

# MorphBPE: Morphology-Aware Tokenization for Efficient LLM Training

Anonymous ACL submission

## Abstract

001 Tokenization fundamentally shapes NLP per- 043  
002 formance, affecting both efficiency and lin- 044  
003 guistic fidelity. While Byte Pair Encoding 045  
004 (BPE) underpins most Large Language Mod- 046  
005 els (LLMs), its frequency-driven merges often 047  
006 disregard morpheme boundaries, yielding in- 048  
007 consistent and semantically opaque segmenta- 049  
008 tions in morphologically rich languages. We 050  
009 introduce MorphBPE, a simple extension of 051  
010 BPE that constrains merge operations during 052  
011 tokenizer training to respect morpheme bound- 053  
012 aries, while leaving inference unchanged and 054  
013 fully compatible with existing LLM pipelines. 055  
014 We evaluate tokenization quality using two in- 056  
015 trinsic metrics: Morphological Consistency F1, 057  
016 which measures whether shared morphemes are 058  
017 assigned consistent token representations, and 059  
018 Morphological Edit Distance, which quantifies 060  
019 alignment with morpheme boundaries. We then 061  
020 train 300M and 1B parameter decoder-only 062  
021 LMs from scratch across four typologically di- 063  
022 verse languages: English, Russian, Hungarian, 064  
023 and Arabic, under identical vocabulary sizes 065  
024 and training settings. Across all languages, 066  
025 MorphBPE consistently improves intrinsic mor- 067  
026 phological coherence and reduces language 068  
027 model cross-entropy; moreover, token length 069  
028 statistics indicate that these gains are not at- 070  
029 tributable to materially shorter tokens. Finally, 071  
030 on the Belebele multilingual reading compre- 072  
031 hension benchmark, MorphBPE yields signif- 073  
032 icant improvements in morphologically rich 074  
033 languages such as Russian and Arabic. 075

034 **Availability:** The *MorphBPE* codebase, 076  
035 datasets, and tokenizer playground will be re- 077  
036 leased upon publication. 078

## 037 1 Introduction 079

038 Tokenization is a central design choice in modern 080  
039 NLP systems and a critical bottleneck for multilin- 081  
040 gual Large Language Models (LLMs). By mapping 082  
041 raw text into discrete units such as bytes (Gillick 083  
042 et al., 2016), characters (Al-Rfou et al., 2019),

subwords (Sennrich et al., 2016), or words, tok-  
enization directly determines vocabulary size, se-  
quence length, and the granularity at which lin-  
guistic regularities can be learned. Errors or in-  
consistencies introduced at this stage propagate  
through the entire modeling pipeline and can sub-  
stantially affect both training efficiency and down-  
stream performance (Sajjad et al., 2017; Adel et al.,  
2018). Despite growing interest in tokenization-  
free or character-level alternatives (Clark et al.,  
2022; Deiseroth et al., 2024), nearly all state-of-the-  
art LLMs, including Gemma (Team et al., 2024),  
LLaMA (Touvron et al., 2023), DeepSeek (Bi et al.,  
2024), and OpenAI’s GPT family, rely on Byte Pair  
Encoding (BPE) or closely related variants due to  
their favorable trade-offs between efficiency and  
coverage.

BPE is a frequency-driven algorithm that iter-  
atively merges common symbol pairs, making it  
well suited for concatenative morphology, as in  
English, where morphemes are typically formed  
by linear affixation. However, this same mecha-  
nism leads to systematic failures in languages with  
richer or more complex morphological systems. In  
non-concatenative languages such as Arabic and  
Hebrew, meaning is expressed through root-and-  
pattern morphology rather than linear affixation,  
and frequent substrings do not necessarily corre-  
spond to meaningful units (Khaliq and Carroll,  
2013). Agglutinative languages such as Hungar-  
ian, Turkish, and Korean further challenge BPE,  
as long sequences of productive affixes result in a  
large space of related word forms that BPE frag-  
ments inconsistently (Hakkani-Tür et al., 2000). As  
a result, standard BPE tokenizations often fail to  
align with true morpheme boundaries, producing  
subword units that are neither linguistically inter-  
pretable nor stable across related word forms. In  
practice, BPE segmentations in morphologically  
rich languages frequently introduce semantic am-  
biguity by reusing frequent substrings across un-

related words. For example, in Arabic the word الرحمن (Al-Rahman, “The Merciful”) may be segmented into من (min, “whom”), ال (al, “the”), and رح, even though من is semantically unrelated to the original word. Such segmentations force the model to disentangle spurious token-level associations during training, increasing the burden on representation learning. Similar phenomena arise in agglutinative languages, where morphemes expressing tense, number, or case are split inconsistently across tokens, undermining the model’s ability to generalize across inflected forms.

A natural solution is to incorporate morphological information into tokenization. However, purely morphology-driven segmentation has been shown to conflict with corpus statistics and can lead to inefficient vocabularies or brittle behavior when applied naively (Durrani et al., 2019; Marco and Fraser, 2024). This highlights a key open challenge for multilingual NLP: how to enforce morphological coherence in tokenization without sacrificing the statistical efficiency and scalability that make BPE attractive for large-scale LLM training.

In this work, we introduce *MorphBPE*, a simple and practical extension of BPE that constrains merge operations during tokenizer training to respect morpheme boundaries. Unlike purely morphology-based segmentation, MorphBPE remains fully data-driven, allowing frequent and productive morphemes to be preferentially captured while avoiding unnecessary fragmentation of rare or uninformative morphemes. The constraint is applied only during training, and the resulting tokenizer behaves identically to standard BPE at inference, introducing no additional runtime cost and remaining fully compatible with existing LLM architectures and deployment pipelines. A central motivation for MorphBPE is *morphological consistency*. Beyond aligning tokens with morpheme boundaries, consistency requires that words sharing morphemes are represented using shared subword units, and that shared tokens correspond to shared morphological content. By integrating morphological constraints into a frequency-based merge objective, MorphBPE yields more stable and interpretable subword representations, facilitating more efficient language model training and improved performance, particularly in morphologically rich languages.

