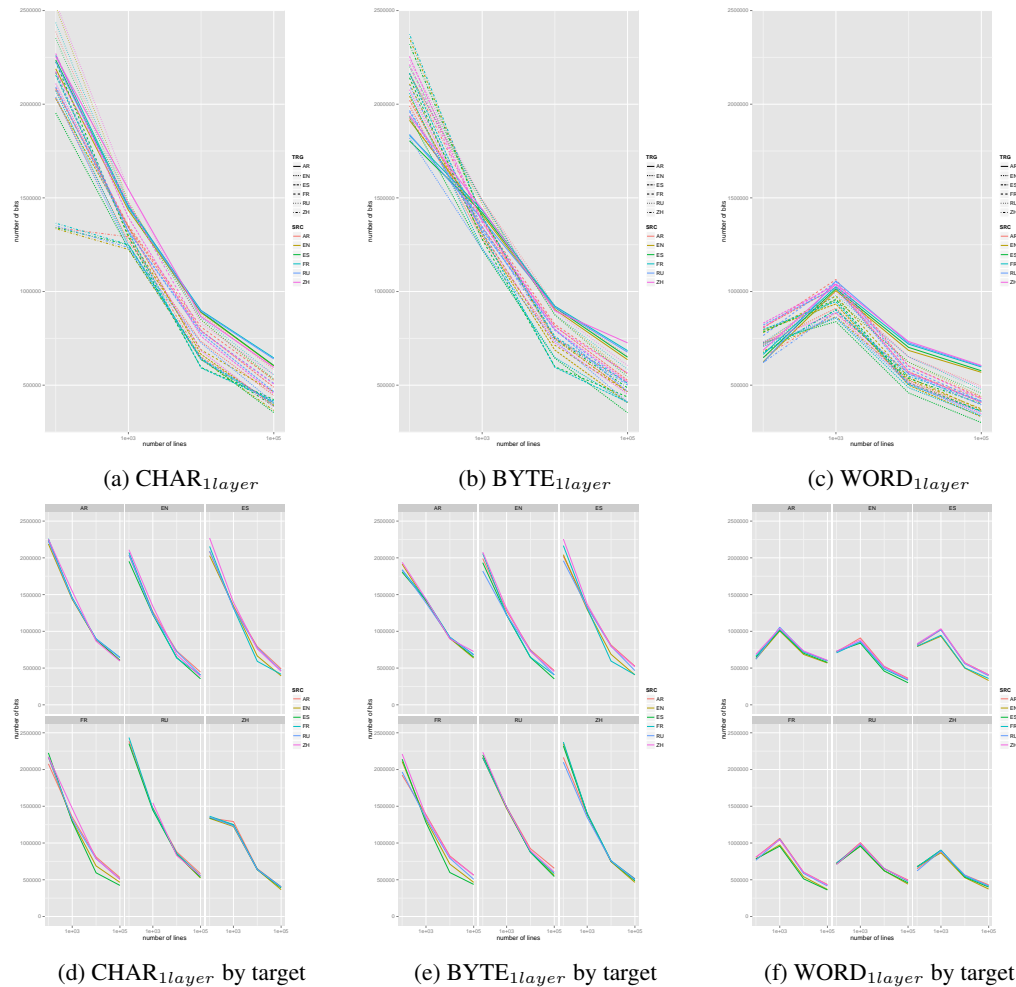


Figure 20: One-layer Transformer models

O PAQs (PREVIOUSLY ASKED QUESTIONS)

O.1 ONE SETTING FOR ALL

Q: Normally, one trains a model with the objective of optimizing based on the training and evaluation data with hyperparameter tuning. The experiments here used one setting for all. Some model configurations might train better and converge close to their optima while other configurations might not reach their full potential. Can this not create a distortion in the results?

A: For conventional engineering practice, we agree that hyperparameter tuning would be a sine qua non. However, the evaluation objective is the relational distance between languages, hence we need to see it in a different light. Here is a loose analogy:

Assume 3 objects in 3 different locations in space.

Relative evaluation from one setting allows one to capture the distance between these objects. It does not matter whether these three objects are in their “best” states.

