

Parity-Aware Byte-Pair Encoding: Improving Cross-lingual Fairness in Tokenization

Anonymous ACL submission

Abstract

Tokenization is the first—and often least scrutinized—step of most NLP pipelines. Standard algorithms for learning tokenizers rely on frequency-based objectives, which favor languages dominant in the training data and consequently leave lower-resource languages with tokenizations that are disproportionately longer, morphologically implausible, or even riddled with <UNK> placeholders. This phenomenon ultimately amplifies computational and financial inequalities between users from different language backgrounds. To remedy this, we introduce Parity-aware Byte Pair Encoding (BPE), a variant of the widely-used BPE algorithm. At every merge step, Parity-aware BPE maximizes the compression gain of the currently worst-compressed language, trading a small amount of global compression for cross-lingual parity. We find empirically that Parity-aware BPE leads to more equitable token counts across languages, with negligible impact on global compression rate and no substantial effect on LM performance in downstream tasks.

1 Introduction

At a time of rapid innovation and constant change in natural language processing (NLP), tokenization continues to be a foundational and comparatively stable component of NLP pipelines. Tokenization is the transformation of raw sequences of bytes¹ into sequences of byte-spans, *i.e.*, subwords; it enables computational efficiency and provides essential inductive biases by defining meaningful textual units. This design choice can have a major impact on various aspects of model performance (Bostrom and Durrett, 2020; Ali et al., 2024; Goldman et al., 2024).

¹Early work considered characters the “base unit” of strings, but raw bytes have become popular because their fixed 256-symbol vocabulary can encode any character from any encoding, eliminating out-of-vocabulary issues.

The predominant tokenization algorithms—Byte Pair Encoding (BPE; Sennrich et al., 2016) and UnigramLM (Kudo, 2018), for example—select the vocabulary by maximizing frequency-based objectives computed over an entire training corpus. In multilingual corpora, this global criterion inevitably favors the languages with the greatest representation. Under vocabulary size constraints, subwords that primarily benefit high-resource languages are preferentially included, often at the expense of those needed for lower-resource languages. This bias has both qualitative and economic consequences. NLP models trained on fragmented or semantically incoherent tokenizations lose valuable inductive biases and tend to perform worse. At the same time, texts in lower-resource languages—often tokenized into more tokens—incur higher computational costs from language-model-based services charging based on token count, which disproportionately burdens users of underrepresented languages and exacerbates existing inequalities.

In the effort to mitigate these inequities, we introduce Parity-aware BPE. The classic version of BPE learns its vocabulary by repeatedly selecting the subword pair with the highest corpus-level co-occurrence count; it adds the concatenation of this pair to the vocabulary and replaces all pair co-occurrences with the new symbol.² Parity-aware BPE is a simple variant of this algorithm, retaining the iterative framework but redefining the merge selection rule: at each step, it computes co-occurrence statistics separately for each language and then uses statistics from the language with the current worst compression rate for selecting the next merge. In other words, instead of greedily maximizing a global objective, Parity-aware BPE performs a “fair-max” update that progressively equalizes string compression

²This process acts as a form of data compression, replacing frequent subword sequences with shorter representations.

076 rates across languages. Notably, this modification
 077 affects only vocabulary learning; the inference
 078 procedure remains the same as in Classical BPE.

079 Empirically, Parity-aware BPE leads to better
 080 token-count parity across languages than Classical
 081 BPE while maintaining comparable global com-
 082 pression. Across 13 multilingual benchmarks, mod-
 083 els trained with a Parity-aware tokenizer match
 084 or exceed downstream performance compared to
 085 those trained with a Classical BPE tokenizer. Fair-
 086 ness metrics improve significantly, reflecting a
 087 more even distribution of token costs without sac-
 088 rificing efficiency or downstream performance. In
 089 short, Parity-aware BPE reduces tokenizer-induced
 090 disparities between languages, enabling more equi-
 091 table resource allocation and balanced performance
 092 across diverse linguistic settings.

093 2 Text Tokenization

094 Text can be decomposed at many granularities:
 095 graphemes, Unicode code points, or multi-
 096 character tokens; but at the most fundamental digi-
 097 tal layer, every string is represented as a sequence
 098 of **bytes**, the foundation on which all other units
 099 are constructed. Let $b \in \mathcal{B} = \{0, \dots, 255\}$ de-
 100 note an individual byte. A finite **byte-string** is
 101 written $\mathbf{b} \in \mathcal{B}^*$, where $\mathbf{b} = b_1 b_2 \dots b_{|\mathbf{b}|}$. Through-
 102 out, bytes are treated as atomic symbols. Note that
 103 all subsequent definitions apply if one substitutes
 104 bytes with another finite alphabet of atomic units.

105 2.1 Byte-level Tokenizers

106 In plain terms, **tokenization** is the act of mapping
 107 raw byte-strings (sequences of bytes) to sequences
 108 of subwords. A **tokenizer** specifies the rules that
 109 perform this mapping—as well as the rules that con-
 110 vert sequences of subwords back into byte-strings.
 111 The following components can define a tokenizer:

- 112 • **Vocabulary** $\mathcal{V} \subset \mathcal{B}^+$: a finite set of non-empty
 113 byte-spans, often called *subwords*.³
- 114 • **Tokenization function** $\tau : \mathcal{B}^* \rightarrow \mathcal{V}^*$: a mapping
 115 from byte-strings \mathbf{b} to sequences of tokens $\mathbf{v} =$
 116 v_1, v_2, \dots
- 117 • **Detokenization function** $\perp : \mathcal{V}^* \rightarrow \mathcal{B}^*$: a map-
 118 ping from sequences of subwords to byte-strings.
 119 This operation is often just simple string concat-
 120 enation (denoted as \circ) of subwords’ corresponding
 121 byte spans: $\perp(v_1, \dots, v_n) = v_1 \circ v_2 \circ \dots \circ v_n$

³To guarantee representability of any byte-string, we as-
 sume \mathcal{V} includes all singleton bytes: $\mathcal{B} \subseteq \mathcal{V}$.

Pre-tokenization and Normalization. Many to-
 kenization algorithms include a *pre-tokenization*
 (and often a normalization) step that segments or
 rewrites raw byte strings according to deterministic
 criteria. Pre-tokenization can encompass several
 operations, including Unicode normalization or
 splitting on whitespace. Notably, pre-tokenization
 determines subword boundaries, and thus also de-
 termines the set of possible candidates for the vo-
 cabulary as well as the attainable compression rate.
 For example, if whitespace is used as a subword
 boundary, languages without explicit whitespace
 (e.g., Chinese, Japanese) or with rich morphology
 may have the potential for higher compression. For
 simplicity, we assume this step is baked into τ .

137 2.2 Text Compression

138 *Why does mapping a raw byte-string to a sequence*
 139 *of larger subword tokens help a language model*
 140 *(LM)?* The precise inductive biases this process
 141 imbues remain an open research question (Zouhar
 142 et al., 2023a; Schmidt et al., 2024), but one plau-
 143 sible explanation is the *compression* it provides:
 144 a good tokenizer tends to map each input \mathbf{b} to a
 145 shorter sequence of tokens, reducing the length of
 146 the model’s effective input and, potentially making
 147 learning easier. At the very least, it can significantly
 148 reduce model-side computations.

149 For a fixed tokenizer $(\mathcal{V}, \tau, \perp)$ we define the **com-**
 150 **pression rate** of a byte-string \mathbf{b} as

$$151 \text{CR}(\mathbf{b}; \tau) \stackrel{\text{def}}{=} \frac{|\mathbf{b}|_u}{|\tau(\mathbf{b})|} \quad (1)$$

152 where $|\mathbf{b}|_u$ denotes the length of \mathbf{b} in terms of
 153 a given **normalization unit** u (e.g., characters,
 154 words, lines, or simply bytes). In words, $\text{CR}(\mathbf{b}; \tau)$
 155 measures the multiplicative factor by which our
 156 original sequence length is reduced after tokeniza-
 157 tion. A higher CR indicates stronger compression.

158 We are generally interested in a tokenizer’s av-
 159 erage compression, which can be estimated from a
 160 corpus \mathcal{D} :

$$161 \text{CR}(\mathcal{D}; \tau) \stackrel{\text{def}}{=} \frac{1}{|\mathcal{D}|} \sum_{\mathbf{b} \in \mathcal{D}} \frac{|\mathbf{b}|_u}{|\tau(\mathbf{b})|} \quad (2)$$

162 This quantity provides an estimate of the average
 163 number of tokens that an autoregressive LM must
 164 process per unit u of raw text. Tokenizers differ
 165 in how small they can make $\text{CR}(\mathcal{D}; \tau)$ while
 166 remaining lossless. Further, this rate can vary
 167 across strings from different languages, which

motivates our last definition: language-specific compression rate.

Let $\mathcal{L} = \{\ell^{(r)}\}_{r=1}^R$ be our set of languages and $\mathcal{M} = \{(\mathbf{b}^{(s)}, \ell^{(s)})\}_{s=1}^S$ be a labeled multilingual corpus, *i.e.*, a corpus where each byte-string $\mathbf{b}^{(s)}$ is labeled with its respective language. For a fixed tokenizer, we define a language’s compression rate as

$$\text{CR}(\ell; \tau) \stackrel{\text{def}}{=} \text{CR}(\mathcal{D}_\ell; \tau) \quad (3)$$

where $\mathcal{D}_\ell = \{\mathbf{b}^{(s)} : (\mathbf{b}^{(s)}, \ell^{(s)}) \in \mathcal{M}, \ell^{(s)} = \ell\}$.

Compression Rate as a Notion of Fairness.

Many commercial technology services offer APIs that bill per token; services’ processing speeds also scale with the number of tokens in the input and output. Token counts thus dictate both the economics and the latency of these services. For a given byte sequence \mathbf{b} , the number of tokens produced by tokenization function τ is determined by the tokenizer’s compression rate $\text{CR}(\mathbf{b}; \tau)$. $\text{CR}(\ell; \tau)$ —the expected compression rate over sequences from language ℓ —is thus a direct proxy for the average per-language cost (and expected latency) of using one of these services. Variance in $\text{CR}(\ell; \tau)$ across languages thus implies different user costs purely as a function of language choice (Petrov et al., 2023), and whether $\text{CR}(\ell; \tau)$ is comparable across languages is therefore one way of measuring tokenizer fairness. Byte Pair Encoding, discussed next, optimizes for compression rate across a corpus without regard for language. Our subsequent adjustment to Classical BPE adds an auxiliary objective of equalizing $\text{CR}(\ell; \tau)$ across languages.

2.3 Byte Pair Encoding

Here we provide a detailed description of the classic version of BPE; readers already familiar with the algorithm can skip this subsection.

Byte Pair Encoding (BPE; Sennrich et al., 2016) is one popular algorithm for creating a tokenizer adapted from the byte-pair compression scheme of Gage (1994). In short, BPE tokenizes text by iteratively merging adjacent tokens whose token-types (*i.e.*, subwords) were observed to co-occur frequently in the training data.