**Contributions:** (i) We propose *MorphBPE*, a morphology-aware extension of BPE that con-

strains merge operations using morpheme boundaries during tokenizer training, while leaving inference unchanged and fully compatible with existing LLM pipelines. (ii) We introduce two intrinsic metrics for evaluating morphological quality of tokenizers: **Morphological Consistency F1**, which quantifies consistency across words sharing morphemes, and **Morphological Edit Distance**, which measures alignment with morpheme boundaries. (iii) We conduct controlled language model experiments with 300M and 1B parameter models trained from scratch across four typologically diverse languages—English, Russian, Hungarian, and Arabic—showing that MorphBPE consistently reduces cross-entropy under identical vocabulary sizes. (iv) We demonstrate that these intrinsic improvements translate into measurable gains on the Belebele multilingual reading comprehension benchmark, with statistically significant improvements in morphologically rich languages such as Russian and Arabic.

## 2 Background and Related Work

**Subword Tokenization in Language Models:** Subword tokenization has become a foundational component of modern neural language models, enabling a balance between open-vocabulary coverage and computational efficiency. Early approaches explored character-level modeling (Al-Rfou et al., 2019) and byte-level representations (Gillick et al., 2016), which offer robustness across scripts but often incur longer sequences and higher computational cost. Subword-based methods, most notably Byte Pair Encoding (BPE) (Sennrich et al., 2016), WordPiece, and the SentencePiece Unigram LM (Kudo and Richardson, 2018), emerged as practical compromises, allowing frequent patterns to be captured while decomposing rare words into reusable units.

BPE, originally proposed as a text compression algorithm (Shibata et al., 1999), was adapted for neural machine translation in 2016 and rapidly became the de facto standard in NLP and Large Language Models (LLMs). Its popularity stems from its simplicity, deterministic behavior, and effectiveness in controlling vocabulary size while handling out-of-vocabulary words. Numerous extensions have been proposed to mitigate its shortcomings, including BPE-Dropout (Provilkov et al., 2020), which injects stochasticity to improve generalization, sampling-based and probabilistic vari-

ants (Asgari et al., 2020), byte-level adaptations for improved robustness across scripts (Wang et al., 2020), and multilingual BPE schemes designed to encourage cross-lingual token sharing (Liang et al., 2023). Despite these advances, most BPE-based approaches remain purely frequency-driven and largely agnostic to linguistic structure.

**Morphology-Aware Tokenization:** The limitations of frequency-based tokenization in morphologically rich languages have motivated a growing body of work on morphology-aware tokenization. Early approaches often rely on explicit morphological analyzers or pre-tokenization strategies, segmenting words into morphemes prior to applying a standard subword tokenizer (Otani et al., 2020; Nzeyimana and Niyongabo Rubungo, 2022). Other methods incorporate morphological dictionaries or multi-view segmentation signals (Park et al., 2021), aiming to expose models to linguistically meaningful units.

More recent work has explored hybrid approaches that seek to balance morphological structure with statistical efficiency. For example, morpheme-aware or linguistically informed tokenizers encourage alignment between subwords and morpheme boundaries (Jabbar, 2023; Marco and Fraser, 2024), while analytical studies examine how subword segmentation interacts with morphological complexity and model performance (Weller-Di Marco and Fraser, 2024). However, many of these methods introduce additional preprocessing stages, require runtime morphological analysis, rely on stochastic objectives, or depart substantially from standard BPE training and inference pipelines. As a result, their adoption in large-scale LLM training remains limited.

**Positioning of MorphBPE:** MorphBPE is designed to address these limitations while preserving the practical advantages that have made BPE ubiquitous. Importantly, MorphBPE is *not* a runtime morphological analyzer, nor does it require morphological lookup or rule-based processing at inference time. All morphological information is used solely during tokenizer training to constrain merge operations, and the resulting tokenizer behaves identically to standard BPE at inference. What distinguishes MorphBPE from prior work is the point of integration and the simplicity of the constraint. Rather than pre-segmenting text (pre-tokenization) or post-processing tokenizations, MorphBPE injects morphological structure directly into the BPE merge objective by disallowing merges that cross

known morpheme boundaries. This yields a deterministic tokenizer, avoids the memory and sampling overhead of probabilistic approaches such as the Unigram LM, and maintains full compatibility with existing LLM architectures and training workflows. MorphBPE bridges the linguistic motivation and the data-driven thinking in a way that prior approaches have not fully achieved.

### 3 MorphBPE Method

MorphBPE is guided by four core design principles motivated by practical deployment requirements for large-scale multilingual LLMs. **(i) Minimal intervention:** Rather than redesigning tokenization from scratch, MorphBPE introduces a single, targeted modification to the standard BPE training procedure. The goal is to preserve the empirical strengths of BPE while correcting its systematic failures on morphological structure. **(ii) Determinism:** MorphBPE is fully deterministic. Unlike stochastic tokenization schemes such as BPE-Dropout or Unigram LM sampling, identical inputs and hyperparameters always produce the same tokenizer. This property is important for reproducibility and large-scale training stability. **(iii) Training-time supervision only:** Morphological information is used exclusively during tokenizer training to guide merge decisions. At inference time, MorphBPE behaves identically to standard BPE and does not require morphological annotation, lookup tables, or runtime analysis. **(iv) Pipeline compatibility:** The resulting tokenizer is a standard BPE model and integrates seamlessly into existing LLM training and inference pipelines without architectural or infrastructural changes.