For example, if one were to use a camera to capture these 3 objects and one does not adjust the setting (using just one random aperture, shutter speed, and focus), i.e. no tuning to capture any of these 3 specifically, nor does one try to model these 3 to their individual bests separately, what would result could be a picture that captures one of these 3 objects more favorably than the others, or it could be that all of these would be blurred. But either way, there is a degree of blurriness to be measured, giving us an idea of the relative distance between the objects. Such relative measurement is the evaluation strategy that our paper adopts.

Now, to add to the camera analogy, say one of the objects is running water, which was extra blurry [erraticity]: we suggest freezing the water, so even from the one arbitrary angle, it could be captured better. And it worked.

Also, while one might generally like to have a “pretty” photo, one that is e.g. taken with sub-optimal lighting, say, overexposure, can have a telling effect as it can bring out details in something dark, like a black box.

Alternatively, one can tune hyperparameters for each model individually such that each model would be a more optimized one and then compare these models. In that case, one would be interpreting the differences between language in terms of hyperparameters, and the paper would be one that is algorithm-centric. That is of course also a possibility. Our approach, however, is a data-centric one. We would, first of all, like to understand the nature of language data, i.e. what it is about language, if there is anything at all, that makes it a different data type than other data, and what kind of structural constraints, if any, that we need to take into consideration. Then with findings from this data perspective, we try to relate back to the algorithm and make connections so to create a more holistic picture.

O.2 TRANSLATIONESE / WORD ORDER

Q: Multitexts are parallel texts or translations with the same meaning. There is little to no variation in word order, hence they are just “Translationese” (Gellerstam, 1986). That is why they turn out to be the same, with no performance disparity.

A: Our findings do show that when the semantics is properly controlled, such as in multitexts, the factors influencing performance are statistical properties related to sequence length and vocabulary, e.g. $|V|$ or TTR, and the languages tested can be different. Semantic equivalence is also not a reason why we should expect neutralization of source language instances, as that would mean we should expect equal results across target languages.

We agree that faithfulness is often a priority in producing good translations. Whether the translations are produced by humans or machines, only a single best translation can surface as the translation of choice. There may be many other competing hypotheses, but regardless of whether it is done through an automatic ranking algorithm by a machine or through a human expert, the purpose of

translation is the same. However, *styles* and preferences in translations can vary. While faithfulness is generally preferred in the translations of legal texts, more freedom with skillful rearrangement of and play on words (or rather, character or sub-character sequences) or sounds being a criterion for literary texts could be appreciated by certain readers. We agree that it could be very interesting and necessary to model these variations, and we understand that languages can surface in many multimodal forms beyond the confines of texts as well. But with a data-driven perspective, to model this broader variation in language, we need corresponding datasets — we suggest contrast sets where the difference in e.g. sequential order is explicit. And for evaluation, we would require an even more systematic meta evaluation, one that spans different datasets.

But the argument that language or data *could be* different beyond how it appears in one dataset is irrelevant in the evaluation of experiments involving said dataset.

P UNDERSTANDING THE PHENOMENA WITH ALTERNATE REPRESENTATIONS (EXTENDED VERSION)

[Appendix P is an extended version of § 4.]

To understand why some languages show different results than others, we carried out a secondary set of control experiments with representations targeting the problematic statistical properties of the corresponding target languages.

Character level On the character level, it is well known that ZH differs from the other languages in its high $|V|$, in this study it has an averaged mean \pm std of 2550 ± 1449 ¹² across all 5 data sizes from all 3 datasets compared to 170 ± 87 from all other 5 languages combined, may these be in Latin or Cyrillic alphabet or the Abjad script. But what is often not known is that the character sequence length of logographic languages such as ZH is typically short (think and compare the sequence length of the Ancient Egyptian hieroglyphs or the Demotic script with that of the Greek script on the Rosetta Stone). Here in our case, the averaged mean sequence length in characters for ZH is 35 ± 19 , compared to 129 ± 71 from the other 5 languages. Heuristics to mitigate high $|V|$ often involve decomposition, which automatically resolve the problem of short sequence length. We tried 2 methods to lower character $|V|$ with representations in ASCII characters — Pinyin and Wubi. The former is a romanization of ZH characters based on their pronunciations and the latter is an input algorithm that decomposes character-internal information into stroke shape and ordering and matches these to 5 classes of radicals (Lunde, 2008). We replaced the ZH data with these formats *only on the target side* and reran the experiments involving ZH as a target language ($ZH_{tr.g}$) on the character level.