The notion of a **merge** lets us formalize this procedure. A merge is defined as an ordered pair $m = (v, v')$ with $v, v' \in \mathcal{V}$. The application of a merge to a token sequence replaces every bigram token v, v' by a single token $v \circ v'$. Each bigram token replacement shortens the token sequence by exactly one token, thereby *compressing* the sequence. To

tokenize a piece of text with a BPE tokenizer, we start from its representation as a byte-string, *i.e.*, a sequence of base bytes—all of which necessarily appear in our tokenizer’s vocabulary. We then iteratively apply a given list of merges to that sequence. Note that because the merge list is fixed in advance, the encoding is deterministic. Intuition for the merge procedure is perhaps best acquired by a small example:

Example 2.1 (Example of the iterative application of merge sequence \mathbf{m} to byte sequence \mathbf{b}).

$$\mathbf{m} = [(b, a), (ba, b)]; \quad \mathbf{b} = \text{babab}$$

$$\mathbf{v}_0 = b, a, b, a, b$$

$$\text{Step 1: } b, a \rightarrow ba \implies \mathbf{v}_1 = ba, ba, b$$

$$\text{Step 2: } ba, b \rightarrow bab \implies \mathbf{v}_2 = ba, bab$$

In terms of our earlier tokenizer notation, a BPE tokenizer is defined as follows:

- $\mathcal{V} = \mathcal{B} \cup \{v \circ v' : (v, v') \in \mathbf{m}\}$
- $\tau_{\mathbf{m}}$ carries out the procedure described above, *i.e.*, it applies each m_k to an input byte-string \mathbf{b} in the prescribed order. In example 2.1, $\tau_{\mathbf{m}}(\mathbf{b}) = ba, bab$.
- $\perp(v_1, \dots, v_n) = v_1 \circ v_2 \circ \dots \circ v_n$

We use the \mathbf{m} subscript here to make explicit the tokenization function’s dependence on \mathbf{m} .

Learning \mathbf{m} . The BPE algorithm seeks the merge list \mathbf{m}^* (subject to a size constraint K) that maximizes the compression rate of the given corpus:

$$\mathbf{m}^* = \max_{\mathbf{m}: |\mathbf{m}|=K} \text{CR}(\mathcal{D}; \tau_{\mathbf{m}}) \quad (4)$$

BPE takes a greedy approach to choosing \mathbf{m} , finding an approximate solution to eq. 4 (Zouhar et al., 2023b). It starts with the singleton-byte vocabulary $\mathcal{V}_0 = \mathcal{B}$ and repeatedly greedily enlarges the vocabulary. At each of K steps, the current tokenizer $\tau_{\mathbf{m}_{<k}}$ is applied to the entire training corpus, and the algorithm counts how often every adjacent pair of tokens occurs. The subword-type pair with the highest count, which we denote as (v^*, v'^*) , is deemed the most “compressive.” Its concatenation $v^* \circ v'^*$ is added to the vocabulary, the merge $m_k = (v^*, v'^*)$ is recorded, and every occurrence of the bigram (v^*, v'^*) in the corpus is replaced by the new token so the next iteration works with updated token sequences. Repeating this process K times yields the ordered list $\mathbf{m} = [m_1, \dots, m_K]$ and the final vocabulary \mathcal{V}_K . When encoding a new text, $\tau_{\mathbf{m}}$ simply applies these merges in the same order. The pseudocode for the algorithm is provided in Alg. 1 in §A.

3 Parity-aware Byte Pair Encoding

Classical BPE chooses merges that maximize a *global* frequency objective, implicitly favoring the compression of languages with a larger presence in the training corpus. Here we introduce **Parity-aware BPE**, which replaces this global objective with a *max-min* criterion: at every step, it selects the merge that most improves the language currently suffering the poorest compression rate. In this section, we formalize the objective and describe the resulting algorithm.

3.1 Greedy min-max objective

Our adjustment to the Classical BPE objective (eq. 4) explicitly encodes our earlier notion of tokenizer fairness: equality across per-language compression rates. Formally, parity-aware BPE seeks a merge list $\mathbf{m} = [m_1, \dots, m_K]$ that maximizes the *minimum* compression rate across languages:

$$\mathbf{m}^* = \max_{\mathbf{m}:|\mathbf{m}|=K} \min_{\ell} \text{CR}(\ell; \tau_{\mathbf{m}}) \quad (5)$$

This min-max objective trades a small amount of global compression for fairness across languages.

3.2 Algorithm

Parity-aware BPE retains the greedy iterative framework of Classical BPE but changes *which* statistics are inspected each time a merge is added. At merge step k , it identifies the language with the worst compression under the tokenizer defined by the merge list thus far ($\mathbf{m}_{<k}$)

$$\ell^* = \arg \min_{\ell \in \mathcal{L}} \text{CR}(\ell; \tau_{\mathbf{m}_{<k}}) \quad (6)$$

To choose the next merge, it uses the same maximum pair count criterion as Classical BPE, albeit with pair counts computed over only \mathcal{D}_{ℓ^*} —the portion of the corpus corresponding to ℓ^* . The rest of the algorithm follows the Classical BPE procedure: the chosen merge is applied to all texts (*i.e.*, across $\mathcal{D}_{\ell} \forall \ell$)⁴ and the procedure is repeated for $k = 1, \dots, K$, yielding the final merge list \mathbf{m} . We provide pseudocode in Alg. 2.

Cross-lingual Compression Rate Comparison.

Parity-aware BPE relies on the comparison of $\text{CR}(\ell; \tau_{\mathbf{m}})$ across different ℓ . The choice of normalization unit u has a large impact on the measured $\text{CR}(\ell; \tau_{\mathbf{m}})$ and even when u is held constant

⁴Crucially, this is what distinguishes our algorithm from a combination of monolingual merge lists (Petrov et al., 2023), allowing us to find more “compressive” merges.

across measurements for different languages, if not considered carefully, the choice can introduce bias into the comparison. As concrete examples, certain normalization units are more appropriate in some languages than in others, *e.g.*, whitespace-delimited “words” are ill-defined in many languages; although principled and universal, even normalizing by number of bytes can skew perceived compression because scripts differ greatly in average number of bytes per character (*e.g.*, ASCII vs. UTF-8 CJK). Parallel corpora provide a principled solution: computing compression rates over aligned texts (sentences, lines, or documents) normalizes by content, making cross-language comparisons more meaningful. We therefore recommend the use of a parallel corpus for computing eq. 6. Notably, this evaluation corpus need not be the same one used for computing subword pair frequency statistics, for which a larger corpus with only language annotation is necessary. For generality, we thus differentiate between the corpora used to compute frequency statistics and in computing eq. 6, referring to them as our training and development datasets, respectively. Alg. 2 makes this difference explicit. We present experimental results both with a separate, parallel development set and using a single (not parallel) multilingual dataset for all computations.

Complexity and Data Requirements. Relative to Classical BPE, Parity-aware BPE incurs only a $O(|\mathcal{L}|)$ overhead per-merge from recomputing the language-specific compression rates on the dev set. Parity-aware BPE retains the same asymptotic complexity as Classical BPE, requiring only some modest additional bookkeeping. The need for a parallel multilingual corpus can at first seem prohibitive, but several pragmatic design choices can reduce the burden of this requirement. A small, aligned dataset suffices to drive the max-min decision in eq. 6, and the training dataset need not have this level of annotation. In addition, automatic language ID tools or script heuristics can help provide language labels when none are readily available. Also note that only the BPE learning phase differs; there is no algorithmic change to the tokenization function itself.

3.3 Algorithmic Variants

Preliminary experiments have shown several challenges with parity-aware BPE, which we address by introducing two variants.

Hybrid parity-aware BPE. Model developers may want to include and tokenize data for which parallel data is not available or where this concept does not even apply, such as programming code. Also, they may not want to guarantee full parity, but still give a high weight to global compression. We support these goals with a hybrid learning algorithm that uses the global objective of Classical BPE (eq. 4) for the first K merges, then switches to the parity-aware objective (eq. 5) for another J merges. K and J can be chosen by model developers to trade off global compression and fairness according to their priorities.

Moving-window balancing. There may be a point where the compression in a language no longer or barely improves, even if it is repeatedly chosen for the next merge. This could happen if the development dataset (the dataset used to choose the language) is too small or does not match the domain or language variant of the training dataset,⁵ or if $|\tau(\mathbf{b})|$ approaches the length of the pre-tokenized sequence. To prevent our algorithm from being “stuck” selecting the same language exclusively, we track the W most recent languages selected in eq. 6, and do not select a language if it occurs more than $\alpha \frac{W}{|Z|}$ times in this moving window.

4 Experimental Setup

We conduct experiments to evaluate the effectiveness of Parity-aware BPE, comparing it against baseline tokenization methods: Classical BPE and UnigramLM. All tokenizers are byte-level.

4.1 Tokenizer Training

Training Data. We train tokenizers using the multilingual C4 (mC4) corpus (Xue et al., 2021; Raffel et al., 2020).⁶ For choosing the focus language at each merge step when training Parity-aware BPE tokenizers (*i.e.*, computing eq. 6), we use the dev portion of FLORES+ (NLLB Team et al., 2024) as our multilingual development corpus—except for the *no-dev* systems, for which the training corpus is used to measure compression rate with bytes as normalization unit. To investigate how the number of languages, their linguistic diversity, and the variety of writing systems influence tokenizers, we consider two language sets: one with 30 languages (*30-lang*) and

another with 60 languages (*60-lang*). For each of these language sets, we create two dataset versions: one with uniform quantities of data per language (*balanced*) and one with per-language quantities proportional to amounts in the mC4 dataset (*unbalanced*). We present results for the *unbalanced* datasets here, as this is arguably the more realistic setting, with results for the *balanced* setting shown in §D. To enable tokenizer analyses as a function of different dataset qualities, we categorize the languages in each set based on the amount of training data available and the script family, performing some of our analyses by these categories. Languages with $> 1\text{M}$ examples are considered high-resource; those with 500k – 1M examples are medium-resource; and those with fewer than $< 500\text{k}$ examples are classified as low-resource. The full list of languages included in each set and the script family groupings are presented in Table 8.

Hyperparameter Settings. We look at vocabulary sizes $128k$ and $256k$. For *hybrid* systems, we learn half of the merges using the global strategy, and the second half using the parity-aware strategy. For systems with moving-window balancing (*window*), we use a window size of 100, and $\alpha = 2$. More information is provided in Appendix §C.

4.2 Evaluation

Our evaluations consist of task-independent tokenizer properties (intrinsic metrics) and downstream model performance (extrinsic metrics).

4.2.1 Intrinsic Metrics

We measure a variety of intrinsic tokenizer metrics on the devtest portion of FLORES+. All metrics are computed both globally and per-language to capture language-specific tokenization behavior. Normalization units differ across metrics, both in the effort to control for confounding factors and to tailor the metric to the tokenizer quality it is trying to measure. For example, morphologically motivated units can better reflect linguistic structure, and using character or document-level units can partially normalize the large differences in average bytes-per-character observed across writing systems (*e.g.*, Latin vs. UTF-8–encoded CJK⁷ scripts). For the sake of space, we provide brief metric descriptions here; more detailed definitions and formulae for all metrics can be found in §B.

⁵Kreutzer et al. (2022) discuss possible quality issues such as wrong or ambiguous language codes.

⁶<https://huggingface.co/datasets/allenai/c4>

⁷Chinese, Japanese and Korean scripts.