#### 3.1 Algorithm

Byte Pair Encoding (BPE) begins with a character-level vocabulary and iteratively merges the most frequent adjacent symbol pairs until a target vocabulary size is reached (Sennrich et al., 2016). Merge decisions are purely frequency-based and unconstrained by linguistic structure.

MorphBPE modifies this process by introducing a single constraint. During tokenizer training, merges that cross known morpheme boundaries are disallowed. All other aspects of BPE, including frequency computation, merge ranking, and inference behavior, remain unchanged. Formally, let a word be segmented into a sequence of morphemes according to available morphological annotations.

---

**Algorithm 1** Morphology-Aware Byte Pair Encoding (MorphBPE)

---

```
1: Initialize vocabulary with individual characters
2: Obtain morpheme boundaries for the training corpus
3: while number of merges < target vocabulary size do
4:   Compute frequencies of adjacent symbol pairs
5:   Select the most frequent pair that does not cross a
   morpheme boundary
6:   Merge the selected pair and update the vocabulary
7: end while
```

---

286 During BPE training, a candidate merge between  
287 symbols  $x$  and  $y$  is permitted only if both symbols  
288 belong to the same morpheme span. If the merge  
289 would combine symbols originating from different  
290 morphemes, it is skipped in favor of the next most  
291 frequent valid pair.

292 Algorithm 1 summarizes the procedure.

293 Because constraints are applied only during  
294 merge selection, the learned merges define a stan-  
295 dard BPE tokenizer. No additional metadata or  
296 morphology-specific logic is required at inference.

### 3.2 Practical Properties

297 MorphBPE has the same asymptotic time complex-  
298 ity as standard BPE. The additional cost of check-  
299 ing morpheme boundary constraints is linear in the  
300 number of candidate merges and negligible rela-  
301 tive to frequency computation. Memory usage is  
302 unchanged, as morpheme boundaries are required  
303 only during tokenizer training and are not stored in  
304 the final tokenizer.  
305

306 From a systems perspective, MorphBPE is drop-  
307 in compatible with existing tokenization libraries  
308 and LLM training frameworks. Models trained  
309 with MorphBPE require no changes to architec-  
310 tures, batching strategies, or inference code. This  
311 makes MorphBPE suitable for large-scale training  
312 regimes where reproducibility, efficiency, and sim-  
313 plicity are critical.

## 4 Evaluation Framework

314 We evaluate MorphBPE at three complementary  
315 levels in order to isolate the effects of morphology-  
316 aware tokenization and avoid confounding factors.  
317 First, we assess intrinsic properties of the tokenizer  
318 itself, independent of any language model. Sec-  
319 ond, we evaluate intrinsic language model behav-  
320 ior during training using controlled cross-entropy  
321 comparisons. Finally, we measure downstream  
322 task performance to determine whether intrinsic  
323 improvements translate into practical gains. This  
324 layered evaluation design makes the logic of our  
325

experiments explicit and aligns with best practices  
for analyzing tokenization methods in LLMs.

### 4.1 Tokenizer-Level Evaluation (Intrinsic)

326 Tokenizer-level evaluation focuses on properties of  
327 the segmentation itself, without involving language  
328 model training. These metrics directly quantify  
329 how well a tokenizer aligns with morphological  
330 structure and how efficiently it represents text.  
331

332 **Fertility:** Fertility ( $\phi$ ) measures the average num-  
333 ber of subword tokens produced per word relative  
334 to a whitespace-based baseline (Rust et al., 2021).  
335 Lower fertility indicates higher compression effi-  
336 ciency and potentially longer effective context win-  
337 dows. However, fertility must be interpreted cau-  
338 tiously, as morphologically rich languages such as  
339 Hungarian and Arabic naturally require higher fer-  
340 tility to encode productive inflectional and deriva-  
341 tional processes. We therefore report fertility along-  
342 side morphology-sensitive metrics rather than treat-  
343 ing it as a standalone indicator of quality.  
344

345 **Morphological Edit Distance:** We introduce Mor-  
346 phological Edit Distance ( $\mu_e$ ), an intrinsic metric  
347 that quantifies alignment between tokenizer output  
348 and gold morpheme boundaries. Using a dynamic  
349 programming alignment that preserves token or-  
350 der, this metric measures the minimum number of  
351 insertions, deletions, and substitutions required to  
352 transform a token sequence into a morpheme se-  
353 quence. Lower values indicate better adherence to  
354 morphological structure and greater interpretability.  
355 We report raw edit distances to reflect the average  
356 number of boundary mismatches.  
357

358 **Morphological Consistency F1:** Morphological  
359 Consistency F1 ( $\mu_c$ ), inspired by Marco and Fraser  
360 (2024), measures whether segmentation decisions  
361 are consistent across words that share morphemes.  
362 Recall captures whether words with shared mor-  
363 phemes receive shared tokens, while precision mea-  
364 sures whether shared tokens correspond to shared  
365 morphemes. Their harmonic mean yields  $\mu_c$ . To  
366 ensure scalability, we cluster words using  $k$ -means  
367 ( $k = 100$ ), sample  $C = 50$  word pairs per clus-  
368 ter, and estimate scores with  $N = 10$  bootstrap  
369 resamples.

### 4.2 Language Model Evaluation (Intrinsic)

370 Intrinsic language model evaluation assesses how  
371 tokenization affects training dynamics and repre-  
372 sentational efficiency when all other factors are  
373 held constant.  
374

**Token-Level Cross-Entropy.** We train decoder-only language models from scratch using either BPE or MorphBPE and compare token-level cross-entropy loss during training. Cross-entropy reflects both convergence speed and the quality of learned representations and is closely related to perplexity while providing finer-grained resolution. Comparisons are performed only between models with identical architectures, training data, and vocabulary sizes to ensure fairness.

