

Phonemes to the Rescue: Multilingual Tokenization Based on International Phonetic Alphabet

Anonymous ACL submission

Abstract

Multilingual language models often exhibit performance disparities across languages that can arise as early as the tokenization stage. Widely-used subword tokenization approaches favor high-resource languages, and tokenizer-free methods still yield longer sequences for scripts with a higher bytes-per-character ratio. To address these shortcomings, we propose to use the International Phonetic Alphabet (IPA) as a language-agnostic input representation for multilingual tokenizers. IPA provides a compact symbol inventory, greater cross-lingual character overlap, and a more balanced byte-per-character distribution across languages. We train matched pairs of text vs. IPA subword tokenizers across 24 languages and 14 scripts and demonstrate that IPA tokenizers consistently improve tokenization quality, especially for non-Latin scripts, and generalize more effectively to unseen languages and scripts.

1 Introduction

Despite their growing global impact, multilingual language models (MLMs) continue to exhibit performance disparities across languages. Prior work shows that these differences originate as early as the tokenization stage, with direct consequences for downstream model performance (Petrov et al., 2023; Lotz et al., 2025; Arnett and Bergen, 2025; Rust et al., 2020). The main challenge of multilingual tokenization is how to represent a large number of diverse languages in a common fixed-size vocabulary. Currently predominant subword tokenization algorithms (BPE, Sennrich et al. 2016; UnigramLM, Kudo 2018) learn this vocabulary by optimizing objectives defined over corpus statistics. By design, this process favors languages that are more prevalent in the training set. As a result, underrepresented languages, especially those with rich morphologies or those using non-Latin scripts, often require more tokens for encoding the

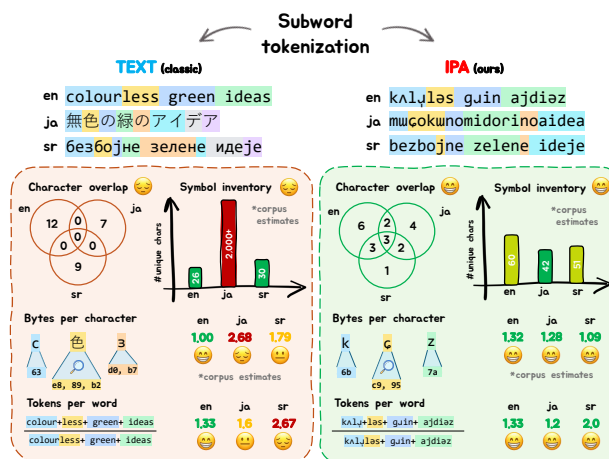


Figure 1: Benefits of IPA for multilingual tokenization. We show how our trained **Text** and **IPA** tokenizers segment the same input sequence in English (en), Japanese (ja), and Serbian (sr). Switching to IPA maps all languages into a shared phonetic alphabet, which (i) increases cross-lingual character overlap, (ii) reduces the symbol inventory needed to represent the corpus, and (iii) mitigates byte-per-character differences. Together, this results in better overall tokenization quality, as shown for instance by reduced tokens per word (what we later refer to as word fertility).

same semantic content (Ahia et al., 2023). This leads to increased computational cost, slower inference times, and diminished fluency and output quality in those languages (Qin et al., 2025). A prominent alternative to subword tokenization in recent years is *tokenization-free* byte-level modeling, which represents each input byte as a token. While this approach avoids many of the issues associated with traditional tokenizers, it does not fully eliminate cross-lingual disparities: characters in some scripts require up to 4× more bytes to encode than others, leading to longer input sequences—and thus, higher computational costs—for those languages (Mielke et al., 2021; Limisiewicz et al., 2024). This indicates that tokenization quality across languages depends not only on the tokenization method (or lack thereof), but also on the input encoding on

059 which it operates.

060 Guided by this observation, in this work we
061 consider an alternative input representation that
062 is offered by the International Phonetic Alphabet
063 (IPA, [International Phonetic Association 1999](#)) in-
064 stead of using standard orthography. IPA maps lan-
065 guages into a shared (phonetic) alphabet, increasing
066 character-level overlap across languages that oth-
067 erwise use disjoint scripts. Because the symbol
068 inventory is small (≈ 200 characters), a fixed-size
069 subword vocabulary can be used more efficiently,
070 with less capacity wasted on script-specific vari-
071 ants of the same words or morphemes. Finally, IPA
072 is encoded mostly with 1–2 byte UTF-8 charac-
073 ters, which mitigates script-driven differences in
074 bytes-per-character and therefore reduces dispari-
075 ties in per-language encoding cost. An overview
076 of IPA’s benefits is illustrated in Fig. 1 and we pro-
077 vide empirical results to support these claims in
078 Appendix I. These benefits inspired IPA’s usage
079 in NLP at various stages of the modeling pipeline,
080 from pre-training tasks ([Gale et al., 2023](#)) to prompt-
081 ing ([Nguyen et al., 2024](#)), or even integrating it into
082 multimodal pipelines ([Matsuhira et al., 2023](#)). How-
083 ever, to the best of our knowledge, we are the first
084 to systematically evaluate the effects of IPA on mul-
085 tilingual tokenization quality across a diverse set of
086 configurations and languages.

087 Concretely, we train subword tokenizers on mul-
088 tilingual data spanning 24 languages and 14 distinct
089 scripts in two input formats—standard orthography
090 and IPA—while varying the tokenization algorithm,
091 vocabulary size, and data sampling strategy. We
092 evaluate each tokenizer on a suite of intrinsic tok-
093 enization quality metrics measuring compression,
094 token frequency distribution properties, vocabu-
095 lary usage, and cross-lingual equity. We further
096 train GPT-2 models on a set of selected tokenizers
097 and evaluate them on two widely used multilin-
098 gual downstream tasks: XNLI and PAWS-X. We
099 find that IPA tokenizers consistently outperform
100 standard text tokenizers on intrinsic metrics, under
101 their respective best configurations, without sacri-
102 ficing downstream task performance. Overall, our
103 approach preserves MLM capabilities, while yield-
104 ing more equitable treatment across languages.