- **Fertility** measures the average number of tokens produced per normalization unit by a tokenizer; whitespace-delimited words are often the unit of interest (and are the units used in our computations). In this case, fertility quantifies how many tokens (on average) a word is broken up into.
- **Compression Rate (CR)** (as defined in §2) is a measure of the degree to which a unit of text has been shrunk after applying the given tokenizer (higher is better). We use documents as the normalization unit.
- **Vocabulary Utilization** is the fraction of the tokenizer’s vocabulary that actually appears in the evaluation corpus. Low utilization for a language signals wasted capacity or—when there are large differences across languages—biased vocabulary allocation.
- **Tokenizer Fairness Gini** (Meister, 2025) adapts the Gini coefficient to the per-language tokenization cost distribution (*e.g.*, tokens per line (document) in a parallel corpus). Values near 0 mean equal cost across languages; values closer to 1 indicate inequality.
- **MorphScore** (Arnett et al., 2025) measures how well token boundaries align with true morpheme boundaries, computed as morpheme-level precision/recall (and F1). High scores mean tokens respect morphological structure; low precision implies over-segmentation, while low recall may suggest under-segmentation.

For completeness, we also track Type–Token Ratio and Average Token Rank (vocabulary diversity; Limisiewicz et al., 2023) as well as Rényi entropies (distributional concentration; Zouhar et al., 2023a), evaluated using the TokEval suite (Meister, 2025).

4.2.2 Extrinsic Metrics.

For extrinsic evaluation, we train models using different tokenizers and assess their performance across a range of downstream tasks.

Model Architecture and Pretraining Data. We train decoder-only Transformer models (Vaswani et al., 2017) following the LLaMA architecture (Touvron et al., 2023) with 3 billion (3B) parameters. Full details on model configurations and training parameters are provided in §F. Models are trained on the FineWeb2 corpus (Penedo et al., 2025). We adopt temperature sampling with $\tau = 3.3$, following recommendations from prior work (Raffel et al., 2020; Conneau et al., 2020). We use the total 100B tokens to train each model.

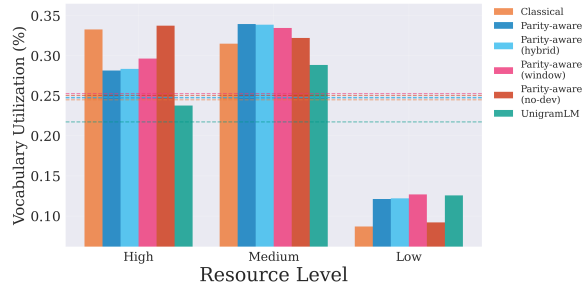


Figure 1: Vocabulary utilization for 128k tokenizers on the *unbalanced 30-lang* grouped by dataset resource levels. — lines indicate macro-averages across groups.

Benchmarks. We evaluate the models using perplexity on a held-out validation set from the respective pretraining datasets. In addition, we assess downstream performance on a suite of multilingual benchmarks. Results are aggregated per language to produce a score for each model-language pair. A full list of benchmarks and aggregation procedures is provided in the §G.

5 Results and Analysis

We present and analyze results on both intrinsic and extrinsic metrics. Within a set of comparisons, we fix the data distribution (*balanced* or *unbalanced*) and vocabulary size (128k or 256k).

5.1 Intrinsic Evaluation

Results in Table 1 show that all variants of Parity-aware BPE outperform Classical BPE in terms of the Gini coefficient, indicating more equitable token costs per document across languages. Among these variants, the base Parity-aware BPE emerges as the “fairest” tokenizer. Notably, Classical BPE and the Parity-aware BPE variants attain almost identical compression and Rényi entropies; we take this as evidence that the parity-aware variants match global efficiency while redistributing it more evenly. In addition, Parity-aware variants reduce fertility and increase MorphScore and vocabulary utilization, signaling better alignment with morphological boundaries and fairer vocabulary allocation. As could perhaps be expected, the no-dev variant of Parity-aware BPE performs most similarly to Classical BPE, closing only part of the Gini gap but matching it almost exactly on every other metric. This observation demonstrates the importance of thoughtful consideration of normalization units for the effectiveness of the algorithm; future work could address this by explicitly compensating for differences in cross-language length statistic, *e.g.*, by introducing a per-language

Tokenizer	Type-Token Ratio	Fertility	Compression Rate	Rényi Entropy ($\alpha=2.5$)	Gini Coefficient	MorphScore Precision	MorphScore Recall
Classical	0.0743	4.260 \pm 0.049	0.0303 \pm 0.0001	8.13	0.064	0.412 \pm 0.051	0.456 \pm 0.049
UnigramLM	0.0475	4.612 \pm 0.042	0.0228 \pm 0.0001	4.68	0.094	0.153 \pm 0.037	0.268 \pm 0.053
Parity-aware	0.0765	4.204 \pm 0.049	0.0300 \pm 0.0001	8.12	0.011	0.407 \pm 0.051	0.457 \pm 0.049
Parity-aware (hybrid)	0.0770	4.191 \pm 0.049	0.0303 \pm 0.0001	8.10	0.018	0.412 \pm 0.051	0.457 \pm 0.049
Parity-aware (window)	0.0788	4.219 \pm 0.050	0.0302 \pm 0.0001	8.11	0.013	0.405 \pm 0.049	0.453 \pm 0.047
Parity-aware (window+hybrid)	0.0794	4.203 \pm 0.050	0.0305 \pm 0.0001	8.09	0.022	0.416 \pm 0.049	0.460 \pm 0.047
Parity-aware (no-dev)	0.0772	4.310 \pm 0.050	0.0303 \pm 0.0001	8.12	0.059	0.423 \pm 0.051	0.466 \pm 0.049

Table 1: Intrinsic evaluation of 128k tokenizers on the (*unbalanced*) 30-lang dataset. Values are global statistics across the parallel corpus, except for MorphScore, which is macro-averaged across available languages.

Language	Classical BPE	Parity-aware (hybrid)	Parity-aware (window+hybrid)	Random
Arabic	38.19 \pm 2.90	39.04 \pm 2.89	38.84 \pm 2.90	32.00
Bengali	24.95 \pm 3.09	23.54 \pm 2.98	23.91 \pm 3.01	25.00
German	32.92 \pm 3.14	34.78 \pm 3.66	36.82 \pm 4.04	30.62
Greek	41.95 \pm 3.18	42.55 \pm 3.22	43.16 \pm 3.25	37.50
Spanish	37.53 \pm 2.66	38.83 \pm 2.71	39.30 \pm 2.75	32.77
Persian	42.80 \pm 5.39	39.15 \pm 5.27	39.15 \pm 5.27	25.00
French	38.67 \pm 3.90	36.59 \pm 2.84	37.10 \pm 2.82	32.00
Hindi	33.92 \pm 2.25	33.92 \pm 2.24	33.86 \pm 2.24	30.62
Indonesian	38.95 \pm 2.62	40.55 \pm 2.66	40.46 \pm 2.66	35.00
Italian	32.82 \pm 2.86	35.01 \pm 3.00	34.62 \pm 2.98	27.22
Japanese	37.43 \pm 2.39	37.45 \pm 2.39	37.43 \pm 2.39	34.00
Korean	33.00 \pm 5.22	33.00 \pm 5.22	34.33 \pm 5.29	25.00
Polish	29.75 \pm 2.50	31.14 \pm 2.60	28.97 \pm 2.49	23.75
Portuguese	33.63 \pm 2.81	33.15 \pm 2.77	33.06 \pm 2.77	27.50
Russian	36.36 \pm 2.27	36.21 \pm 2.26	36.57 \pm 2.28	32.77
Tamil	31.32 \pm 2.81	32.25 \pm 2.90	32.19 \pm 2.90	31.25
Telugu	32.73 \pm 2.61	33.52 \pm 2.61	33.26 \pm 2.61	30.00
Turkish	39.04 \pm 2.89	38.46 \pm 2.83	37.89 \pm 2.75	35.00
Vietnamese	33.69 \pm 2.31	33.87 \pm 2.27	33.80 \pm 2.30	29.50
Chinese	38.43 \pm 2.11	38.58 \pm 2.11	38.32 \pm 2.10	35.00
English	43.04 \pm 1.84	44.15 \pm 1.85	43.74 \pm 1.85	35.50
Thai	40.76 \pm 1.62	40.96 \pm 1.63	41.06 \pm 1.63	37.50

Table 2: Average downstream performance (accuracy %) across 13 multilingual benchmarks (tokenizers trained on the (*unbalanced*) 30-lang dataset with 128k vocab size). The **Random** column shows the expected accuracy of a random classifier. Best performance per language is bolded. Benchmark details for each language are provided in §G and Table 9, respectively.

multiplicative factor.⁸ By contrast, the hybrid and window variants land between Parity-aware and Classical BPE on many metrics: they recover a small slice of global compression while reducing the Gini inequity coefficient of Classical BPE by roughly three-quarters and achieving the lowest fertility of all runs. Taken together, the variants outline a smooth fairness–efficiency frontier, allowing practitioners to select the point that best suits their resource constraints and fairness targets.

Fig. 1 presents vocabulary utilization grouped by resource tier, as defined in §4.1. For high-resource languages, the results indicate that parity-aware BPE (no-dev) performs comparably to Classical BPE in terms of vocabulary utilization, while the other variants provide worse vocabulary utilization. In contrast, for low- and medium-resource languages, the hybrid and window variants achieve higher vocabulary utilization, highlighting their

⁸For example, without a multi-parallel corpus, developers can estimate the desired compression rate from a number of parallel corpora.

effectiveness at providing fairer vocabulary allocation across languages. We show per-language results for compression rate and vocabulary utilization in Fig. 2 and 3. We observe that Parity-aware tokenizers attain substantially more uniform compression across languages. While they also generally achieve higher vocabulary utilization rates across most languages, the level of benefit varies by language.

Vocabulary Size. Repeating the experiment with a 256k vocabulary (Table 7 and Fig. 4) yields the same conclusion: parity-aware BPE tokenizers consistently outperform Classical BPE on cross-lingual fairness metrics, with minimal impact on performance-oriented measures.

5.2 Extrinsic Evaluation

Table 2 presents the performance of LMs trained with three different 128k tokenizers: Classical BPE, Hybrid Parity-aware BPE, and Hybrid Parity-aware BPE (moving-window), evaluated on the 30-lang set. For each language, we report mean performance, standard errors, and a random baseline to account for varying benchmark counts (Table 9). The evaluation assesses Parity-aware BPE’s impact on downstream performance, particularly in languages where Classical BPE is efficient. Results indicate that Parity-aware BPE maintains performance across languages: accuracy changes relative to Classical BPE are small. Models trained with the hybrid variant show a median per-language change in accuracy of +0.19 percentage points, with 14 languages improving and 6 declining; the window+hybrid leads to very similar changes in accuracy. These results confirm that parity-aware tokenizers can handle diverse languages without compromising LM performance. We show per-language perplexity results in Fig. 6 in §D. Here, we see that models trained with parity-aware tokenizers show much more uniform perplexity across languages, whereas Classical BPE yields a handful of languages with markedly higher perplexity.

6 Related Work

Multilingual Tokenization. Despite their popularity, BPE and related subword tokenization methods often underperform in multilingual settings due to limited handling of spelling variation and morphological complexity (Bostrom and Durrett, 2020). A tokenizer’s tokenization parity and fertility directly affect computational cost and model performance. Prior work has explored vocabulary allocation strategies: Zhang et al. (2022) show that larger vocabularies improve NMT robustness across scripts, while Gowda and May (2020) demonstrate that tuning BPE merges can mitigate sequence length issues. Rust et al. (2021) find that integrating specialized monolingual tokenizers into multilingual systems can boost performance; however, recent results indicate that optimal vocabulary size depends on the task and model (Dagan et al., 2024). For multilingual vocabulary construction, Chung et al. (2020) explore clustering-based sharing of subword units across languages, and Limisiewicz et al. (2023) propose an explicit tokenizer-merging algorithm to combine per-language vocabularies. Finally, tokenization-free models such as CANINE (Clark et al., 2022) and ByT5 (Xue et al., 2022) offer an alternative approach for improved multilingual handling.