**Rationale for Fair Comparison.** Vocabulary size directly affects branching factors and loss values, making unequal vocabularies incomparable. We therefore fix vocabulary sizes per language and show empirically that MorphBPE and BPE produce nearly identical token-length distributions, ruling out explanations based on trivial token shortening (see Appendix B for detailed analysis). Under these controlled conditions, differences in cross-entropy can be attributed to segmentation quality rather than capacity or compression artifacts.

### 4.3 Downstream Evaluation (Extrinsic)

Extrinsic evaluation measures whether intrinsic improvements translate into gains on real language understanding tasks. We use the **LM Evaluation Harness** (Gao et al., 2024) to conduct evaluations under a standardized zero-shot protocol. We evaluate models on the Belebele multilingual reading comprehension benchmark (Bandarkar et al., 2024), which consists of multiple-choice questions derived from Flores-200 passages across 122 language variants. We report results for English, Russian, Hungarian, and Arabic, enabling direct comparison across morphological typologies.

#### Evaluation Protocol and Significance Testing:

For each example, models score candidate answers by conditional log-probability and select the most likely option. Accuracy is computed over all examples. To assess statistical significance, we apply McNemar’s test on paired predictions, combined with  $10^6$  bootstrap resamples of accuracy differences. Multiple comparisons are corrected using the Benjamini–Hochberg false discovery rate procedure with  $\alpha = 0.05$ . This protocol allows us to determine whether observed gains are robust rather than due to sampling noise.

## 5 Experimental Setup

This section details the data resources, tokenizer training procedure, and language model training

configuration used in our experiments. The goal is to ensure reproducibility and make explicit the controls used to enable fair comparisons between standard BPE and MorphBPE.

### 5.1 Morphological Resources

Morphological supervision is required only during tokenizer training. We use manually annotated and high-confidence automatically generated morpheme segmentations covering four typologically diverse languages: English, Russian, Hungarian, and Arabic.

For English, Russian, and Hungarian, morphological segmentations are obtained from the SIGMORPHON 2022 Shared Task on Morpheme Segmentation (Batsuren et al., 2022). These datasets provide high-quality gold annotations covering both inflectional and derivational morphology. For Arabic, which exhibits non-concatenative templatic morphology, we combine multiple complementary resources: the Arabic Treebank (ATB) (Taji et al., 2017), the Dialectal Segmentation Dataset (Darwish et al., 2018), and the Quranic Morphology dataset (Dukes and Habash, 2010). To increase coverage of frequent surface forms, we additionally include one million high-confidence Arabic segmentations generated using Farasa (Darwish and Mubarak, 2016).

All datasets are normalized to a unified segmentation format and deduplicated. Manually annotated resources are split into 80% training, 10% validation, and 10% test sets. Automatically generated Arabic segmentations are used only for tokenizer training and excluded from intrinsic evaluation. Dataset statistics are summarized in Table 1.

### 5.2 LLM Training Data

For language model training, we use the FineWeb2 corpus (Penedo et al., 2024), a large-scale multilingual web dataset covering over 1,000 languages. FineWeb2 provides sufficient data volume and linguistic diversity to support controlled monolingual training while adhering to the Chinchilla scaling law (Hoffmann et al., 2022).

For each language, we extract language-specific subsets using FineWeb2 metadata and train models on a fixed token budget. This ensures that differences in model behavior arise from tokenization rather than data scale. Using the same underlying text for both BPE and MorphBPE further isolates the effect of segmentation.

Table 1: Morphological segmentation dataset statistics used for BPE and *MorphBPE* training and evaluation across languages.

Language	Morphology Type	# of Words	Avg. Morphemes per Word
English	Fusional (low complexity)	571,495	2.33
Russian	Fusional (moderate complexity)	784,212	3.84
Hungarian	Agglutinative (high complexity)	930,312	3.22
Arabic	Templatic (high complexity)	1,395,835	2.50

### 5.3 Tokenizer Training Setup

Both BPE and MorphBPE tokenizers are trained using identical corpora and target vocabulary sizes. MorphBPE differs only in that morpheme boundary constraints are applied during merge selection, as described in Section 3. All other aspects of tokenizer training, including initialization, frequency computation, and merge ordering, remain unchanged.

Vocabulary sizes are selected separately for each language following a principled procedure. We train tokenizers with vocabulary sizes ranging from 8K to 96K in 8K increments and compute morphological edit distance on development sets. The smallest vocabulary size beyond which improvements are not statistically significant is selected using a paired *t*-test. This yields 24K for Hungarian, 64K for Russian, and 96K for English and Arabic. These sizes are then fixed for all tokenizer-level and language model experiments to ensure fair comparison.

### 5.4 Language Model Training Setup

We train decoder-only Transformer language models at two scales: a 300M-parameter model and a 1B-parameter model. All models use identical architectures, optimization settings, and training schedules, and differ only in the tokenizer used. Training is implemented using the LLaMA-Factory framework (Zheng et al., 2024).

The 300M models are trained on approximately 6B tokens, while the 1B models are trained on approximately 20B tokens, consistent with Chinchilla-optimal scaling. Optimization uses AdamW with cosine learning rate decay, identical batch sizes, and the same random seeds across tokenizer variants. All experiments are conducted on H100 GPUs, with total compute on the order of several thousand GPU-hours.

By tightly controlling model architecture, data, vocabulary size, and training budget, this setup ensures that any observed differences in training dynamics or downstream performance can be at-

tributed to the tokenizer rather than confounding factors.