105 **2 Background and Related Work**

106 **Disparities in tokenization quality.** Prior work
107 has shown that standard multilingual subword tok-
108 enizers systematically treat languages inequitably

([Petrov et al., 2023](#); [Ahia et al., 2023](#); [Ali et al., 2024](#), inter alia). Resulting frequency-based vocabu-
109 larities allocate more capacity to high-resource lan-
110 guages, so that lower-resource languages are split
111 into many more tokens per word. This is especially
112 amplified in case of morphologically more complex
113 languages, as well as those using non-Latin scripts.
114 [Petrov et al. \(2023\)](#) quantify this disparity, show-
115 ing a difference in the amount of tokens needed
116 to encode the same semantic content of up to $15\times$
117 between languages. This over-segmentation has di-
118 rect consequences, as the impacted languages will
119 require much higher training and inference costs.¹
120

121 These findings have inspired several recent works
122 that aim to tackle these tokenization challenges.
123 [Petrov et al. \(2023\)](#) suggest training monolingual to-
124 kenizers for each language and then merging them
125 based on tokenization parity. Inheriting the idea of
126 tokenization parity, [Foroutan et al. \(2025\)](#) introduce
127 Parity-Aware BPE, optimizing the training objec-
128 tive in classic BPE such that merges are guided by
129 per-language compression levels, instead of simple
130 frequency. [Feher et al. \(2025\)](#) propose a dynamic
131 tokenization approach, which merges frequently co-
132 occurring subwords on the fly to reduce the token-
133 count inflation. The main alternative to subword to-
134 kenization in recent years have been different flavors
135 of byte-level approaches, which eliminate learned
136 segmentation by treating each byte in the input as
137 a separate token ([Xue et al., 2022](#); [Kallini et al., 2024](#);
138 [Yu et al., 2023](#), inter alia). However, cross-
139 lingual disparities still remain, as languages dif-
140 fer in the number of bytes-per-character, as well
141 as the number of symbols needed to encode the
142 same semantic content ([Arnett et al., 2024](#)). [Ahia et al. \(2024\)](#)
143 attempt to tackle this issue by integrat-
144 ing adaptive gradient-based tokenization into the
145 model through language- and script-specific pre-
146 dictors that learn to place token boundaries such
147 that they equalize segmentation granularity across
148 scripts. While our method could in principle be
149 applied to byte-level approaches, we focus on cur-
150 rently predominant subword tokenization in this
151 work. We provide a summary of different tokeniza-
152 tion methods in Appendix B.
153

154 **IPA in Natural Language Processing.** IPA has
155 seen a growing range of applications in NLP in
156 recent years, particularly in settings that require

¹[Lundin et al. \(2025\)](#) refer to this effect as *token tax*, show-
ing that doubling in tokens can increase training cost by a
factor of four.

Stage	LATN										ARAB			CYRL		DEVA	JPAN	HANG	HANI	THAI	LAOO	MYMR	GREK	HEBR	ETHI
	DE	EN	ES	FI	FR	IT	PL	SW	TR	AR	FA	UR	RU	SR	HI	JA	KO	ZH	TH	LO	MY	EL	HE	AM	
Tokenizer training	✓	✓	-	✓	✓	✓	-	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	-	-	-	-		
GPT-2 pre-training	✓	✓	-	✓	✓	✓	-	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	-	-	-	-		
Intrinsic eval	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓		
XNLI	✓	✓	✓	-	✓	-	-	✓	✓	-	✓	✓	-	✓	-	-	✓	✓	-	-	✓	-	-		
PAWS-X	✓	✓	✓	-	✓	-	-	-	-	-	-	-	-	-	✓	✓	✓	-	-	-	-	-	-		

Table 1: Language coverage at different stages of our experiments. Columns show languages (ISO 639-1 codes) grouped by script (ISO 15924 codes). We provide the language/script-to-ISO mapping in Appendix C.1. ✓ indicates that a language is included in the corresponding stage. For tokenizer and GPT-2 training we use a sample from the CulturaX dataset (Nguyen et al., 2023). For intrinsic evaluation we use WikiPron (Lee et al., 2020) and FLORES+ (NLLB Team et al., 2024). We perform downstream evaluation on XNLI (Conneau et al., 2018a) and PAWS-X (Yang et al., 2019a).

phonological awareness. Leong and Whitenack (2022) and Sohn and Mortensen (2025) show benefits of IPA representations for improved performance in named entity recognition (NER) tasks. Gale et al. (2023) introduce the BORT model, an extension of BART (Lewis et al., 2020) with added self-supervised pre-training pronunciation tasks by leveraging IPA representations. Nguyen et al. (2024) propose phonemic prompting as a way to enhance the multilingual capabilities of LLMs. Matsuhira et al. (2023) developed IPA-CLIP, by incorporating IPA representations into a multimodal setting, extending the CLIP model. Taken together, these studies highlight the versatility and effectiveness of IPA in enhancing cross-lingual generalization, representation learning, and multilingual performance. Most closely related to our work are Goriely et al. (2024), who also pre-train GPT-2 models on IPA input, testing both character-level and BPE tokenization. However, their study is limited to a monolingual setting (English) and does not analyze the impact of IPA on tokenization quality.

3 Data and Languages

3.1 Languages

Table 1 summarizes language coverage at each stage of our experiments. For tokenizer training (Section 4.3) and GPT-2 pre-training (Section 4.5) we consider 18 languages spanning 10 scripts. We evaluate the intrinsic tokenization quality (Section 4.4) on the same language set together with six additional languages (AM, EL, ES, HE, MY, PL) that introduce four new scripts (GREK, HEBR, ETHI, MYMR) and allow us to test zero-shot generalization to languages/scripts unseen during tokenizer training. In total, this results in 24 languages and 14 scripts considered in this work.

Our language selection aimed to cover a diverse sample across orthographic, typological, and resource profiles. For instance, it includes high-resource Indo-European languages (e.g., EN, DE, FR, IT, RU), as well as comparatively lower-resource languages (e.g., SW, LO, MY, AM), which are often underrepresented in multilingual pre-training pipelines. We also consider morphologically rich languages (e.g., FI, TR) which often exhibit worse tokenization quality (Raj et al., 2024) and those with more complex orthographies (e.g., ZH, JA, KO, TH, LO) which are typically more expensive to encode due to their Unicode ranges. We also include a language written in multiple scripts, JA, spanning Katakana, Hiragana, and Kanji (collectively labeled as JPAN).^{2,3}

3.2 Datasets

CulturaX. For tokenizer training (Section 4.3) and GPT-2 pre-training (Section 4.5), we use CulturaX (Nguyen et al., 2023), a large-scale multilingual web corpus derived from filtered mC4 (Xue et al., 2021) and OSCAR (Suárez et al., 2019, 2020). We apply additional post-processing to remove residual non-linguistic artifacts and reduce cross-language contamination on our sample of languages (Appendix N). We sample a total of 1GB of data distributed across our 18 languages using four strategies: (i) **byte-uniform**, allocating an equal number of bytes to each language; (ii) **semantic-uniform**, allocating equal semantic content by compensating for language-specific byte premiums (Arnett et al., 2024); (iii) **data-proportional**, sampling in proportion to each language’s share

²Although, in our experiments we only consider Katakana and Hiragana, due to limitations in EpiTran’s Kanji support

³Serbian is also a multi-script language (LATN and CYRL), but our sample only contained Cyrillic text.

in CulturaX; and (iv) `data-smoothed`, applying temperature sampling ($\alpha=0.3$) to upweight lower-resource languages (as done in mT5, Xue et al. (2021); mBERT, Devlin (2018); XLM-R, (Conneau et al., 2020)). The first two strategies provide a *balanced* sample of languages (under different notions of balance), while the latter two result in *unbalanced*, but arguably more realistic distributions in practice. More details and distribution visualizations are in Appendix C.2.