Tokenization Bias and Recent Advances. Recent research highlights biases from tokenization in LLMs. While Wan (2022) argues that character- and byte-level representations are intrinsically fair,⁹ other studies (Petrov et al., 2023; Ahia et al., 2023) show that tokenization differences across languages; even at character and byte levels; affect costs, latency, and contextual understanding. This has spurred efforts like Aya (Aryabumi et al., 2024) and methods to mitigate tokenization unfairness (Fujii et al., 2024; Abboud and Oz, 2024; Limisiewicz et al., 2024). Although newer character- and byte-level models use compression techniques such as entropy-based patching (Pagnoni et al., 2025), cross-lingual parity of these representations remains unstudied.

7 Discussion and Conclusion

Tokenizers optimized using standard algorithms and data can lead to disparities in users’ costs and experiences as a result of their choice of language.

⁹They report more random performance with these representations, which we do not view as “fairer.”

Parity-optimized tokenization can remedy this by explicitly balancing compression across languages, enabling fairer treatment of users of low-resource languages. Parity-aware BPE implements this idea: it was designed to improve cross-lingual tokenization parity, and our experiments confirm that it does. On the unbalanced 30-language set, the Gini coefficient of per-line token costs falls from 0.064 with Classical BPE to 0.011 with our parity-aware variant while compression ratios for most variants stay competitive with Classical BPE, often improving when looking at averages across languages on the whole. Crucially, this fairness gain does not come at the expense of downstream quality: across 13 multilingual benchmarks, models trained with parity-aware tokenizers either outperform or stay within a single standard error of the Classical BPE baseline for every tested language (Table 2).

Overall, the trade-offs required for using parity-aware BPE are minimal. From the model-developer’s perspective, parity-aware BPE is a drop-in replacement: it requires no architectural changes and minimal changes to the tokenizer pipeline. Concretely, inference remains exactly the same as in Classical BPE. During the learning stage, the algorithm adds only an $\mathcal{O}(\mathcal{L})$ pass per merge for recomputing language-level compression rates on a dev corpus, leaving the asymptotic complexity identical to Classical BPE. Further, we observe empirically that a small, sentence-aligned development set is sufficient to drive the “fair-max” decision (§C.3). When resource or domain mismatches make full equality undesirable, hybrid and moving-window variants further let practitioners trade off global compression versus strict parity; our empirical results validate that these variants perform well in practice.

Making the tokenization step of the NLP pipeline more equitable is therefore not just desirable but feasible. Parity-aware BPE offers a clear avenue towards this goal by building fairness into the tokenizer itself. With a simple modification—selecting merges that benefit the worst-compressed language—it substantially reduces the hidden “token tax” imposed on speakers of low-resource languages without sacrificing compression or task accuracy. Future work can push this agenda further by extending parity objectives to alternative tokenization schemes and other modalities such as speech and vision, as well as by developing benchmarks and metrics for fairness assessments in tokenization beyond compression parity.

696 Limitations

697 Our study uses parallel corpora to estimate
698 per-language costs; in domains where aligned docu-
699 ments are unavailable or difficult to obtain, using
700 unaligned corpora and alternative normalization
701 units for making the language choice may intro-
702 duce bias. While we consider 60 languages and
703 two vocabulary sizes, the interplay between tok-
704 enization parity and model scaling still needs to be
705 explored for much larger models and for code or
706 multimodal inputs. Finally, fairness here is defined
707 purely in terms of token counts. While we measure
708 other potential quantification of fairness (*e.g.*, mor-
709 phological alignment), there are still other notions
710 that are unaccounted for. We leave optimization for
711 these metrics during tokenizer learning to future
712 work.

713 References

714 Khadige Abboud and Gokmen Oz. 2024. Towards eq-
715 uitable natural language understanding systems for
716 dialectal cohorts: Debiasing training data. In *Pro-
717 ceedings of the 2024 Joint International Conference
718 on Computational Linguistics, Language Resources
719 and Evaluation, LREC/COLING 2024, 20-25 May,
720 2024, Torino, Italy*, pages 16487–16499. ELRA and
721 ICCL.

722 Orevaoghene Ahia, Sachin Kumar, Hila Gonen, Jungo
723 Kasai, David R. Mortensen, Noah A. Smith, and
724 Yulia Tsvetkov. 2023. Do all languages cost the
725 same? tokenization in the era of commercial lan-
726 guage models. In *Proceedings of the 2023 Confer-
727 ence on Empirical Methods in Natural Language Pro-
728 cessing, EMNLP 2023, Singapore, December 6-10,
729 2023*, pages 9904–9923. Association for Computa-
730 tional Linguistics.

731 Mehdi Ali, Michael Fromm, Klaudia Thellmann,
732 Richard Rutmann, Max Lübbering, Johannes Lev-
733 eling, Katrin Klug, Jan Ebert, Niclas Doll, Jasper
734 Buschhoff, Charvi Jain, Alexander Weber, Lena Ju-
735 rkschat, Hammam Abdelwahab, Chelsea John, Pedro
736 Ortiz Suarez, Malte Ostendorff, Samuel Weinbach,
737 Rafet Sifa, Stefan Kesselheim, and Nicolas Flores-
738 Herr. 2024. Tokenizer choice for LLM training: Neg-
739 ligible or crucial? In *Findings of the Association
740 for Computational Linguistics: NAACL 2024*, pages
741 3907–3924, Mexico City, Mexico. Association for
742 Computational Linguistics.

743 Catherine Arnett, Marisa Hudspeth, and Brendan
744 O’Connor. 2025. Evaluating Morphological Align-
745 ment of Tokenizers in 70 Languages. In *Proceedings
746 of the ICML 2025 Tokenization Workshop (TokShop)*.

747 Viraat Aryabumi, John Dang, Dwarak Talupuru,
748 Saurabh Dash, David Cairuz, Hangyu Lin, Bharat

Venkitesh, Madeline Smith, Kelly Marchisio, Se-
bastian Ruder, et al. 2024. Aya 23: Open weight
releases to further multilingual progress. *CoRR*,
abs/2405.15032. 749
750
751
752

Lucas Bandarkar, Davis Liang, Benjamin Muller, Mikel
Artetxe, Satya Narayan Shukla, Donald Husa, Naman
Goyal, Abhinandan Krishnan, Luke Zettlemoyer, and
Madian Khabsa. 2024. The Belebele benchmark: a
parallel reading comprehension dataset in 122 lan-
guage variants. In *Proceedings of the 62nd Annual
Meeting of the Association for Computational Lin-
guistics (Volume 1: Long Papers)*, pages 749–775,
Bangkok, Thailand. Association for Computational
Linguistics. 753
754
755
756
757
758
759
760
761
762

Kaj Bostrom and Greg Durrett. 2020. Byte pair encod-
ing is suboptimal for language model pretraining. In
*Findings of the Association for Computational Lin-
guistics: EMNLP 2020*, pages 4617–4624, Online.
Association for Computational Linguistics. 763
764
765
766
767

Michael Chen, Mike D’Arcy, Alisa Liu, Jared Fer-
nandez, and Doug Downey. 2019. CODAH: An
adversarially-authored question answering dataset
for common sense. In *Proceedings of the 3rd Work-
shop on Evaluating Vector Space Representations for
NLP*, pages 63–69. Association for Computational
Linguistics. 768
769
770
771
772
773
774

Hyung Won Chung, Dan Garrette, Kiat Chuan Tan, and
Jason Riesa. 2020. Improving multilingual models
with language-clustered vocabularies. In *Proced-
ings of the 2020 Conference on Empirical Methods in
Natural Language Processing, EMNLP 2020, Online,
November 16-20, 2020*, pages 4536–4546. Associa-
tion for Computational Linguistics. 775
776
777
778
779
780
781

Jonathan H. Clark, Dan Garrette, Iulia Turc, and John
Wieting. 2022. CANINE: Pre-training an efficient
tokenization-free encoder for language representa-
tion. *Trans. Assoc. Comput. Linguistics*, 10:73–91. 782
783
784
785

Alexis Conneau, Kartikay Khandelwal, Naman Goyal,
Vishrav Chaudhary, Guillaume Wenzek, Francisco
Guzmán, Edouard Grave, Myle Ott, Luke Zettle-
moyer, and Veselin Stoyanov. 2020. Unsupervised
cross-lingual representation learning at scale. In *Pro-
ceedings of the 58th Annual Meeting of the Associa-
tion for Computational Linguistics*, pages 8440–8451.
Association for Computational Linguistics. 786
787
788
789
790
791
792
793

Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina
Williams, Samuel Bowman, Holger Schwenk, and
Veselin Stoyanov. 2018. XNLI: Evaluating cross-
lingual sentence representations. In *Proceedings of
the 2018 Conference on Empirical Methods in Nat-
ural Language Processing*, pages 2475–2485, Brus-
sels, Belgium. Association for Computational Lin-
guistics. 794
795
796
797
798
799
800
801

Gautier Dagan, Gabriel Synnaeve, and Baptiste Roz-
ière. 2024. Getting the most out of your tokenizer
for pre-training and domain adaptation. In *Proced-
ings of the 41st International Conference on Machine
Learning, ICML’24*. JMLR.org. 802
803
804
805
806

A Pseudocode

Algorithm 1: Algorithm for learning \mathbf{m} using Classical BPE.

Input: Corpus \mathcal{D} ; number of merges K
Output: Vocabulary \mathcal{V}_K ; merge sequence \mathbf{m}_K

```

1  $\mathcal{V}_0 \leftarrow \mathcal{B}$ 
2  $\mathbf{m}_0 \leftarrow \langle \rangle$ 
3 for  $k \leftarrow 1$  to  $K$  do
    // Count all adjacent token pairs
4    $Pairs \leftarrow \{ \}$ 
5   foreach occurrence of consecutive
     tokens  $v v'$  in  $\mathcal{D}$  where  $v, v' \in \mathcal{V}_{k-1}$  do
6      $\lfloor Pairs[(v, v')] \leftarrow Pairs[(v, v')] + 1$ 
7    $(v^*, v'^*) \leftarrow \arg \max_{(v, v')} Pairs[(v, v')]$ 
8    $w^* \leftarrow v^* \circ v'^*$ 
   // Update vocabulary and merge
   sequence
9    $\mathcal{V}_k \leftarrow \mathcal{V}_{k-1} \cup \{w^*\}$ 
10   $\mathbf{m}_k \leftarrow \mathbf{m}_{k-1} + \langle (v^*, v'^*) \rangle$ 
   // Replace all occurrences in corpus
11  foreach occurrence of  $v^* v'^*$  in  $\mathcal{D}$  do
12     $\lfloor$  Replace  $v^* v'^*$  with  $w^*$ 
13 return  $\mathcal{V}_K, \mathbf{m}_K$ 

```

B Intrinsic Tokenizer Evaluation Metrics

We provide detailed descriptions of the intrinsic tokenizer metrics used in §5, grouped by the general tokenizer characteristic the metric aims to assess. Metric formulae are defined in terms of our definition of a tokenizer $T = (\mathcal{V}, \tau, \perp)$ given in §2. For this tokenizer, we denote the empirical unigram frequency distribution of tokens $v \in \mathcal{V}$ as X_T , which is computed on our evaluation corpus.