## 6 Results

### 6.1 Tokenizer-Level Results

Figure 1 and Table 2 summarize intrinsic tokenizer results across English, Russian, Hungarian, and Arabic. Across all languages, MorphBPE achieves consistently lower Morphological Edit Distance ( $\mu_e$ ) and higher Morphological Consistency F1 ( $\mu_c$ ) than standard BPE, while incurring only a marginal increase in fertility ( $\phi$ ). The improvements are largest for morphologically rich languages, especially Hungarian and Arabic, where respecting morpheme boundaries prevents frequent but linguistically spurious merges.

For completeness, we also compare Morphological Consistency F1 ( $\mu_c$ ) against SentencePiece Unigram LM (Kudo and Richardson, 2018) under identical vocabulary sizes and corpora (Table 2). MorphBPE achieves higher  $\mu_c$  across all four languages, indicating more stable and morphologically coherent tokenization.

Overall, these results confirm that constraining merges at morpheme boundaries produces tokenizations that better preserve morphological structure and reduce segmentation ambiguity. This effect increases with morphological complexity, supporting the motivation for morphology-aware tokenization in agglutinative and templatic languages.

### 6.2 Language Model Results (Intrinsic)

Figure 2 reports training cross-entropy curves for the 300M and 1B models across the four languages. Under identical vocabulary sizes, architectures, and training data, MorphBPE consistently yields lower cross-entropy than BPE. The trend is stable across both scales and persists throughout training, with Figure 2 showing a representative window of approximately 14B tokens for readability.

The cross-entropy reductions indicate more efficient learning dynamics when token boundaries align with morphological structure. Improvements

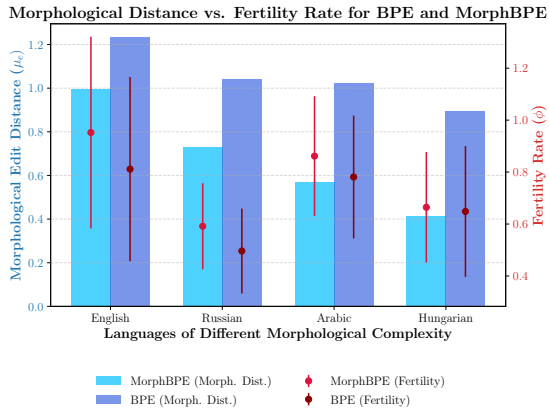


Figure 1: Comparison of morphological distance and fertility rate for BPE and *MorphBPE* across four languages. Lower fertility is generally preferred, and lower morphological distance indicates better alignment with morpheme boundaries.

are most pronounced for morphologically rich languages such as Hungarian and Arabic, where *MorphBPE* provides clearer subword regularities and reduces the burden of learning morphology from fragmented and inconsistent substrings. Importantly, these gains arise under controlled settings, suggesting they are attributable to segmentation quality rather than vocabulary capacity or data scale.

### 6.3 Downstream Results (Extrinsic)

We evaluate downstream effects using the Belebele multilingual reading comprehension benchmark (Bandarkar et al., 2024). Because our goal is to isolate tokenization effects, we report zero-shot performance for the 300M and 1B models without supervised fine-tuning. Absolute accuracy is therefore modest, but relative differences remain informative and can be assessed statistically.

Across the four languages, *MorphBPE* yields an average accuracy gain of approximately 1% over BPE. Gains are consistent for morphologically rich languages, Arabic, Russian, and Hungarian, and negligible for English, where morphology is less productive. Using McNemar’s paired test with  $10^6$  bootstrap resamples and Benjamini–Hochberg correction at  $\alpha = 0.05$ , improvements for Arabic and Russian are statistically significant, while the Hungarian gain is not.

These results suggest that morphology-aware tokenization can translate into measurable downstream benefits for morphologically complex languages, complementing intrinsic tokenizer metrics

and cross-entropy improvements. Together, the intrinsic and extrinsic results support the claim that improving morphological alignment at the tokenization stage strengthens representation learning and downstream comprehension.

## 7 Analysis and Discussion

**When and Why MorphBPE Helps:** *MorphBPE* is most effective in languages with rich and productive morphology, where surface word forms encode substantial grammatical and semantic information through affixation or non-concatenative processes. In agglutinative languages such as Hungarian and templatic languages such as Arabic, standard BPE frequently fragments words into statistically frequent but linguistically incoherent substrings. By constraining merges to respect morpheme boundaries, *MorphBPE* produces more interpretable and stable subword units that align with meaningful linguistic structure. This improved alignment manifests in higher morphological consistency and lower morphological edit distance, indicating that related word forms are segmented in a more systematic manner. From a representation learning perspective, such consistency reduces ambiguity in token semantics and allows the language model to reuse parameters across morphologically related forms more effectively. The resulting representations are therefore easier to learn and generalize, which is reflected in faster convergence and lower cross-entropy loss during training. In languages with relatively simple or weak morphology, such as English, the benefits of morphology-aware tokenization are naturally limited. English word formation relies less on productive inflection and more on fixed lexical items, and standard BPE already performs reasonably well at capturing frequent subword patterns. As a result, *MorphBPE* yields only modest improvements in intrinsic tokenizer metrics and negligible gains in downstream accuracy for English.

These findings are expected and highlight an important property of *MorphBPE*. The method does not degrade performance in low-morphology settings, but its advantages emerge primarily when morphological structure plays a central role in word formation. This behavior suggests that *MorphBPE* is a targeted improvement rather than a universally transformative change to tokenization.