WikiPron and FLORES+. For evaluating intrinsic tokenization quality we use two datasets: (i) WikiPron (Lee et al., 2020), a multilingual collection of word lists paired with IPA transcriptions extracted from Wiktionary⁴ and (ii) FLORES+ (NLLB Team et al., 2024), a dataset of human-translated parallel sentences across more than 200 languages. We use WikiPron for word-level metrics (word fertility and proportion of continued words; see Section 4.4) and FLORES+ for all the remaining metrics, computed at the sentence level. Details can be found in Appendix C.

XNLI and PAWS-X. We evaluate cross-lingual transfer on XNLI (Conneau et al., 2018b) and PAWS-X (Yang et al., 2019b), two standard multilingual natural language understanding (NLU) benchmarks. These two benchmarks cover 13 and 7 languages from our selection, respectively, as shown in Table 1.

4 Methodology

4.1 Overview

Our experiments compare subword tokenization trained on two input representations: standard orthographic text (`Text`) and its phonemic transcription in the International Phonetic Alphabet (`IPA`). Starting from the same multilingual corpora, we (i) build an IPA version of each dataset with a multilingual grapheme-to-phoneme (G2P) pipeline, (ii) train matched pairs of `Text` vs. `IPA` tokenizers on different configurations, (iii) evaluate intrinsic tokenization quality across languages, and (iv) pre-train GPT-2 models with selected tokenizers and evaluate cross-lingual transfer by fine-tuning on English and testing zero-shot on other languages.

4.2 Grapheme-to-phoneme conversion

To train and evaluate IPA tokenizers and IPA-based language models, we require IPA text at scale.

⁴<https://www.wiktionary.org/>

While some resources provide curated IPA transcriptions (e.g., WikiPron (Lee et al., 2020) or Universal Dependencies (UD; Nivre et al., 2020)), these are often limited in size and language coverage. In this work, we utilize Epitran (Mortensen et al., 2018), an open-source rule-based grapheme-to-phoneme (G2P) toolkit designed specifically for multilingual applications. Epitran supports more than 60 languages across a variety of scripts and is being actively extended. It takes standard orthographic text as input, together with the source language and script code (e.g. `tokenization_eng-Latn`), and returns its IPA equivalent (e.g. `towkənəzɛjʃən`). We applied minor language-specific fixes and extensions, as well as efficiency improvements to Epitran’s conversion pipeline, to adapt it best to our needs (see Appendix D for a detailed discussion). For languages that are not supported by Epitran (Greek and Hebrew), we use Phonemizer (Bernard and Titeux, 2021) instead.

4.3 Tokenizer training

To isolate how the input representation affects subword tokenization we vary a set of configurations outlined below. For each configuration we train a paired set of tokenizers: one on `Text` and one on `IPA`. We train all tokenizers with SentencePiece (Kudo and Richardson, 2018) to allow comparisons across subword algorithms under consistent training conditions and because it avoids whitespace-based pre-tokenization, which is important for languages in our study without explicit word boundary markers. We detail hyper-parameter settings and representation-specific adjustments (such as Unicode normalization and character coverage) in Appendix E.1.

Experimental grid. We vary three factors in our experiments: (1) subword tokenization algorithm (`BPE`, `UnigramLM`), (2) vocabulary size (40k, 80k, 100k, 200k), and (3) data sampling strategy (`byte-uniform`, `semantic-uniform`, `data-proportional`, `data-smoothed`; defined in Section 4.3). This yields $2 \times 4 \times 4 = 32$ configurations, and thus 32 paired `Text/IPA` tokenizers (64 tokenizers in total).

4.4 Intrinsic evaluation metrics

We assess tokenization quality using ten metrics grouped into four categories—compression, token frequency distribution shape, vocabulary usage, and cross-lingual equity—and refer to Appendix F for formal definitions.