B.1 Vocabulary Usage

Vocabulary Utilization and Type-Token Ratio. Vocabulary utilization measures the proportion of a tokenizer’s full vocabulary that is actively used when processing a given corpus. For tokenizer T on corpus \mathcal{D} , we compute it as:

$$\text{VocabUtil}(T) = \frac{|\{v : v \in \tau(\mathbf{b}), \mathbf{b} \in \mathcal{D}\}|}{|\mathcal{V}|} \quad (7)$$

Here, the numerator counts the number of distinct tokens observed across the tokenization of

Algorithm 2: Algorithm for learning \mathbf{m} using Parity-aware Byte Pair Encoding with separate training and development sets.

Input: $\{\mathcal{D}_\ell\}_{\ell \in \mathcal{L}}$ (multilingual training corpus);
 $\{\mathcal{D}_\ell^{\text{dev}}\}_{\ell \in \mathcal{L}}$ (multilingual development corpus);
 K (number of merges)

Output: \mathcal{V}_K (vocabulary); \mathbf{m}_K (merge list)

```

1  $\mathcal{V}_0 \leftarrow \mathcal{B}; \mathbf{m}_0 \leftarrow \langle \rangle$ 
2 for  $k \leftarrow 1$  to  $K$  do
    // Calculate compression rate for each
    language
3   foreach language  $\ell \in \mathcal{L}$  do
4      $\left[ \begin{array}{l} \text{CR}(\mathcal{D}_\ell^{\text{dev}}, \tau_{\mathbf{m}_{<k}}) \leftarrow \\ \frac{\sum_{\mathbf{b} \in \mathcal{D}_\ell^{\text{dev}}} |\mathbf{b}|_u}{\sum_{\mathbf{b} \in \mathcal{D}_\ell^{\text{dev}}} |\tau_{\mathbf{m}_{<k}}(\mathbf{b})|} \end{array} \right.$ 
5    $\ell^* \leftarrow \arg \min_{\ell \in \mathcal{L}} \text{CR}(\mathcal{D}_\ell^{\text{dev}}, \tau_{\mathbf{m}_{<k}})$ 
   // Consider token pairs only in  $\mathcal{D}_{\ell^*}$ 
6    $Pairs \leftarrow \{ \}$ 
7   foreach occurrence of consecutive
     tokens  $v v'$  in  $\mathcal{D}_{\ell^*}$  where  $v, v' \in \mathcal{V}_{k-1}$  do
8      $\lfloor Pairs[(v, v')] \leftarrow Pairs[(v, v')] + 1$ 
9    $(v^*, v'^*) \leftarrow \arg \max_{(v, v')} Pairs[(v, v')]$ 
10   $w^* \leftarrow v^* \circ v'^*$ 
   // Update vocabulary and merge list
11   $\mathcal{V}_k \leftarrow \mathcal{V}_{k-1} \cup \{w^*\}$ 
12   $\mathbf{m}_k \leftarrow \mathbf{m}_{<k} + \langle (v^*, v'^*) \rangle$ 
   // Apply merge across all languages
13  foreach language  $\ell \in \mathcal{L}$  do
14    foreach occurrence of  $v^* v'^*$  in  $\mathcal{D}_\ell$ 
      and  $\mathcal{D}_\ell^{\text{dev}}$  do
15       $\lfloor$  Replace  $v^* v'^*$  with  $w^*$ 
16 return  $\mathcal{V}_K, \mathbf{m}_K$ 

```

all strings in the corpus. The type-token ratio quantifies lexical diversity by measuring the proportion of unique tokens (types) relative to the total number of tokens produced by a tokenizer:

$$\text{TTR}(T) = \frac{|\{v : v \in \tau(\mathbf{b}), \mathbf{b} \in \mathcal{D}\}|}{\sum_{\mathbf{b} \in \mathcal{D}} |\tau(\mathbf{b})|} \quad (8)$$

where $|\tau(\mathbf{b})|$ is the number of tokens produced by tokenizer T for input \mathbf{b} . In words, the numerator counts distinct token types and the denominator counts total tokens across the corpus.

High vocabulary utilization and type-token ratio indicate efficient use of the learned vocabulary; low

values of these metrics for a particular language may suggest tokenizer bias, as only a small portion of the tokenizer’s vocabulary is used/applicable for that language.

Average Token Rank. Average token rank (Limisiewicz et al., 2023) measures the typical position of tokens in a tokenized text within the frequency-ordered vocabulary. In more detail, we compute the rank of each token (denoted as $\text{rank}(v)$) in our unigram frequency distribution X_T ; rank 1 corresponds to the most frequent token. We compute average token rank across tokens in the evaluation corpus as:

$$\text{AvgRank}(T) = \frac{\sum_{\mathbf{b} \in \mathcal{D}} \sum_{v \in \tau(\mathbf{b})} \text{rank}(v)}{\sum_{\mathbf{b} \in \mathcal{D}} |\tau(\mathbf{b})|} \quad (9)$$

This metric can be seen as another measure of the proportion of the vocabulary used by a tokenizer. Lower average ranks indicate that the tokenizer predominantly uses a small set of tokens, while higher averages suggest more diverse token usage, including rare vocabulary items. When computed per language (*i.e.*, when ranks are computed using the language’s respective frequency distribution), systematic differences in average token rank across languages reveal vocabulary allocation bias.

B.2 Information-theoretic Metrics

Compression Rate. We evaluate compression rate—as defined in eq. 1—across a parallel corpus. As discussed in §3.2, this enables us to use lines (documents) as our normalization unit. Recall that higher compression rates are generally desirable for computational efficiency in downstream tasks. In multilingual corpora, compression ratio disparities across languages indicate systematic tokenizer bias, where certain languages achieve better compression efficiency than others, potentially leading to unequal computational costs.

Rényi Entropy. We compute Rényi entropy of order α over the empirical unigram frequency distribution X_T for a given tokenizer T to capture different aspects of token distribution:

$$H_\alpha(X_T) = \frac{1}{1-\alpha} \log_2 \left(\sum_{v \in \mathcal{V}} p(v)^\alpha \right) \quad (10)$$

for $0 < \alpha < \infty$; at $\alpha = 1$, the definition of Shannon entropy is typically adopted. Rényi entropy provides a parametric family of measures that emphasize different aspects of the distribution: H_1

(Shannon entropy), H_2 (collision entropy), and H_∞ (min-entropy). Rényi efficiency is Rényi entropy normalized by the size of the support, which is helpful for comparing tokenizers with different vocabulary sizes (Zouhar et al., 2023a). As all of our comparisons are between tokenizers of the same vocabulary size, we omit this normalization step and compare entropies directly.

B.3 Morphological and Multilingual Fairness Metrics

Fertility. Fertility measures the average number of tokens produced per unit (word, character, or byte) by a tokenizer; the unit of interest for fertility is often the *word*, in which case, fertility quantifies how many tokens (on average) a word is broken up into. We use words as our normalization unit in our computations, as determined by the HuggingFace Whitespace Pretokenizer. We formally define tokenizer fertility for a given corpus \mathcal{D} as:

$$\text{Fertility}(T) = \frac{\sum_{\mathbf{b} \in \mathcal{D}} |\tau(\mathbf{b})|}{\sum_{\mathbf{b} \in \mathcal{D}} |\mathbf{b}|_u} \quad (11)$$

This metric can give a sense for the computational efficiency imbued by a tokenizer, as well as for sequence length estimates for downstream modeling tasks.

MorphScore. MorphScore (Arnett et al., 2025) evaluates tokenizer quality through morpheme-level precision and recall, measuring how well tokenizers preserve morphological information during segmentation. We point the reader to the original work. Differences in cross-language MorphScore reveal how consistently a tokenizer’s sub-token boundaries align with true morpheme boundaries. A higher score in one language than another indicates that the tokenizer preserves that language’s morphological structure more faithfully. MorphScore provides a notion of both precision and recall (we point the reader to the original work for the exact description of the computation). Low precision indicates tokenizer oversegmentation; low recall is suggestive of *under* segmentation.

Tokenizer Fairness Gini Coefficient. We use an adaptation of the Gini coefficient—often used as a measure of economic inequality—to encapsulate tokenizer fairness across languages (Meister, 2025). Formally, let $c_1 \leq c_2 \leq \dots \leq c_n$ be the “costs” under a given tokenizer T for languages

$\mathcal{L} = \{l_1, l_2, \dots, l_n\}$. Here, we quantify cost as the average number of tokens it takes to encode the unit of interest (*e.g.*, a byte, word or line);¹⁰ when using a parallel corpus, this can be cost per line (document), which controls for discrepancies between average character byte lengths across different scripts. The Gini coefficient for tokenizer T is then:

$$\text{Gini}(T) = \frac{1}{n} \left(n + 1 - 2 \frac{\sum_{i=1}^n (n + 1 - i)c_i}{\sum_{i=1}^n c_i} \right) \quad (12)$$

Values range from 0 (completely equal costs across languages) to 1 (maximum inequality). This metric condenses multilingual tokenizer fairness into a single number by measuring the degree of inequality in computational costs across languages; lower Gini coefficients indicate more equitable tokenizer compression across languages, while higher values suggest systematic bias toward certain languages.

C Hyperparameter Selection

C.1 Hybrid parity-aware BPE

For the hybrid parity-aware tokenizer, we set the number of parity-aware merges K equal to the number of classical BPE merges J (*i.e.*, $K = J$, with each set to either $64k$ or $128k$). This choice represents a midpoint between classical BPE ($J > 0, K = 0$) and fully parity-aware BPE ($J = 0, K > 0$). Alternative settings allow practitioners to trade off efficiency on high-resource languages (*e.g.*, English) against cross-lingual fairness. Because no single operating point is objectively optimal across use cases, we do not claim a universally optimal choice for this trade-off.

C.2 Moving-Window Balancing

The moving-window balancing mechanism is designed to prevent degenerate cases in which a single language repeatedly consumes merge operations despite diminishing compression gains. The window size parameter ranges from 1 (enforcing uniform language selection across merges) to the number of languages N (equivalent to no constraint). In a 30-language setting, this implies that even if compression improvements for a given language stagnate, it can consume at most approximately $1/15$ of merge operations. Increasing the window size causes the behavior to increasingly resemble the unbalanced setting.

¹⁰This is equivalent to fertility, or the inverse of the compression rate.

Because this mechanism primarily guards against pathological merge allocation rather than directly optimizing compression, we do not expect the window size to be a sensitive hyperparameter.

To validate this intuition, we conduct ablation experiments varying the window size in the moving-window balancing variant. Table 3 reports intrinsic tokenizer metrics for window sizes ranging from 50 to 200. Average compression rate, morphological plausibility, and other intrinsic metrics remain largely stable across window sizes. Empirically, the windowing mechanism prevents repeated selection of the same language caused by development-training set mismatch, without introducing sensitivity to the specific window size. These results suggest that the moving-window approach can be applied robustly without careful tuning.

C.3 Parallel Development Set Size

We further study the effect of the size of the parallel development set used in cross-lingual compression rate comparisons for language selection (Eq. 7). Table 4 reports results for development set sizes of 100, 300, and 1000 examples.

The results indicate that relatively small development sets suffice to capture most of the fairness improvements. A parallel development set of 100 examples achieves a 67% reduction in the Gini Inequality Coefficient relative to the classical baseline, compared to a 72% reduction when using the full 1000-example development set. Global compression rates remain essentially unchanged across settings, suggesting that incorporating a parallel development set is feasible without imposing substantial data or computational burdens.