**Tokenization, Fertility, and Vocabulary Size:** A recurring assumption in tokenizer evaluation is that

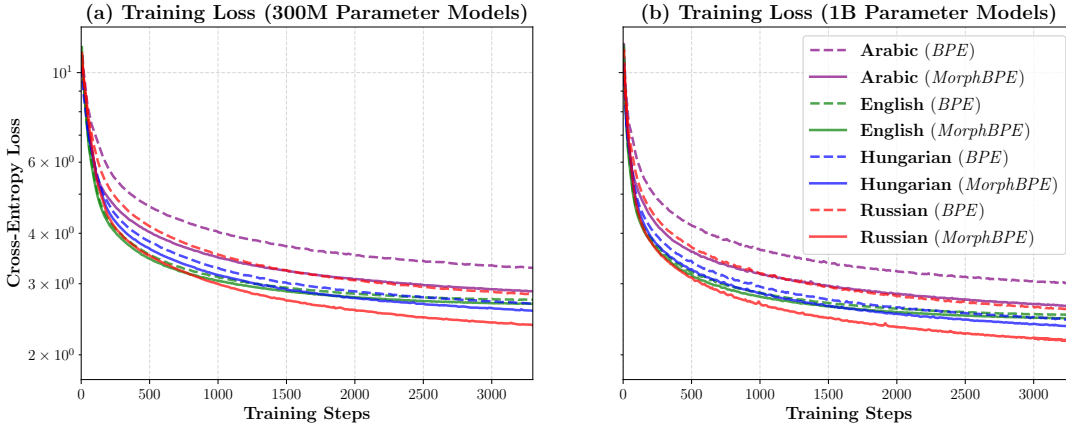


Figure 2: Training cross-entropy loss comparison between BPE and *MorphBPE* across English, Russian, Hungarian, and Arabic for both the 300M and 1B models (lower values indicate better performance).

Table 2: Morphological consistency evaluation for BPE, *MorphBPE*, and SentencePiece Unigram LM across four languages. Precision, recall, and F1-score ( $\mu_c$ ) are reported as mean  $\pm$  standard deviation over multiple resamples of the test sets. Higher F1-scores indicate greater consistency in segmenting words that share or differ in morphemes.

Model	Precision (Mean $\pm$ Std)	Recall (Mean $\pm$ Std)	Morph. F1 ( $\mu_c$ )
<b>English (96K)</b>			
BPE	0.00 $\pm$ 0.00	0.03 $\pm$ 0.02	0.00
<i>MorphBPE</i>	0.31 $\pm$ 0.44	0.34 $\pm$ 0.09	<b>0.32</b>
SentencePiece Unigram LM	0.20 $\pm$ 0.41	0.44 $\pm$ 0.07	0.28
<b>Russian (64K)</b>			
BPE	0.10 $\pm$ 0.32	0.06 $\pm$ 0.01	0.07
<i>MorphBPE</i>	0.69 $\pm$ 0.48	0.33 $\pm$ 0.06	<b>0.45</b>
SentencePiece Unigram LM	0.63 $\pm$ 0.49	0.22 $\pm$ 0.06	0.33
<b>Hungarian (24K)</b>			
BPE	0.08 $\pm$ 0.25	0.29 $\pm$ 0.04	0.13
<i>MorphBPE</i>	0.98 $\pm$ 0.03	0.78 $\pm$ 0.07	<b>0.87</b>
SentencePiece Unigram LM	0.93 $\pm$ 0.17	0.81 $\pm$ 0.10	0.87
<b>Arabic (96K)</b>			
BPE	0.00 $\pm$ 0.00	0.08 $\pm$ 0.03	0.00
<i>MorphBPE</i>	0.89 $\pm$ 0.31	0.53 $\pm$ 0.05	<b>0.66</b>
SentencePiece Unigram LM	0.73 $\pm$ 0.27	0.49 $\pm$ 0.04	0.58

638 lower fertility or shorter token sequences directly  
639 indicate better tokenization. Our results challenge  
640 this view. While *MorphBPE* sometimes produces  
641 slightly higher fertility than BPE, especially in mor-  
642 phologically rich languages, it consistently yields  
643 better morphological alignment and improved lan-  
644 guage model performance.

645 These findings indicate that fertility alone is an  
646 insufficient proxy for tokenizer quality. Vocabu-  
647 lary size and compression interact with linguistic  
648 structure in complex ways, and aggressive compres-  
649 sion can obscure systematic morphological patterns  
650 that are beneficial for learning. *MorphBPE* demon-  
651 strates that modest increases in token count can be  
652 offset by gains in interpretability, consistency, and  
653 representation stability, leading to more efficient

learning overall.

## 8 Conclusion

654  
655  
656 We introduced *MorphBPE*, a morphology-aware  
657 extension of Byte Pair Encoding that integrates lin-  
658 guistic structure into subword tokenization while  
659 preserving the efficiency, determinism, and com-  
660 patibility of standard BPE. Across four typologi-  
661 cally diverse languages, *MorphBPE* consistently  
662 improves morphological alignment, tokenizer con-  
663 sistency, and language model training dynamics,  
664 with downstream gains for morphologically rich  
665 languages.

666 Our results demonstrate that respecting mor-  
667 pheme boundaries during tokenizer training leads  
668 to more stable and interpretable representations,  
669 challenging the view that compression alone de-  
670 termines tokenizer quality. *MorphBPE* provides a  
671 practical and scalable way to incorporate lingu-  
672 stic insight into LLM training pipelines, and we hope  
673 it encourages further exploration of morphology-  
674 aware methods for multilingual language modeling.