323	4.4.1 Compression	
324	We report four compression metrics: (i) Word Fer-	
325	tility (WF) , the average number of subword tokens	
326	per unique word type, where lower values indicate	
327	that words are represented with fewer pieces; (ii)	
328	Proportion of Continued Words (PCW) , the frac-	
329	tion of unique words split into more than one token,	
330	where lower values indicate fewer splits; (iii) Aver-	
331	age Token Length (ATL) , the mean token length	
332	measured in input characters of the representation	
333	(Text or IPA), where higher values indicate tokens	
334	spanning longer segments; and (iv) Compression	
335	Rate (CR) , the average number of input characters	
336	per token at the sentence level, where higher values	
337	indicate stronger compression.	
338	4.4.2 Token frequency distribution shape	
339	We characterize the token frequency profile with	
340	two metrics: (i) Rényi Entropy (RE) , computed	
341	from token frequencies on evaluation data and re-	
342	ported for $\alpha=1$ (Shannon), $\alpha=2$ (collision), and	
343	$\alpha=\infty$ (min-entropy), where higher values indicate	
344	more uniform token usage across the vocabulary;	
345	and (ii) Zipf Deviation (ZipfD) , which measures	
346	how closely the empirical rank–frequency curve	
347	matches an ideal Zipfian distribution, where lower	
348	values indicate closer adherence to Zipf’s law (Lotz	
349	et al., 2025).	
350	4.4.3 Vocabulary usage	
351	We quantify vocabulary usage with two metrics:	
352	(i) Vocabulary Utilization (VU) , the share of the	
353	learned vocabulary that appears at least once in	
354	the evaluation data, where higher values indicate	
355	less unused capacity; and (ii) Type–Token Ratio	
356	(TTR) , the number of distinct token types divided	
357	by the total number of produced tokens, where	
358	higher values indicate more diverse token usage	
359	relative to sequence length.	
360	4.4.4 Cross-lingual equity	
361	We assess cross-lingual equity with two metrics:	
362	(i) Tokenization Parity (TP) , the average ratio of	
363	token counts between English and a target language	
364	computed over parallel sentence pairs, where values	
365	closer to 1 indicate similar segmentation length	
366	(Petrov et al., 2023); and (ii) Tokenization Fairness	
367	Gini (TFG) , the Gini coefficient over language-	
368	level tokenization cost, where cost is defined as	
369	token count normalized by input bytes, and lower	
370	values indicate more equal tokenization cost across	
371	languages (Meister, 2025).	
	4.5 GPT-2 pre-training and fine-tuning	372
	Tokenizer selection. Since training a language	373
	model for all 32 tokenizer settings for both Text	374
	and IPA is computationally expensive, we select a	375
	small set of tokenizers for GPT-2 pre-training exper-	376
	iments using our intrinsic evaluation. We rank the	377
	32 Text and 32 IPA tokenizers by macro-averaging	378
	each intrinsic metric over languages and then ag-	379
	gregating the resulting metric-wise ranks into a sin-	380
	gle mean-rank score. We select the highest ranked	381
	Text and IPA tokenizers, which we refer to as <i>Text</i>	382
	<i>Opt</i> and <i>IPA Opt</i> , respectively. <i>Text Opt</i> tokenizer	383
	is found under the following setting [BPE, 200k,	384
	data-proportional], while for <i>IPA Opt</i> it is [Un-	385
	igramLM, 200k, byte-uniform]. Because these	386
	optima arise under different settings, the compar-	387
	ison is no longer driven only by the input repre-	388
	sentation. Therefore, we also select two control	389
	tokenizers: Text tokenizer with IPA Opt settings	390
	(<i>Text Subopt</i>) and IPA tokenizer with Text Opt set-	391
	tings (<i>IPA Subopt</i>). We therefore pre-train four	392
	GPT-2 models in total. The ranking procedure and	393
	the full ranking results are reported in Appendix G.	394
	Pre-training. We train GPT-2 Small models	395
	(Radford et al., 2018, 2019) from scratch on the	396
	data-smoothed sample of CulturaX for our 18	397
	languages using the four selected tokenizers. Given	398
	that all of our tokenizers use a 200k vocabulary,	399
	we end up with a total of 240M parameter mod-	400
	els (rather than 125M in the original GPT-2 Small,	401
	which uses a 50k vocabulary). All hyperparame-	402
	ters are fixed across models and their values are	403
	reported in Appendix H.	404
	Fine-tuning and evaluation. We fine-tune each	405
	pre-trained model on English-only training data	406
	from XNLI or PAWS-X and evaluate zero-shot on	407
	all other languages, as well as on a held-out portion	408
	of the English split. We report per-benchmark fine-	409
	tuning hyperparameters in Appendix H.	410
	5 Results and Analysis	411
	5.1 Intrinsic tokenization quality across	412
	languages and configurations	413
	We begin with two complementary diagnostic	414
	views: the configuration-level heatmaps (Fig. 2),	415
	which compare intrinsic performance across all to-	416
	kenizer settings for Text vs. IPA , and a detailed	417
	per-language evaluation (Fig. 3) of the optimal tok-	418
	enizers (with mean statistics for suboptimal ones).	419

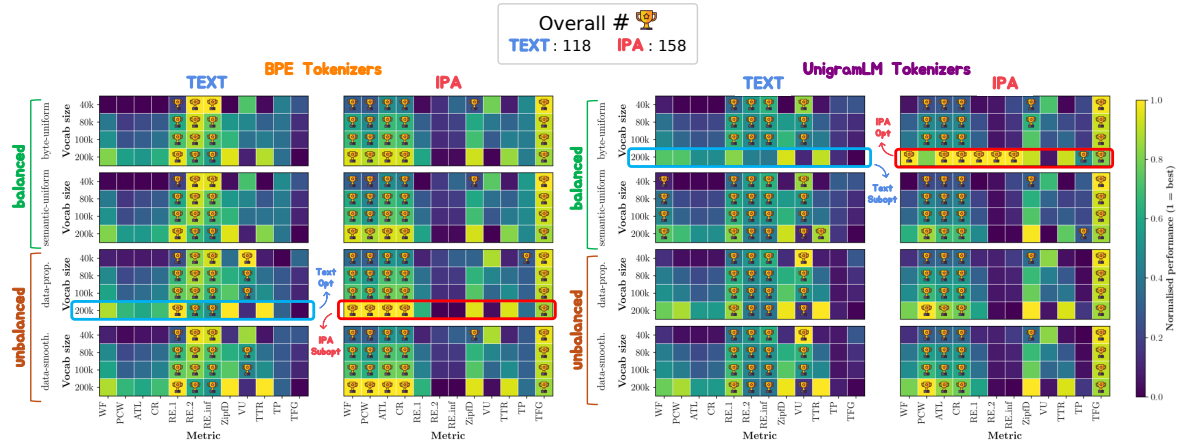


Figure 2: Intrinsic performance across all tokenizer configurations for **Text** and **IPA** representations. Each heatmap cell reports a metric value macro-averaged across languages and min–max normalized to $[0, 1]$ (higher is better) for a configuration defined by vocabulary size, subword algorithm, and data-sampling strategy. Cells marked with a trophy indicate a >0.1 advantage over the paired configuration in the other representation (**Text** vs. **IPA**). Outlined cells denote the selected *Opt/Subopt* tokenizers for GPT-2 training. **Takeaway:** Across configurations, **IPA** is consistently stronger on compression (WF, PCW, ATL, CR) and cross-lingual equity (TFG), while **Text** yields higher token-frequency uniformity as measured by Rényi entropies ($\alpha \in \{1, 2, \infty\}$). Overall, **IPA** wins more often (158 trophies) than **Text** (118 trophies).

IPA consistently improves tokenization quality.

Under their respective optimal settings, IPA improves overall tokenization quality across languages compared to Text (Fig. 3). Across all tested configurations, IPA wins more often than Text (158 vs. 118 notable wins; $\Delta > 0.1$ in the normalized metric scale), as shown in Fig. 2, with especially consistent advantages on compression metrics (WF, PCW, ATL, CR) and cross-lingual equity (TFG), whereas Text more often achieves higher Rényi entropies. Fig. 5 further summarizes paired effect sizes d_z per metric, aggregated across languages ($d_z > 0 \Rightarrow$ IPA is better than Text; see Appendix M), corroborating that IPA’s gains are systematic rather than driven by a single language or setting.