D Additional Results and Ablation Studies

In this section, we present the results of our ablation studies. Table 5 reports the intrinsic evaluation of tokenizers with a $128k$ vocabulary size on the (unbalanced) *60-lang* dataset. Table 6 shows the corresponding results for the (balanced) *30-lang* dataset, also with a $128k$ vocabulary size. Finally, Table 7 presents the intrinsic evaluation of tokenizers with a $256k$ vocabulary size on the (*unbalanced*) *30-lang* dataset. Together, these results demonstrate the effectiveness of Parity-aware BPE across different language settings, vocabulary sizes, and data distributions.

Compression Ratio by Tokenizer

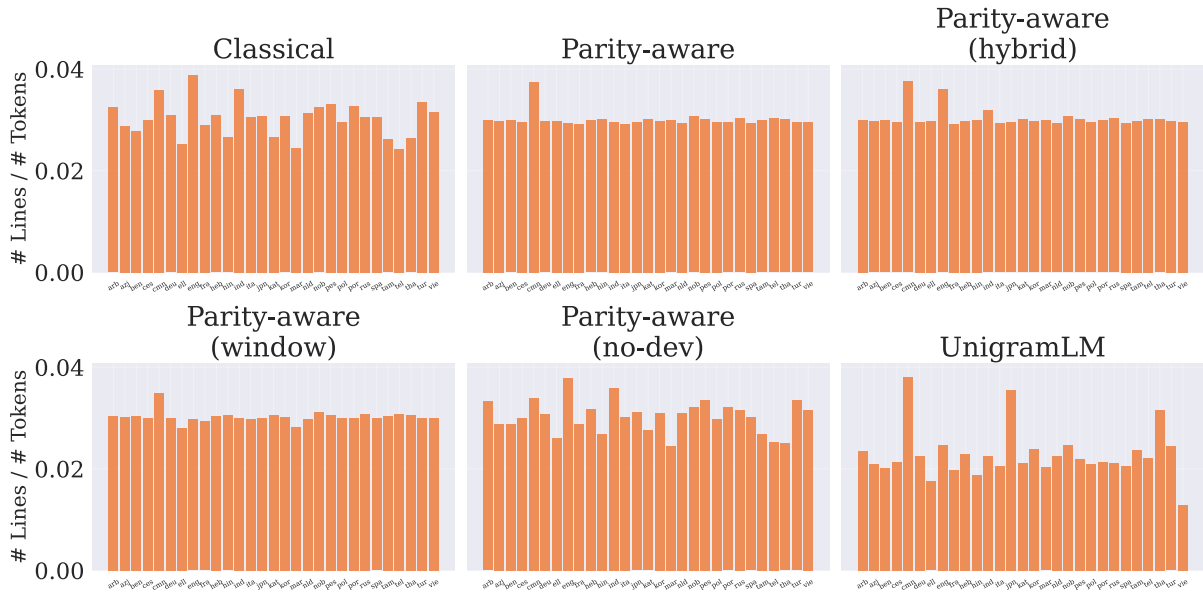


Figure 2: Compression rate of 128k tokenizers on the (*unbalanced*) 30-lang per language.

Vocabulary Utilization by Tokenizer

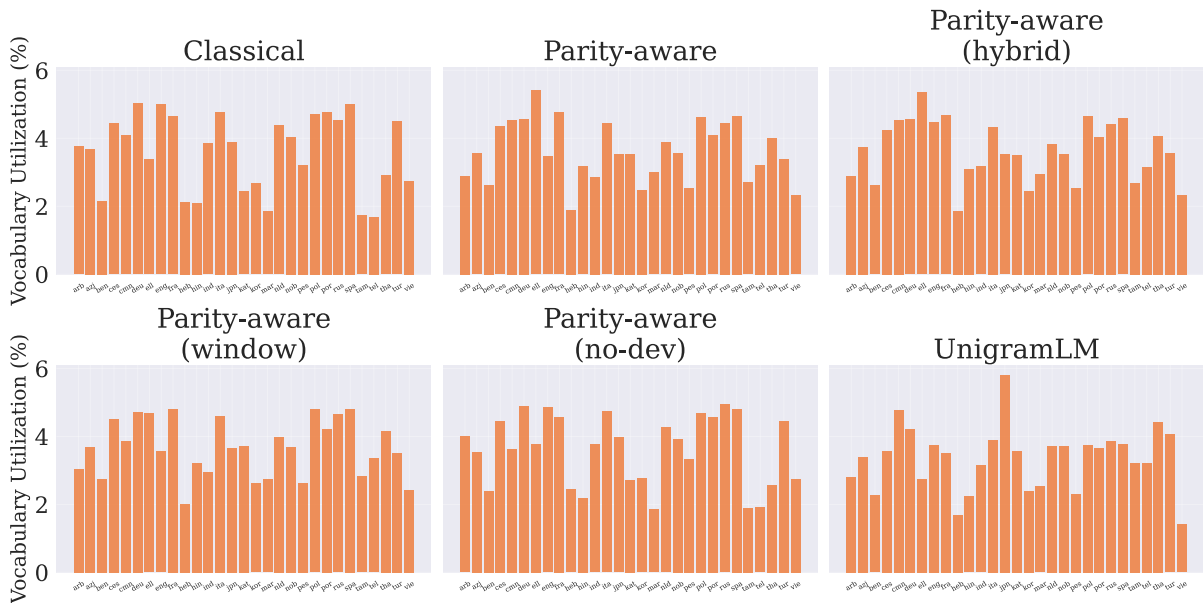


Figure 3: Vocabulary utilization of 128k tokenizers on the (*unbalanced*) 30-lang per language.

Tokenizer	Comp. Rate	TTR	Fertility	Rényi ($\alpha = 2.5$)	Gini	MorphScore
Classical	0.0283	0.0706	4.508	7.56	0.061	0.412
Parity-aware (hybrid)	0.0284	0.0731	4.439	7.53	0.017	0.412
Hybrid + window 50	0.0285	0.0751	4.444	7.52	0.020	0.412
Hybrid + window 100	0.0286	0.0754	4.451	7.52	0.021	0.416
Hybrid + window 150	0.0286	0.0755	4.453	7.52	0.022	0.418
Hybrid + window 200	0.0286	0.0754	4.451	7.52	0.021	0.416

Table 3: Intrinsic tokenizer evaluation metrics for ablations over the moving-window size. Results indicate minimal sensitivity to the window size parameter.

Tokenizer	Comp. Rate	TTR	Fertility	Rényi ($\alpha = 2.5$)	Gini	MorphScore
Classical	0.0283	0.0706	4.508	7.56	0.061	0.412
Hybrid (dev=100)	0.0283	0.0724	4.400	7.54	0.020	0.411
Hybrid (dev=300)	0.0283	0.0727	4.426	7.54	0.018	0.412
Hybrid (dev=1000)	0.0284	0.0731	4.439	7.53	0.017	0.412

Table 4: Intrinsic tokenizer evaluation metrics for ablations over the size of the parallel development set. Smaller development sets retain most fairness gains.

Training Data Distribution. To assess the sensitivity of the Parity-aware algorithm to training data distribution, we also analyze results for 128k tokenizers trained on the *balanced* version of the dataset. The results in Table 6 indicate that parity-aware BPE tokenizers perform similarly to Classical BPE across most metrics. However, in terms of fertility, Classical BPE outperforms the parity-aware variants. This suggests that parity-aware tokenizers are particularly beneficial in unbalanced settings, where low-resource languages are more disadvantaged, whereas in balanced scenarios their advantage diminishes. We also interestingly see in Fig. 4—again for the (*balanced*) 30-lang setting—that parity-aware tokenizers yield the largest absolute increases in vocabulary utilization for high-resource languages. Low- and medium-resource languages also improve, though to a smaller extent. One logical conclusion from this result is that the effect of parity-aware tokenizer is better described as balancing utilization across languages rather than directly compensating for data scarcity.

Script Analysis. Fig. 5 illustrates vocabulary utilization across languages for 128k tokenizers on the (unbalanced) 30-lang dataset. For Latin, Arabic, Hebrew, and Cyrillic scripts, the Parity-aware BPE (no-dev) variant outperforms all other parity-aware BPE versions, with Classical BPE ranking second. In contrast, for the CJK scripts, Classical BPE leads, while all parity-aware BPE variants perform similarly and closely follow. For the remaining scripts, the base Parity-aware BPE or the window-balanced variant significantly outperform other tokenizers. These findings suggest that dif-

ferent scripts benefit differently from parity-aware BPE approaches.

Language Model Perplexities. We report language model perplexities on the FineWeb2 validation set in Fig. 6. Results are shown per language. We see a noticeably larger cross-lingual spread in perplexity for language models trained using the Classical BPE tokenizer than for those trained using Parity-aware variants. The Parity-aware tokenizers seem to eliminate the long tail present under Classical BPE while maintaining comparable mean perplexity across languages. Note that we normalize by number of bytes in the text rather than by number of tokens to account for differences in tokenization lengths.

E Balancing the tokenizer training set

Prior work has addressed representational imbalance in subword tokenization primarily through corpus-level reweighting or balancing strategies, such as equalizing the number of documents or bytes per language prior to training a BPE or unigram tokenizer. While such approaches ensure that low-resource languages contribute to the learned vocabulary, they optimize the tokenizer with respect to an artificial distribution that differs from the natural data distribution encountered during downstream training and inference. This mismatch can result in inefficient allocation of vocabulary capacity, over-representation of rare patterns, and increased tokenization length for high-frequency languages. Moreover, corpus balancing constitutes a coarse-grained intervention: the choice of balancing criterion substantially affects outcomes, and the

Tokenizer	Type-Token Ratio	Fertility	Compression Rate	Rényi Entropy ($\alpha=2.5$)	Gini Coefficient	MorphScore Precision	MorphScore Recall
Classical	0.0388	3.374 \pm 0.027	0.0277 \pm 0.0000	8.16	0.086	0.324 \pm 0.031	0.407 \pm 0.029
Parity-aware	0.0321	3.533 \pm 0.029	0.0260 \pm 0.0000	8.25	0.022	0.273 \pm 0.030	0.379 \pm 0.028
Parity-aware (hybrid)	0.0334	3.453 \pm 0.028	0.0269 \pm 0.0000	8.20	0.040	0.283 \pm 0.030	0.379 \pm 0.028
Parity-aware (window)	0.0392	3.438 \pm 0.028	0.0270 \pm 0.0000	8.17	0.030	0.305 \pm 0.029	0.393 \pm 0.028
Parity-aware (window+hybrid)	0.0409	3.362 \pm 0.028	0.0278 \pm 0.0000	8.12	0.044	0.317 \pm 0.029	0.400 \pm 0.028
Parity-aware (no-dev)	0.0401	3.412 \pm 0.027	0.0276 \pm 0.0000	8.16	0.080	0.334 \pm 0.031	0.418 \pm 0.029

Table 5: Intrinsic evaluation of 128k tokenizers on the (*unbalanced*) 60-lang dataset. Values are global statistics, except for MorphScore, which is macro-averaged across available languages.