## 9 Limitations

675 While *MorphBPE* offers significant advantages for  
676 morphologically rich languages, we acknowledge  
677 certain limitations that frame directions for future  
678 research: First, *MorphBPE* requires access to mor-  
679 phological segmentation data during the tokenizer  
680 training phase. It is important to clarify that this  
681 does not involve inference-time lookup tables, mor-  
682 phological lexicons, or rule-based systems. Instead,  
683 morphology-derived boundaries guide BPE merge  
684 decisions during training, resulting in a fully data-  
685 driven, standard BPE tokenizer. While high-quality  
686

resources exist for many morphologically rich languages (e.g., via UniMorph, SIGMORPHON, and MorphyNet, covering  $\sim 100$  languages), coverage remains limited for some low-resource and under-documented languages. However, the set of languages with reliable segmentation resources often aligns with those possessing sufficient corpora for meaningful LLM training.

Second, our current experiments focus on monolingual models to isolate the effects of tokenization across distinct typologies. Extending MorphBPE to multilingual models with joint vocabularies introduces new challenges, particularly in balancing morphological representation across diverse language families within a shared subword space. Investigating the interaction between morphology-aware constraints and cross-lingual vocabulary sharing is an open question.

Finally, our current study involved training 16 models from scratch across four languages and two scales (300M and 1B parameters), utilizing approximately 2,000 H100 GPU-hours. This represents a substantial effort to ensure rigorous, controlled comparison. While scaling further to larger model sizes and broader downstream tasks like instruction following or reasoning would require industrial-scale infrastructure, our findings are consistent with established links between improved perplexity and downstream performance (Wei et al., 2024), suggesting that the benefits of MorphBPE are likely to scale.

## References

Heike Adel, Ehsaneddin Asgari, and Hinrich Schütze. 2018. Overview of character-based models for natural language processing. In *Computational Linguistics and Intelligent Text Processing: 18th International Conference, CICLing 2017, Budapest, Hungary, April 17–23, 2017, Revised Selected Papers, Part I 18*, pages 3–16, Cham. Springer International Publishing.

Rami Al-Rfou, Dokook Choe, Noah Constant, Mandy Guo, and Llion Jones. 2019. Character-level language modeling with deeper self-attention. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 3159–3166.

Ehsaneddin Asgari, Masoud Jalili Sabet, Philipp Dufter, Christopher Ringlstetter, and Hinrich Schütze. 2020. Subword sampling for low resource word alignment. *arXiv preprint arXiv:2012.11657*.

Lucas Bandarkar, Davis Liang, Benjamin Muller, Mikel Artetxe, Satya Narayan Shukla, Donald Husa, Naman

Goyal, Abhinandan Krishnan, Luke Zettlemoyer, and Madian Khabza. 2024. *The belebele benchmark: a parallel reading comprehension dataset in 122 language variants*. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 749–775, Bangkok, Thailand. Association for Computational Linguistics.

Khuyagbaatar Batsuren, Gábor Bella, Aryaman Arora, Viktor Martinovic, Kyle Gorman, Zdeněk Žabokrtský, Amarsanaa Ganbold, Šárka Dohnalová, Magda Ševčíková, Kateřina Pelegrinová, Fausto Giunchiglia, Ryan Cotterell, and Ekaterina Vylomova. 2022. *The SIGMORPHON 2022 shared task on morpheme segmentation*. In *Proceedings of the 19th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology*, pages 103–116, Seattle, Washington. Association for Computational Linguistics.

Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiusi Du, Zhe Fu, et al. 2024. Deepseek llm: Scaling open-source language models with longtermism. *arXiv preprint arXiv:2401.02954*.

Jonathan H Clark, Dan Garrette, Iulia Turc, and John Wieting. 2022. Canine: Pre-training an efficient tokenization-free encoder for language representation. *Transactions of the Association for Computational Linguistics*, 10:73–91.

Kareem Darwish and Hamdy Mubarak. 2016. *Farasa: A new fast and accurate Arabic word segmenter*. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 1070–1074, Portorož, Slovenia. European Language Resources Association (ELRA).

Kareem Darwish, Hamdy Mubarak, Ahmed Abdelali, Mohamed Eldesouki, Younes Samih, Randah Alharbi, Mohammed Attia, Walid Magdy, and Laura Kallmeyer. 2018. *Multi-dialect Arabic POS tagging: A CRF approach*. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

Björn Deiseroth, Manuel Brack, Patrick Schramowski, Kristian Kersting, and Samuel Weinbach. 2024. *T-FREE: Subword tokenizer-free generative LLMs via sparse representations for memory-efficient embeddings*. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 21829–21851, Miami, Florida, USA. Association for Computational Linguistics.

Kais Dukes and Nizar Habash. 2010. *Morphological annotation of Quranic Arabic*. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)*, Valletta, Malta. European Language Resources Association (ELRA).

794	Nadir Durrani, Fahim Dalvi, Hassan Sajjad, Yonatan Belinkov, and Preslav Nakov. 2019. <a href="#">One size does not fit all: Comparing NMT representations of different granularities</a> . In <i>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)</i> , pages 1504–1516, Minneapolis, Minnesota. Association for Computational Linguistics.	851
795		852
796		853
797		854
798		
799		855
800		856
801		857
802		858
		859
		860
803	Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2024. <a href="#">The language model evaluation harness</a> .	861
804		862
805		863
806		864
807		865
808		866
809		867
810		
811	Dan Gillick, Cliff Brunk, Oriol Vinyals, and Amarnag Subramanya. 2016. <a href="#">Multilingual language processing from bytes</a> . In <i>Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1296–1306, San Diego, California. Association for Computational Linguistics.	868
812		869
813		870
814		871
815		872
816		873
817		874
818		
819	Dilek Z. Hakkani-Tür, Kemal Oflazer, and Gökhan Tür. 2000. <a href="#">Statistical morphological disambiguation for agglutinative languages</a> . In <i>COLING 2000 Volume 1: The 18th International Conference on Computational Linguistics</i> .	875
820		876
821		877
822		878
823		879
824	J. Hoffmann, S. Borgeaud, M. Arthur, E. Buchatskaya, T. Cai, E. Rutherford, D. Casas, L. Hendricks, J. Welbl, A. Clark, T. Hennigan, E. Noland, K. Millican, G. Driessche, B. Damoc, A. Guy, S. Osindero, K. Simonyan, E. Elsen, J. Rae, O. Vinyals, and L. Sifre. 2022. <a href="#">training compute-optimal large language models</a> . <i>Arxiv</i> .	880
825		881
826		882
827		883
828		884
829		
830		885
831	Haris Jabbar. 2023. <a href="#">Morphpiece: Moving away from statistical language representation</a> . <i>arXiv preprint arXiv:2307.07262</i> .	886
832		887
833		888
834		889
835		890
836		
837		891
838		892
839		893
840		894
841		895
842		896
843		897
844		898
845		899
846		
847		900
848		901
849		902
850		903
		904
		905
		906
		907
		908