Sensitivity to different configurations. Independent of representation, several configuration trends are stable across metrics: larger vocabularies generally improve compression and yield lower Zipf deviation (best ZipfD typically at 200k), while Rényi entropies show a trade-off for different orders (H_2 and H_∞ strongest at small vocabularies (40k) but H_1 tending to peak at the largest (200k)). Vocabulary usage metrics show a similar pattern: with VU highest at 40k and TTR at 200k vocabularies. Cross-lingual equity prefers smaller vocabularies under the balanced sampling strategies, but shifts this preference towards larger vocabularies under unbalanced mixtures. We provide a more detailed interpretation of these effects in Appendix L. Fi-

nally, the optimal Text and IPA tokenizers emerge under different configurations, as outlined in Section 4.5 and highlighted in Fig. 2.

5.1.1 Metric-level analysis

IPA brings largest gains in compression and cross-lingual equity. Figure 5 shows that IPA’s largest metric improvements are in compression-related measures (ATL and CR) and cross-lingual equity (TFG), followed by consistent gains in token-frequency distribution (Rényi entropies across different orders). Compression gains indicate that IPA achieves more compact segmentations (fewer tokens for comparable content), while improved equity indicates more uniform tokenization quality across languages. In contrast, improvements are smaller for vocabulary utilization (VU) and continued-word behavior (PCW). The smaller improvements in VU are likely a combined result of our very large vocabularies (200k) and small test sets (WikiPron and FLORES+). With a large vocabulary, many tokens sit in a long tail that never appears in the evaluation sample. Smaller improvements in PCW compared to other compression metrics suggest that while words are overall split into less tokens, the IPA tokenizers still rarely store entire words in their vocabularies.

5.1.2 Language- and script-level analysis

Gains concentrate on non-Latin and unseen scripts. To characterize how different languages

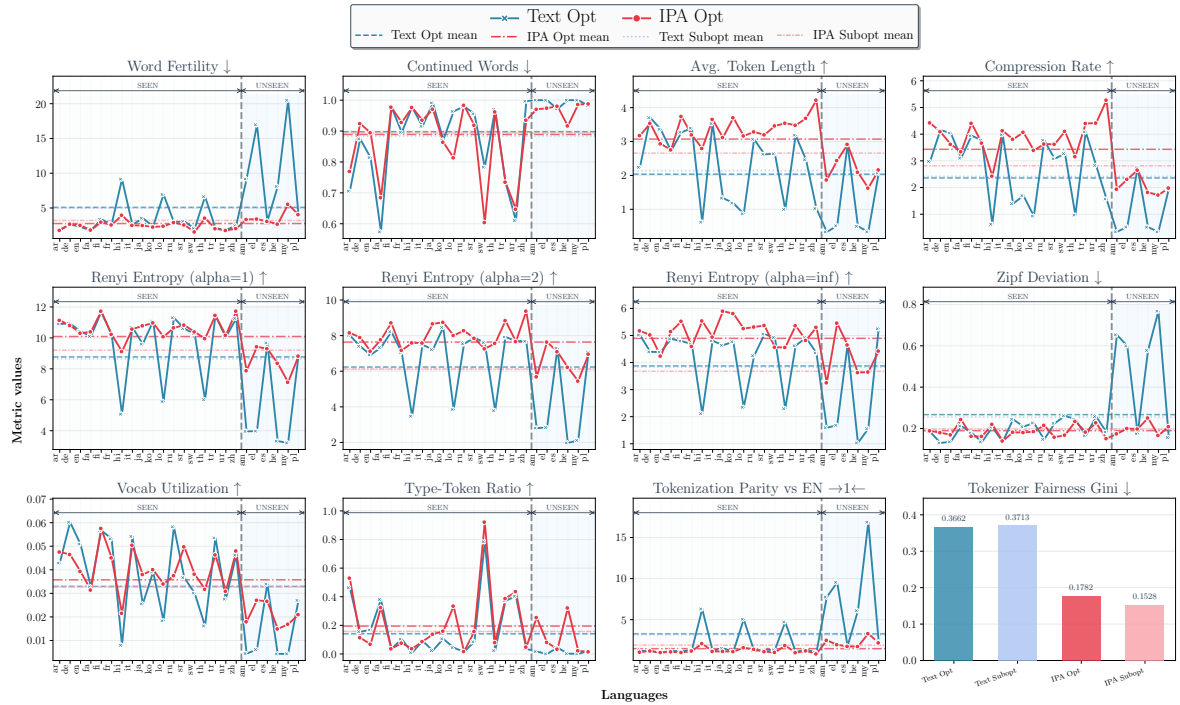


Figure 3: Per-language intrinsic tokenization metrics: **Text** vs. **IPA**. Solid lines show *Text Opt* and *IPA Opt* per language; dashed lines are means. *Subopt* variants shown as mean-only (per-language in Appendix J). Arrows: \uparrow = higher is better, \downarrow = lower is better, $\rightarrow 1 \leftarrow$ = closer-to-1 is better; TFG is a single global tokenizer metric (not per-language), hence shown as bars. SEEN/UNSEEN mark languages included in vs. absent from tokenizer training data. **Takeaway:** *IPA (Opt)* shows consistently better tokenization quality than *Text (Opt and Subopt)* across metrics, with differences particularly pronounced for unseen languages.

are impacted, we use two complementary summaries over intrinsic metrics: *win rate* (WR), the fraction of metrics where IPA outperforms Text for a given language/script, and the *mean z-score* (\bar{z}), which summarizes the average standardized magnitude of IPA–Text differences (both described in more detail in Appendix M). The language- and script-level rankings (Fig. 4a, Fig. 4b) show that the largest improvements concentrate on non-Latin scripts, with the strongest gains on unseen scripts (GREK, ETHI, MYMR, HEBR). Languages in these scripts (e.g., EL, AM, MY, HE) achieve WR=1, indicating that IPA improves all intrinsic metrics compared to Text. Substantial gains also appear for complex orthographies and segmentation regimes (e.g., LO, TH, JA (WR=1), and ZH (WR=0.9)), morphologically rich languages (e.g., FI (WR=0.9) and TR (WR=0.72)), as well as low-resource languages (e.g., LO (WR=1) and sw (WR=0.81)). These results uncover a particular fragility of Text tokenizers: when a language or script is underrepresented or completely unseen during tokenizer training, subword algorithms often split words into individual characters or even bytes, leading to fragmented segmentation. In contrast, IPA’s shared symbol space

across scripts better preserves tokenization quality in these situations.

Only a small set of high-resource Latin languages shows mild degradation. A small subset of languages—EN, DE, FR, ES, PL, RU—shows $WR < 0.5$ (Fig. 4a), meaning Text wins on a majority of intrinsic metrics for these cases. These languages are predominantly high-resource and written in Latin script, which is best represented during tokenizer training, where Text tokenizers already perform well.⁵ However, these degradations have relatively small magnitudes ($\bar{z} < 0$) compared to the improvements across other languages.