Tokenizer	Type-Token Ratio	Fertility	Compression Rate	Rényi Entropy ($\alpha=2.5$)	Gini Coefficient	MorphScore Precision	MorphScore Recall
Classical	0.0780	4.175 \pm 0.049	0.0307 \pm 0.0001	8.09	0.050	0.409 \pm 0.048	0.454 \pm 0.046
Parity-aware	0.0765	4.207 \pm 0.049	0.0300 \pm 0.0001	8.12	0.011	0.405 \pm 0.052	0.455 \pm 0.049
Parity-aware (hybrid)	0.0767	4.192 \pm 0.049	0.0303 \pm 0.0001	8.11	0.016	0.404 \pm 0.050	0.448 \pm 0.048
Parity-aware (window)	0.0787	4.222 \pm 0.050	0.0302 \pm 0.0001	8.11	0.013	0.407 \pm 0.050	0.455 \pm 0.047
Parity-aware (window+hybrid)	0.0794	4.177 \pm 0.049	0.0305 \pm 0.0001	8.09	0.020	0.413 \pm 0.049	0.457 \pm 0.047
Parity-aware (no-dev)	0.0802	4.234 \pm 0.050	0.0306 \pm 0.0001	8.09	0.047	0.418 \pm 0.048	0.463 \pm 0.046
Parity-aware (hybrid+no-dev)	0.0800	4.231 \pm 0.050	0.0307 \pm 0.0001	8.08	0.048	0.415 \pm 0.048	0.458 \pm 0.046

Table 6: Intrinsic evaluation of 128k tokenizers on the (*balanced*) 30-lang dataset. Values are global statistics, except for MorphScore, which is macro-averaged across available languages.

Tokenizer	Type-Token Ratio	Fertility	Compression Rate	Rényi Entropy ($\alpha=2.5$)	Gini Coefficient	MorphScore Precision	MorphScore Recall
Classical	0.1239	3.767 \pm 0.044	0.0340 \pm 0.0001	7.85	0.052	0.515 \pm 0.053	0.545 \pm 0.051
Parity-aware	0.1205	3.809 \pm 0.045	0.0334 \pm 0.0001	7.87	0.010	0.498 \pm 0.053	0.533 \pm 0.051
Parity-aware (hybrid)	0.1212	3.803 \pm 0.045	0.0336 \pm 0.0001	7.86	0.012	0.506 \pm 0.053	0.538 \pm 0.051
Parity-aware (window)	0.1268	3.781 \pm 0.045	0.0338 \pm 0.0001	7.84	0.013	0.510 \pm 0.052	0.543 \pm 0.050
Parity-aware (window+hybrid)	0.1275	3.772 \pm 0.045	0.0340 \pm 0.0001	7.83	0.017	0.518 \pm 0.052	0.548 \pm 0.051
Parity-aware (no-dev)	0.1272	3.799 \pm 0.044	0.0341 \pm 0.0001	7.84	0.050	0.531 \pm 0.052	0.559 \pm 0.051
Parity-aware (hybrid+no-dev)	0.1271	3.797 \pm 0.044	0.0341 \pm 0.0001	7.84	0.050	0.531 \pm 0.052	0.559 \pm 0.051

Table 7: Intrinsic evaluation of 256k tokenizers on the (*unbalanced*) 30-lang dataset. Values are global statistics, except for MorphScore, which is macro-averaged across available languages.

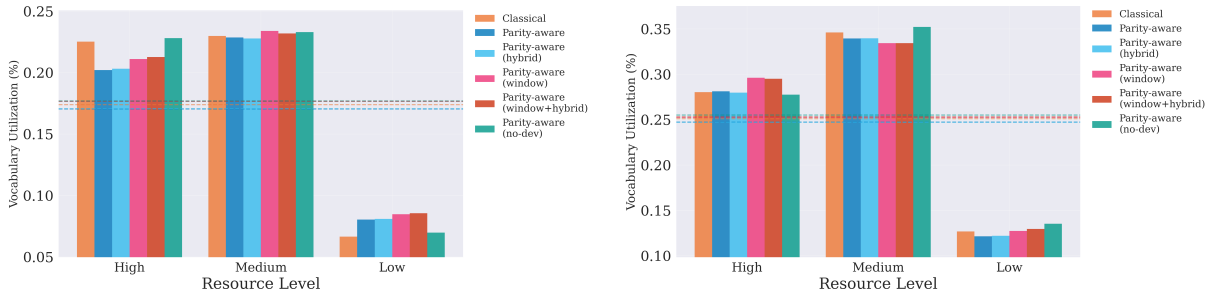


Figure 4: Vocabulary utilization grouped by language resource levels for the 256k tokenizer trained on the (*unbalanced*) 30-lang dataset (left) and 128k tokenizer trained on the (*balanced*) 30-lang dataset (right).

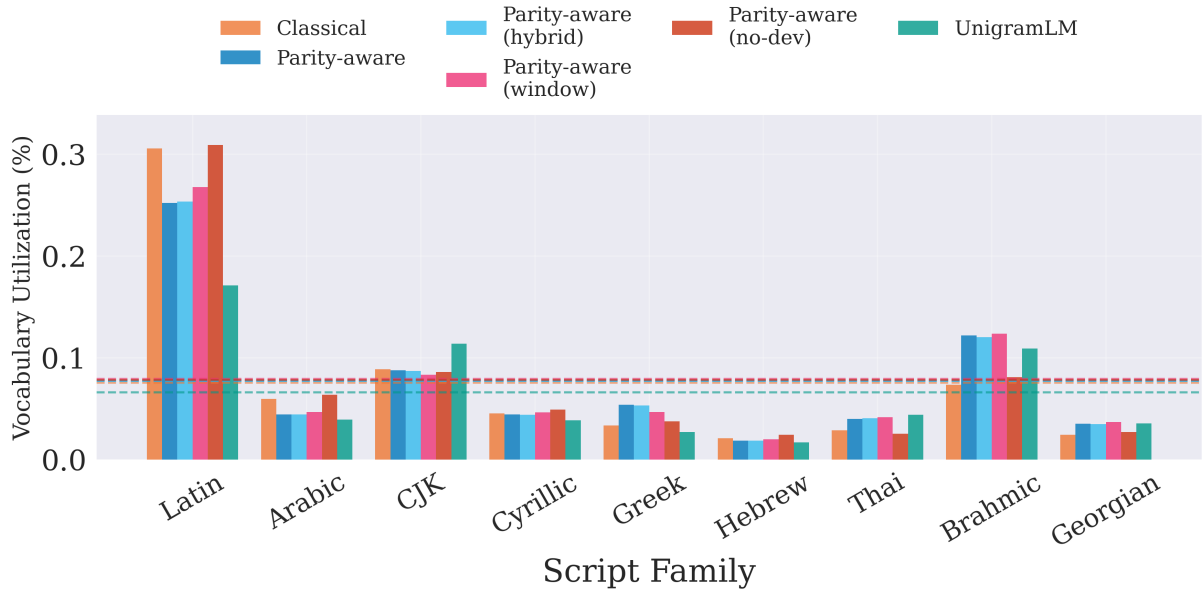


Figure 5: Vocabulary utilization across language scripts for 128k tokenizers on the (*unbalanced*) 30-lang dataset.

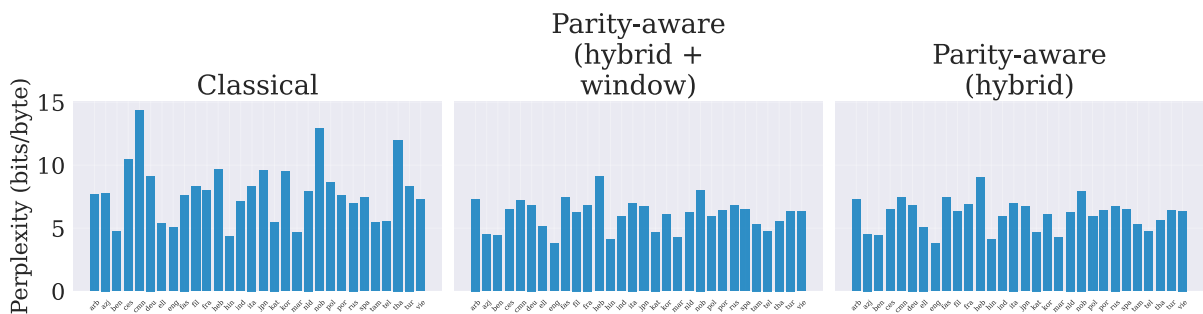


Figure 6: Per-language perplexities normalized by byte of language models trained using the specified tokenizer trained on the (*unbalanced*) 128k 30-lang dataset. Results are computed on the language model validation set at the final checkpoint (see §F for language model details).

approach does not account for heterogeneity within groups or for the local utility of individual merge operations.

Parity-aware BPE instead operates on the unaltered corpus distribution and introduces language-level constraints directly into the merge selection process. By augmenting merge scores with parity-aware penalties or budgets, this approach enables fine-grained control over how subword units are allocated across languages, while preserving the efficiency benefits of frequency-driven merging. In contrast to the balancing approach, parity-aware merges maintain alignment with the downstream data distribution and typically yield lower average tokenization costs at inference, while mitigating systematic over-fragmentation of underrepresented languages. Overall, parity-aware BPE offers a principled mechanism for addressing tokenizer bias at the level of individual merge decisions, complementing existing corpus-level mitigation strategies.

F Language Model Training

Here we provide details about the language models used for evaluating extrinsic tokenizer metrics.

F.1 Model Architecture

Our experiments focus on models with 3 billion parameters (3B), following the LLaMA architecture (Touvron et al., 2023). Model size is determined by adjusting the number of layers, hidden dimensions, and attention heads.

F.2 Training Hyperparameters

We train our models using SwissAI’s fork of Hugging Face’s Nanotron repository.¹¹ The key training hyperparameters are as follows:

- **Learning Rate.** We use a learning rate of 8×10^{-4} with linear warmup over the first 4% of training. A “1-sqrt” decay schedule (Hägele et al., 2024) is applied during the final 20%.
- **Optimizer.** All experiments use AdamW with $\beta = (0.9, 0.95)$ (Loshchilov and Hutter, 2019).
- **Weight Decay.** We set the weight decay parameter to $\lambda = 0.1$ for regularization.
- **Batch Size.** The micro-batch size is fixed at 5 across all runs.

¹¹The codebase: <https://github.com/swiss-ai/nanotron-multilingual>

F.3 Hardware Setup

Training is performed on a large-scale cluster. Each node is equipped with 4 NVIDIA Grace-Hopper H100 GPUs (96 GB memory each). The 3B models are trained on 64 nodes (256 GPUs in total), requiring approximately 18 hours per 100B tokens. This corresponds to a global batch size of 640 examples.

F.4 Sampling Methods

Let \mathcal{L} be the set of languages in the dataset, and let $\pi^{\text{natural}} \in \Delta_{|\mathcal{L}|}$ represent the natural distribution of these languages, defined as:

$$\pi_l^{\text{natural}} = \frac{\omega_l}{\sum_{l' \in \mathcal{L}} \omega_{l'}}$$

where ω_l denotes the number of words (or tokens) for language l in the dataset. In this work, we use the number of words as a proxy for language frequency, a common practice when presenting statistics for highly multilingual datasets (Penedo et al., 2025). We use Temperature Sampling which is defined as:

This method adjusts the natural distribution using a temperature parameter τ to create a less skewed distribution:

$$\pi_l^{\text{temp}, \tau} = \frac{\omega_l^{1/\tau}}{\sum_{l' \in \mathcal{L}} \omega_{l'}^{1/\tau}}$$

By tuning τ , the distribution can be shifted towards uniformity, thereby reducing imbalance among languages.