Results in Figure 2 and Table 1 show that the elimination of disparity on character level is possible if ZH is represented through Pinyin (transliteration), as in Subfigure 2c. But Wubi exhibits erraticity (Subfigure 2a). Wubi in our data has a maximum sequence length of 688 characters. As we shall also show in our byte-level analysis below, there are reasons to attribute length as cause to erraticity.

Decomposition into strokes may seem like a natural remedy analogous to decomposing an EN word into character sequences, but one needs to be mindful of not exceeding an optimal length given finite computation. Considering the ZH in the UN data is represented in simplified characters, decomposing traditional characters would surely complicate the problem. As there are also sub-character semantic and phonetic units (Zhang & Komachi, 2018) that can be exploited for information and aligned with character sequences of other alphabets, qualitative advances in this area can indeed be a new state of the art.

Byte level On the byte level, we observe irregularity for AR and RU. We find minimum sequence length of the target language to be one of the highest metrics correlating positively with the total number of bits ($\rho = 0.60$).¹³ Our data is based on 300 characters as maximum length per line. While we wanted to retain at least 75% of the UN data after length filtering, this length still renders a maximum sequence length that exceeds 100 words (the default maximum length for the word alignment model, GIZA++ (Och & Ney, 2003), in the traditional SMT pipeline). Translated into bytes with UTF-8 encoding, data with 300 characters maximum gives us, e.g. for the 10^6 -line datasets, an averaged mean \pm std of 185 ± 106 in length for AR and 246 ± 142 for RU, considerably larger than that for ZH (94 ± 53) and for EN/ES/FR ($\approx 145.41\pm 77$). With UTF-8 encoding, each character in AR, RU, and ZH contains 2 or more bytes. ZH typically has shorter line length in characters, compensating for the total byte sequence in length, even when most ZH characters are 3 bytes each. However, AR and RU generally have long line length in characters, so when converted to bytes, the sequence length remains long even when most of the characters might be just 2 bytes each. Results from our pairwise comparisons indicate 8 (non-directional) language pairs to be significantly different (see Table 1 under “BYTE”): ES-RU, EN-RU, FR-RU, RU-ZH, AR-RU, AR-EN, AR-ZH, and AR-FR — all involving AR or RU. (Appendix I lists also the language pairs with significant differences for other representations.)

¹²Figures are rounded to whole number. Complete tables of data statistics are provided in Appendix D.

¹³Top-3 correlates for each representation can be found in Appendix F.

Leveraging language-specific code pages can be a useful practical trick, a reminder that there are alternatives to UTF-8 for analyses and back-end processing if data is clean and homogeneous and if success of larger-scale prediction is not a concern. But one more sustainable alternative is to design a more adaptive and flexible character encoding scheme in general, taking into account the statistical profiles such as length (wrt characters and bytes) and sub-character (atomic/elementary/compound) information of all (or as many as possible) of the world's languages.

Word level The main difference between word and character/byte models is the absence of length as a top contributing factor correlating with performance. Instead, what matters more are metrics concerning word vocabulary, with top correlate being OOV token rate in the target language ($\rho = 0.66$). This is understandable as word segmentation neutralizes sequence lengths — the longer lengths in phonetic alphabetic scripts are shortened through multiple-character groupings, while the shorter lengths in logographic scripts (cf. difference in length for the 3 scripts on the Rosetta Stone, logographic scripts are typically shorter than phonetic ones) are lengthened by the insertion of whitespaces. To remedy the OOV problem, we use BPE, which learns a fixed vocabulary of variable-length character sequences (on word level, as it presupposes word segmentation) from the training data. It is more fine-grained than word segmentation and is known for its capability to model subword units for morphologically complex languages (e.g. AR and RU). We use the same vocabulary of 30,000 as specified in Junczys-Dowmunt et al. (2016). This reduced our averaged OOV token rate by 89-100% across the 5 sizes. The number of language pairs with significant differences ($p \leq 0.001$) reduced to 7 from 8 for word models, showing how finer-grained modeling has a positive effect on closing the disparity gap.

Version 1.1 (graphs to be updated, score tables added)