5.2 Downstream task performance and efficiency

IPA models match downstream accuracy while reducing overall inference cost. Figure 6 (top row) reports per-language accuracies for Text and IPA tokenizers (both *Opt* and *Subopt* configurations). On XNLI, IPA models match Text models on mean accuracy; on PAWS-X, IPA models are

⁵One exception is Russian, which is written in Cyrillic script, but is still one of the better represented languages in our training data.

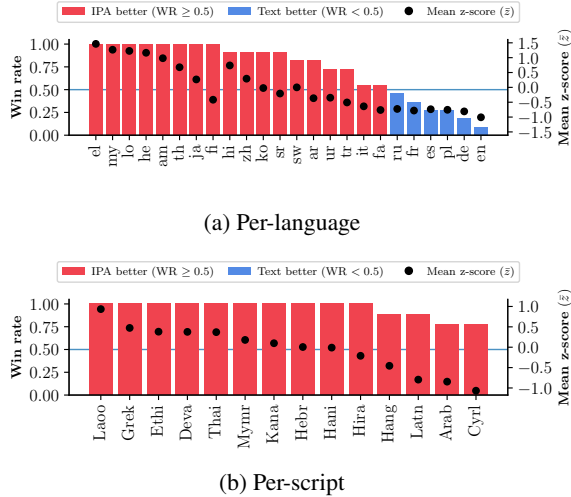


Figure 4: Ranking **IPA** impact on tokenization quality (a) per-language and (b) per-script. Bars show win rate (WR), the fraction of intrinsic metrics on which *IPA Opt* outperforms *Text Opt*: $WR \in [0, 1]$, with $WR=1/0$ indicating that **IPA/Text** wins on all metrics, respectively. **Red** bars mark languages/scripts where **IPA** wins on a majority of metrics ($WR > 0.5$), while **blue** bars mark those where **Text** wins ($WR < 0.5$). Black dots show the **mean z-score** \bar{z} (average standardized IPA–Text delta; $\bar{z} > 0$ favors IPA). **Takeaway:** The largest improvements concentrate on unseen scripts (GREK, ETHI, MYMR, HEBR) and on languages with more complex orthographies (e.g., LO, TH, JA), while only a small set of mostly high-resource Latin languages (FR, ES, PL, DE, EN) and one Cyrillic language (RU) show lower overall quality.

only slightly worse on average. Despite near-parity in accuracy, the accuracy–compression delta plots (Fig. 6 (bottom row)) show that IPA improves overall compression on these two benchmarks. This implies on average fewer tokens per input, which effectively lowers inference cost. Analogously to the intrinsic patterns, the same set of high-resource Latin languages (EN, DE, FR, ES) + RU see mild compression regressions, which are outweighed by more substantial improvements on other languages, resulting in a favorable accuracy–efficiency trade-off: near-par downstream performance with reduced token counts, translating to lower compute and latency at inference.

6 Conclusion

In this work, we show that using IPA as an input representation results in improved and more equitable multilingual tokenization quality. Across 24 languages and 14 scripts, the optimal IPA tokenizer consistently outperforms the optimal Text

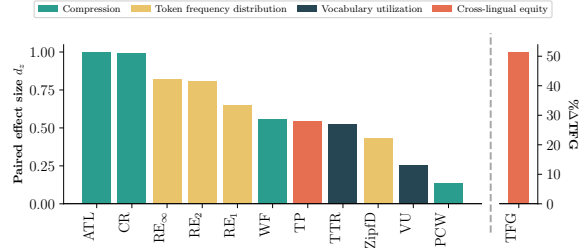


Figure 5: Ranking **IPA** impact on tokenization quality per metric. Bars are grouped and color-coded by metric family (see Sec. 4.4). Left panel reports paired effect size d_z across languages for 11 per-language metrics ($d_z > 0$ favors **IPA**, $d_z < 0$ favors **Text**). TFG is shown separately (right panel) because it is a single global tokenizer statistic rather than a per-language metric; we report its relative change ($\% \Delta TFG_{\text{Text Opt} - \text{IPA Opt}}$; > 0 favors **IPA**). **Takeaway:** IPA improves average intrinsic quality across all metrics, with the largest gains in compression (ATL, CR) and cross-lingual equity (TFG), suggesting better tokenization efficiency and more uniform cross-lingual behavior.

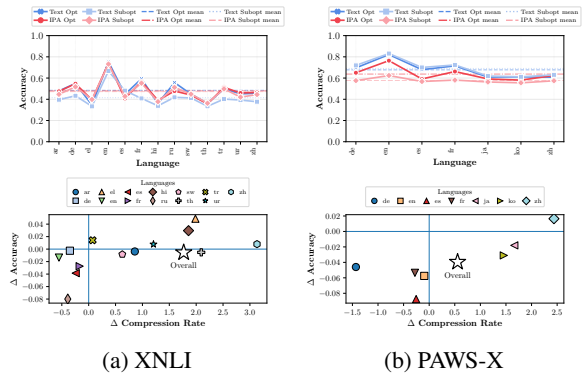


Figure 6: Per-language test accuracy for **Text** vs. **IPA** models (top) and the accuracy–compression trade-off of *IPA Opt* relative to *Text Opt* (bottom) on (a) XNLI and (b) PAWS-X. IPA models are on-par with Text on XNLI and only slightly worse on PAWS-X. In the delta plots, each point is a language with $\Delta \text{Acc} = \text{Acc}_{\text{IPA Opt}} - \text{Acc}_{\text{Text Opt}}$ versus $\Delta \text{CR} = \text{CR}_{\text{IPA Opt}} - \text{CR}_{\text{Text Opt}}$; the large star denotes the macro-average across languages. Points further right indicate better compression and points higher indicate better accuracy. **Takeaway:** Overall, IPA improves compression, resulting in reduced inference cost, with only minimal sacrifices in accuracy.

tokenizer on our suite of 10 intrinsic metrics, with largest gains in compression and cross-lingual equity, and most pronounced improvements for non-Latin and unseen scripts. Models pretrained on IPA match text-based models on downstream performance (XNLI, PAWS-X), while lowering overall inference cost due to improved compression. These establishes IPA as a promising direction for reducing tokenization-driven disparities in MLMs.

7 Limitations

Dependence on G2P quality. Our approach relies on grapheme-to-phoneme (G2P) tools to obtain IPA transcriptions at scale. In practice, G2P quality is not even across all languages, and rule-based or deterministic tools such as Epitran cannot handle pronunciation variation (e.g., British English vs. American English) or homographs, where identical spellings correspond to different pronunciations (e.g., live as a verb vs. adjective). As multilingual G2P systems improve and become more robust, IPA-based pipelines should become increasingly practical for large-scale pretraining.