G Downstream Benchmark Evaluation

We evaluate our models using HuggingFace’s Lighteval codebase (Habib et al., 2023).

G.1 Benchmarks

To assess multilingual performance, we select twelve widely used benchmarks that span diverse downstream tasks, including reading comprehension, commonsense reasoning, semantic similarity, and knowledge-based evaluation:

- **Belebele:** A multilingual reading comprehension dataset containing passages and corresponding questions in many languages. It evaluates models’ ability to understand text and answer related questions (Bandarkar et al., 2024).

- **mTruthfulQA**: A multilingual extension of TruthfulQA, designed to measure whether models generate accurate and non-misleading answers across a broad range of questions (Lin et al., 2022a; Lai et al., 2023).
- **PAWS-X**: A multilingual paraphrase identification benchmark that extends the original PAWS dataset, providing sentence pairs annotated for semantic equivalence (Yang et al., 2019).
- **XCodah**: A multilingual adaptation of CODAH for adversarially-authored commonsense reasoning tasks, testing robustness in natural language understanding (Lin et al., 2021; Chen et al., 2019).
- **XCSQA**: A multilingual version of CommonsenseQA, consisting of multiple-choice questions that require reasoning about everyday concepts and their relations (Lin et al., 2021; Talmor et al., 2019).
- **XNLI**: A cross-lingual natural language inference benchmark, evaluating whether models can perform entailment, contradiction, and neutrality classification across multiple languages (Conneau et al., 2018).
- **XStoryCloze**: A multilingual extension of the StoryCloze Test, where models must choose the most coherent ending to short narratives, testing story comprehension and commonsense reasoning (Mostafazadeh et al., 2017; Lin et al., 2022b).
- **XWinogrande**: A multilingual version of Winogrande, containing sentences with ambiguous pronouns. It measures models' ability to resolve coreference using contextual and commonsense cues (Sakaguchi et al., 2021; Muennighoff et al., 2023; Tikhonov and Ryabinin, 2021).
- **MMMLU**: A multilingual adaptation of MMLU, evaluating model performance across a wide spectrum of tasks and domains (Hendrycks et al., 2021; Lai et al., 2023).
- **INCLUDE**: A large-scale benchmark covering 44 languages, designed to evaluate multilingual LLMs in realistic language environ-

ments with a focus on knowledge and reasoning (Romanou et al., 2025).

- **Exams**: A benchmark of standardized test questions across subjects and educational levels, used to assess reasoning and problem-solving abilities in exam-like conditions (Hardalov et al., 2020).
- **M3Exams**: A multilingual exam-style benchmark that extends Exams across different languages, subjects, and difficulty levels (Zhang et al., 2023).

G.2 Score Aggregations

We aggregate benchmark results to compute a language-specific score for each model. Let \mathcal{T}_l be the set of benchmarks (or tasks) containing a split for language l . The aggregated score for a model m per language l is defined as:

$$s_l^m = \frac{1}{|\mathcal{T}_l|} \sum_{t \in \mathcal{T}_l} s_{t,l}^m$$

where s_l^m is the score of a model m on the split l of a task t . To mitigate biases arising from varying numbers of benchmarks per language, we compute a language-specific random baseline ζ_l . This baseline helps assess whether a given aggregated score significantly outperforms random predictions. Specifically, we calculate the random baseline for each language as the average of the individual random baselines across all tasks that include language l :

$$\zeta_l = \frac{1}{|\mathcal{T}_l|} \sum_{t \in \mathcal{T}_l} \zeta_t$$

Language	Language Family	Script	Resource Level	30-lang	60-lang
English	Indo-European (Germanic)	Latin	High	✓	✓
German	Indo-European (Germanic)	Latin	High	✓	✓
French	Indo-European (Romance)	Latin	High	✓	✓
Italian	Indo-European (Romance)	Latin	High	✓	✓
Russian	Indo-European (Slavic)	Cyrillic	High	✓	✓
Spanish	Indo-European (Romance)	Latin	High	✓	✓
Japanese	Japonic	Kanji & Kana (CJK)	Medium	✓	✓
Polish	Indo-European (Slavic)	Latin	Medium	✓	✓
Portuguese	Indo-European (Romance)	Latin	Medium	✓	✓
Vietnamese	Austroasiatic	Latin	Medium	✓	✓
Turkish	Turkic	Latin	Medium	✓	✓
Dutch	Indo-European (Germanic)	Latin	High	✓	✓
Indonesian	Austronesian	Latin	Medium	✓	✓
Arabic	Afro-Asiatic (Semitic)	Perso-Arabic	Medium	✓	✓
Czech	Indo-European (Slavic)	Latin	Medium	✓	✓
Persian (Farsi)	Indo-European (Iranian)	Perso-Arabic	Medium	✓	✓
Greek	Indo-European (Hellenic)	Greek	Medium	✓	✓
Chinese (Mandarin)	Sino-Tibetan	Hanzi (CJK)	Medium	✓	✓
Hindi	Indo-European (Indo-Aryan)	Devanagari (Brahmic)	Medium	✓	✓
Korean	Koreanic	Hangugeo (CJK)	Medium	✓	✓
Thai	Kra–Dai (Tai)	Thai	Medium	✓	✓
Hebrew	Afro-Asiatic (Semitic)	Hebrew	Medium	✓	✓
Bengali	Indo-European (Indo-Aryan)	Bengali (Brahmic)	Medium	✓	✓
Tamil	Dravidian (Brahmic)	Tamil	Low	✓	✓
Georgian	Kartvelian	Georgian	Low	✓	✓
Marathi	Indo-European (Indo-Aryan)	Devanagari (Brahmic)	Medium	✓	✓
Filipino	Austronesian	Latin	Low	✓	✓
Telugu	Dravidian	Telugu (Brahmic)	Low	✓	✓
Norwegian	Indo-European (Germanic)	Latin	Medium	✓	✓
North Azerbaijani	Turkic	Latin	Low	✓	✓
Swedish	Indo-European (Germanic)	Latin	Medium	-	✓
Romanian	Indo-European (Romance)	Latin	Medium	-	✓
Ukrainian	Indo-European (Slavic)	Cyrillic	Medium	-	✓
Hungarian	Uralic (Ugric)	Latin	Medium	-	✓
Danish	Indo-European (Germanic)	Latin	Medium	-	✓
Finnish	Uralic (Finnic)	Latin	Medium	-	✓
Bulgarian	Indo-European (Slavic)	Cyrillic	Low	-	✓
Slovak	Indo-European (Slavic)	Latin	Low	-	✓
Catalan	Indo-European (Romance)	Latin	Low	-	✓
Malay	Austronesian	Latin	Low	-	✓
Urdu	Indo-European (Indo-Aryan)	Perso-Arabic	Low	-	✓
Belarusian	Indo-European (Slavic)	Cyrillic	Medium	-	✓
Basque	Language Isolate	Latin	Low	-	✓
Tajik	Indo-European (Iranian)	Cyrillic	Medium	-	✓
Sotho (Sesotho)	Niger–Congo (Bantu)	Latin	Low	-	✓
Yoruba	Niger–Congo	Latin	Low	-	✓
Swahili	Niger-Congo (Bantu)	Latin	Low	-	✓
Estonian	Uralic (Finnic)	Latin	Low	-	✓
Latvian	Indo-European (Slavic)	Latin	Low	-	✓
Galician	Indo-European (Romance)	Latin	Low	-	✓
Welsh	Indo-European (Celtic)	Latin	Low	-	✓
Albanian	Indo-European	Latin	Low	-	✓
Macedonian	Indo-European (Slavic)	Cyrillic	Low	-	✓
Malayalam	Dravidian	Malayalam (Brahmic)	Low	-	✓
Burmese	Sino-Tibetan	Mon–Burmese	Low	-	✓
Gujarati	Indo-European (Indo-Aryan)	Gujarati (Brahmic)	Low	-	✓
Afrikaans	Indo-European (Germanic)	Latin	Low	-	✓
Hawaiian	Austronesian	Latin	Low	-	✓
Northern Uzbek	Turkic	Latin	Low	-	✓

Table 8: Details on the languages used to train and evaluate tokenizers.

Language	INCLUDE	Belebele	Exams	M3Exam	MMMLU	mTruthfulQA	PAWS-X	XCodah	XCOPA	XCSQA	XNLI	XStoryCloze	XWinoGrande
English	-	✓	-	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Chinese	✓	✓	-	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Vietnamese	✓	✓	✓	✓	✓	✓	-	✓	✓	✓	✓	-	-
Arabic	✓	✓	✓	-	✓	✓	-	✓	-	✓	✓	✓	-
German	✓	✓	✓	-	✓	✓	✓	✓	-	✓	-	-	-
Spanish	✓	✓	✓	-	✓	✓	-	✓	-	✓	✓	✓	-
French	✓	✓	-	-	✓	✓	✓	✓	-	✓	✓	-	✓
Portuguese	✓	✓	✓	-	✓	✓	-	✓	-	✓	-	-	✓
Hindi	✓	✓	-	-	✓	✓	-	✓	-	✓	✓	✓	-
Russian	✓	✓	-	-	✓	✓	-	✓	✓	✓	✓	✓	✓
Indonesian	✓	-	-	-	✓	✓	-	✓	-	-	-	✓	-
Italian	✓	✓	✓	✓	✓	✓	-	✓	✓	✓	-	-	-
Japanese	✓	-	-	-	-	-	✓	✓	-	✓	-	-	✓
Swahili	-	✓	-	✓	-	-	-	✓	-	✓	✓	✓	-
Tamil	✓	-	-	-	-	-	-	-	✓	-	-	-	-
Telugu	✓	✓	-	-	✓	✓	-	-	-	-	-	✓	-
Thai	-	✓	-	✓	-	-	-	-	✓	-	✓	-	-
Basque	✓	-	-	-	✓	✓	-	-	-	-	-	✓	-
Turkish	✓	✓	✓	-	-	-	-	-	✓	-	✓	-	-
Bulgarian	✓	✓	✓	-	-	-	-	-	-	-	✓	-	-
Albanian	✓	✓	✓	-	-	-	-	-	-	-	-	-	-
Polish	✓	✓	-	-	-	-	-	✓	-	-	-	-	-
Bengali	✓	-	-	-	✓	✓	-	-	-	-	-	-	-
Serbian	✓	-	✓	-	-	✓	-	-	-	-	-	-	-
Estonian	✓	-	-	-	-	-	-	-	✓	-	-	-	-
Macedonian	✓	✓	-	-	-	-	-	-	-	-	-	-	-
Lithuanian	✓	-	✓	-	-	-	-	-	-	-	-	-	-
Greek	✓	-	-	-	-	-	-	-	-	-	✓	-	-
Urdu	✓	-	-	-	-	-	-	-	-	-	✓	-	-
Catalan	-	-	-	-	✓	✓	-	-	-	-	-	-	-
Persian	✓	-	-	-	-	-	-	-	-	-	-	-	-
Finnish	✓	-	-	-	-	-	-	-	-	-	-	-	-
Korean	✓	-	-	-	-	-	-	-	-	-	-	-	-
Quechua	-	-	-	-	-	-	-	-	✓	-	-	-	-
Haitian Creole	-	-	-	-	-	-	-	-	✓	-	-	-	-
Malay	-	-	-	-	-	-	-	-	-	-	-	✓	-

Table 9: Coverage of downstream benchmarks across languages.