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. [Neural machine translation of rare words with subword units](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.

Yusuxke Shibata, Takuya Kida, Shuichi Fukamachi, Masayuki Takeda, Ayumi Shinohara, Takeshi Shinohara, and Setsuo Arikawa. 1999. Byte pair encoding: A text compression scheme that accelerates pattern matching. *Technical Report DOI-TR-161, Department of Informatics, Kyushu University*.

Dima Taji, Nizar Habash, and Daniel Zeman. 2017. [Universal Dependencies for Arabic](#). In *Proceedings of the Third Arabic Natural Language Processing Workshop*, pages 166–176, Valencia, Spain. Association for Computational Linguistics.

Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. 2024. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Changhan Wang, Kyunghyun Cho, and Jiatao Gu. 2020. Neural machine translation with byte-level subwords. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 9154–9160.

Chengwei Wei, Yun-Cheng Wang, Bin Wang, C-C Jay Kuo, et al. 2024. An overview of language models: Recent developments and outlook. *APSIPA Transactions on Signal and Information Processing*, 13(2).

Marion Weller-Di Marco and Alexander Fraser. 2024. [Analyzing the understanding of morphologically complex words in large language models](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 1009–1020, Torino, Italia. ELRA and ICCL.

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, and Zheyang Luo. 2024. [LlamaFactory: Unified efficient fine-tuning of 100+ language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 400–410, Bangkok, Thailand. Association for Computational Linguistics.

## A Impact of Training Data Segmentation

A key advantage of MorphBPE is that it learns morphologically aware subword units directly from raw

text, constrained only during the tokenizer training phase. An alternative approach to incorporating morphology is to pre-segment the language model training data using a morphological analyzer and then apply standard BPE to the pre-segmented text.

To compare these approaches, we trained a BPE tokenizer on pre-segmented Arabic Wikipedia data (using Farasa segmentation) and subsequently trained a language model on this pre-segmented corpus. Figure 3 compares the training loss of this “Pre-segmented BPE” approach against standard BPE on raw text and MorphBPE on raw text.

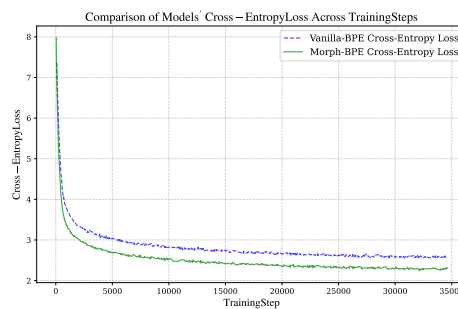


Figure 3: Comparison of Language Model training cross-entropy on Arabic using: (1) Standard BPE on raw text, (2) Standard BPE on pre-segmented text (Farasa), and (3) MorphBPE on raw text.

The results demonstrate that MorphBPE achieves a lower cross-entropy loss than both standard BPE on full text and the word-list approaches.

## B Token-Length Statistics and Fairness

A potential concern when comparing cross-entropy across tokenizers is that models with significantly shorter average token lengths (i.e., higher fertility) might artificially achieve lower per-token cross-entropy. To ensure that the improvements observed with MorphBPE are due to better segmentation quality rather than trivial differences in token granularity, we analyzed token-length distributions derived from large Wikipedia dumps for English, Hungarian, and Russian under the fixed vocabulary sizes used in our experiments.

Table 3 presents descriptive statistics of token lengths. MeanLen is the average token length in characters, and WMeanLen is the frequency-weighted mean token length. The results show that MorphBPE and BPE produce extremely similar token-length characteristics. In all three languages, the difference in mean and weighted mean token lengths between BPE and MorphBPE is miniscule (0.05–0.14 characters). This empirical evidence

1000  
1001  
1002  
1003

confirms that the experimental comparison is fair and that cross-entropy gains are attributable to improved morphological alignment rather than token-length artifacts.

Table 3: Wikipedia-derived token-length and distribution statistics comparing BPE and MorphBPE under fixed vocabulary sizes. The similarity in lengths confirms the fairness of cross-entropy comparisons.

Language / Tokenizer	Tokens	Vocab	MeanLen	WMeanLen	$\mu$	$\sigma$	$\alpha$
English / MorphBPE	7B	16K	5.23	2.87	9.79	2.68	1.01
English / BPE	7B	16K	5.28	2.91	9.87	2.64	1.01
Hungarian / MorphBPE	500M	16K	6.14	2.65	6.83	2.80	1.01
Hungarian / BPE	500M	16K	6.20	2.78	7.17	2.71	1.01
Russian / MorphBPE	3B	16K	6.41	2.42	7.64	3.38	1.01
Russian / BPE	3B	16K	6.27	2.64	8.17	3.06	1.01