Converting from IPA back to Text. Autoregressive models trained purely on IPA cannot directly generate standard orthographic text. IPA is a lossy encoding, and not always is there an unambiguous mapping back to the original writing system, making straightforward decoding a challenge. Moreover, to the best of our knowledge, there currently exist no publicly available automatic multilingual phoneme-to-grapheme (P2G) conversion tools. Since many real-world uses of generative language models require fluent orthographic output, future work could focus on developing accurate P2G tools or finding ways to incorporate the IPA-to-text mapping mechanism directly into the model architecture.

Generalizability. While we cover a diverse set of languages and vary tokenization and data-sampling configurations to improve generalizability of our findings, our pretrained GTP-2 models (240M parameters) are rather small by current standards. It therefore remains an open question whether the same trade-offs hold at larger scales. In addition, we evaluate downstream transfer on only two tasks and report a single run per model. Stronger evidence would come from broader downstream evaluation and reporting average results over multiple runs to improve robustness.

References

Orevaoghene Ahia, Sachin Kumar, Hila Gonen, Valentin Hofmann, Tomasz Limisiewicz, Yulia Tsvetkov, and Noah A Smith. 2024. Magnet: Improving the multilingual fairness of language models with adaptive gradient-based tokenization. *Advances in Neural Information Processing Systems*, 37:47790–47814.

Orevaoghene Ahia, Sachin Kumar, Hila Gonen, Jungo Kasai, David R Mortensen, Noah A Smith, and Yulia

Tsvetkov. 2023. Do all languages cost the same? tokenization in the era of commercial language models. *arXiv preprint arXiv:2305.13707*. 603
604
605

Mehdi Ali, Michael Fromm, Klaudia Thellmann, Richard Rutmann, Max Lübbering, Johannes Levelling, Katrin Klug, Jan Ebert, Niclas Doll, Jasper Buschhoff, and 1 others. 2024. Tokenizer choice for llm training: Negligible or crucial? In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3907–3924. 606
607
608
609
610
611
612

Catherine Arnett and Benjamin Bergen. 2025. Why do language models perform worse for morphologically complex languages? In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 6607–6623. 613
614
615
616
617

Catherine Arnett, Tyler A Chang, and Benjamin Bergen. 2024. A bit of a problem: Measurement disparities in dataset sizes across languages. In *Proceedings of the 3rd Annual Meeting of the Special Interest Group on Under-resourced Languages@ LREC-COLING 2024*, pages 1–9. 618
619
620
621
622
623

Lisa Beinborn and Yuval Pinter. 2023. Analyzing cognitive plausibility of subword tokenization. *arXiv preprint arXiv:2310.13348*. 624
625
626

Mathieu Bernard and Hadrien Titeux. 2021. **Phonemizer: Text to phones transcription for multiple languages in python**. *Journal of Open Source Software*, 6(68):3958. 627
628
629
630

Steven Bird, Edward Loper, and Ewan Klein. 2009. *Natural Language Processing with Python*. O’Reilly Media, Sebastopol, CA. 631
632
633

Kaj Bostrom and Greg Durrett. 2020. Byte pair encoding is suboptimal for language model pretraining. *arXiv preprint arXiv:2004.03720*. 634
635
636

Iaroslav Chelombitko, Egor Safronov, and Aleksey Komissarov. 2024. Qtok: A comprehensive framework for evaluating multilingual tokenizer quality in large language models. *arXiv preprint arXiv:2410.12989*. 637
638
639
640
641

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. 642
643
644
645
646

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th annual meeting of the association for computational linguistics*, pages 8440–8451. 647
648
649
650
651
652
653

Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel Bowman, Holger Schwenk, and Veselin Stoyanov. 2018a. **XNLI: Evaluating cross-lingual sentence representations**. In *Proceedings of* 654
655
656
657

658	<i>the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2475–2485, Brussels, Belgium. Association for Computational Linguistics.	715
659		716
660		717
661	Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel Bowman, Holger Schwenk, and Veselin Stoyanov. 2018b. Xnli: Evaluating cross-lingual sentence representations. In <i>Proceedings of the 2018 conference on empirical methods in natural language processing</i> , pages 2475–2485.	718
662		719
663		720
664		721
665		722
666		723
667	Mathias Creutz and Krista Lagus. 2005. Unsupervised morpheme segmentation and morphology induction from text corpora using morfessor 1.0. In <i>Unsupervised morpheme segmentation and morphology induction from text corpora using Morfessor 1.0</i> . Helsinki University of Technology.	724
668		725
669		726
670		727
671		728
672		729
673	Jacob Devlin. 2018. Multilingual BERT README. https://github.com/google-research/bert/blob/master/multilingual.md . Accessed: 2026-01-05.	730
674		731
675		732
676		733
677	Miguel Domingo, Mercedes García-Martínez, Alexandre Helle, Francisco Casacuberta, and Manuel Heranz. 2019. How much does tokenization affect neural machine translation? In <i>International Conference on Computational Linguistics and Intelligent Text Processing</i> , pages 545–554. Springer.	734
678		735
679		736
680		737
681		738
682		739
683	Darius Feher, Ivan Vulić, and Benjamin Minixhofer. 2025. Retrofitting large language models with dynamic tokenization. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 29866–29883.	740
684		741
685		742
686		743
687		744
688		745
689	Negar Foroutan, Clara Meister, Debjit Paul, Joel Niklaus, Sina Ahmadi, Antoine Bosselut, and Rico Sennrich. 2025. Parity-aware byte-pair encoding: Improving cross-lingual fairness in tokenization. <i>arXiv preprint arXiv:2508.04796</i> .	746
690		747
691		748
692		749
693		750
694	Robert C Gale, Alexandra C Salem, Gerasimos Fergadiotis, and Steven Bedrick. 2023. Mixed orthographic/phonemic language modeling: Beyond orthographically restricted transformers (bort). In <i>Proceedings of the 8th Workshop on Representation Learning for NLP (RepL4NLP 2023)</i> , pages 212–225.	751
695		752
696		753
697		754
698		755
699		756
700	Juan Luis Gastaldi, John Terilla, Luca Malagutti, Brian DuSell, Tim Vieira, and Ryan Cotterell. 2024. The foundations of tokenization: Statistical and computational concerns. <i>arXiv preprint arXiv:2407.11606</i> .	757
701		758
702		759
703		760
704	Zébulon Goriely, Richard Diehl Martinez, Andrew Caines, Paula Buttery, and Lisa Beinborn. 2024. From babble to words: Pre-training language models on continuous streams of phonemes. In <i>The 2nd BabyLM Challenge at the 28th Conference on Computational Natural Language Learning</i> , pages 37–53.	761
705		762
706		763
707		764
708		765
709		766
710	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. <i>arXiv preprint arXiv:2407.21783</i> .	767
711		768
712		769
713		770
714		771
	Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. <i>spaCy: Industrial-strength Natural Language Processing in Python</i> .	715
		716
		717
	Jue Hou, Anisia Katinskaia, Anh-Duc Vu, and Roman Yangarber. 2023. Effects of sub-word segmentation on performance of transformer language models. <i>arXiv preprint arXiv:2305.05480</i> .	718
		719
		720
		721
	International Phonetic Association. 1999. <i>Handbook of the International Phonetic Association: A Guide to the Use of the International Phonetic Alphabet</i> . Cambridge University Press, Cambridge.	722
		723
		724
		725
	Julie Kallini, Dan Jurafsky, Christopher Potts, and Martijn Bartelds. 2025. False Friends are not foes: Investigating vocabulary overlap in multilingual language models. In <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> , pages 21138–21154, Suzhou, China. Association for Computational Linguistics.	726
		727
		728
		729
		730
		731
		732
	Julie Kallini, Shikhar Murty, Christopher D Manning, Christopher Potts, and Róbert Csordás. 2024. Mrt5: Dynamic token merging for efficient byte-level language models. <i>arXiv preprint arXiv:2410.20771</i> .	733
		734
		735
		736
	Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, and 1 others. 2007. Moses: Open source toolkit for statistical machine translation. In <i>Proceedings of the 45th annual meeting of the association for computational linguistics companion volume proceedings of the demo and poster sessions</i> , pages 177–180. Association for Computational Linguistics.	737
		738
		739
		740
		741
		742
		743
		744
		745
	Taku Kudo. 2018. Subword regularization: Improving neural network translation models with multiple subword candidates. <i>arXiv preprint arXiv:1804.10959</i> .	746
		747
		748
	Taku Kudo and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. <i>arXiv preprint arXiv:1808.06226</i> .	749
		750
		751
		752
	Jackson L Lee, Lucas FE Ashby, M Elizabeth Garza, Yeonju Lee-Sikka, Sean Miller, Alan Wong, Arya D McCarthy, and Kyle Gorman. 2020. Massively multilingual pronunciation modeling with wikipron. In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> , pages 4223–4228.	753
		754
		755
		756
		757
		758
	Colin Leong and Daniel Whitenack. 2022. Phone-ing it in: Towards flexible multi-modal language model training by phonetic representations of data. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5306–5315.	759
		760
		761
		762
		763
		764
	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In <i>Proceedings of the 58th annual meeting of</i>	765
		766
		767
		768
		769
		770

771	<i>the association for computational linguistics</i> , pages 7871–7880.	Daniel Zeman. 2020. Universal dependencies v2: An evergrowing multilingual treebank collection. <i>arXiv preprint arXiv:2004.10643</i> .	826
772			827
773	Tomasz Limisiewicz, Terra Blevins, Hila Gonen, Orevaoghene Ahia, and Luke Zettlemoyer. 2024. MYTE: Morphology-driven byte encoding for better and fairer multilingual language modeling . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15059–15076, Bangkok, Thailand. Association for Computational Linguistics.	Marta R. NLLB Team, Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Hefernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, and 3 others. 2024. Scaling neural machine translation to 200 languages . <i>Nature</i> , 630(8018):841–846.	828
774			829
775			830
776			831
777			832
778			833
779			834
780			835
781	Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation. <i>Transactions of the Association for Computational Linguistics</i> , 8:726–742.	OpenAI Team. 2023. Gpt-4 technical report . <i>Preprint</i> , arXiv:2303.08774.	836
782			837
783			838
784			839
785			840
786			841
787	Jonas F Lotz, António V Lopes, Stephan Peitz, Hendra Setiawan, and Leonardo Emili. 2025. Beyond text compression: Evaluating tokenizers across scales. <i>arXiv preprint arXiv:2506.03101</i> .	Artidoro Pagnoni, Ram Pasunuru, Pedro Rodriguez, John Nguyen, Benjamin Muller, Margaret Li, Chunting Zhou, Lili Yu, Jason Weston, Luke Zettlemoyer, and 1 others. 2024. Byte latent transformer: Patches scale better than tokens. <i>arXiv preprint arXiv:2412.09871</i> .	842
788			843
789			844
790			845
791	Jessica M Lundin, Ada Zhang, Nihal Karim, Hamza Louzan, Victor Wei, David Adelani, and Cody Carroll. 2025. The token tax: Systematic bias in multilingual tokenization. <i>arXiv preprint arXiv:2509.05486</i> .	Aleksandar Petrov, Emanuele La Malfa, Philip Torr, and Adel Bibi. 2023. Language model tokenizers introduce unfairness between languages. <i>Advances in neural information processing systems</i> , 36:36963–36990.	846
792			847
793			848
794			849
795	Chihaya Matsuhira, Marc A Kastner, Takahiro Komamizu, Takatsugu Hirayama, Keisuke Doman, Yasutomo Kawanishi, and Ichiro Ide. 2023. Ipa-clip: Integrating phonetic priors into vision and language pretraining. <i>arXiv preprint arXiv:2303.03144</i> .	Shengju Qian, Yi Zhu, Wenbo Li, Mu Li, and Jiaya Jia. 2022. What makes for good tokenizers in vision transformer? <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 45(11):13011–13023.	850
796			851
797			852
798			853
799			854
800	Clara Meister. 2025. Tokeval: A tokenizer analysis suite .	Libo Qin, Qiguang Chen, Yuhang Zhou, Zhi Chen, Yinghui Li, Lizi Liao, Min Li, Wanxiang Che, and Philip S Yu. 2025. A survey of multilingual large language models. <i>Patterns</i> , 6(1).	855
801			856
802	Sabrina J Mielke, Zaid Alyafeai, Elizabeth Salesky, Colin Raffel, Manan Dey, Matthias Gallé, Arun Raja, Chenglei Si, Wilson Y Lee, Benoît Sagot, and 1 others. 2021. Between words and characters: A brief history of open-vocabulary modeling and tokenization in nlp. <i>arXiv preprint arXiv:2112.10508</i> .	Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training . Ms, OpenAI.	857
803			858
804			859
805			860
806			861
807	David R Mortensen, Siddharth Dalmia, and Patrick Littell. 2018. Epitran: Precision g2p for many languages. In <i>Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)</i> .	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. Ms, OpenAI.	862
808			863
809			864
810			865
811			866
812	Hoang H Nguyen, Khyati Mahajan, Vikas Yadav, Julian Salazar, Philip S Yu, Masoud Hashemi, and Rishabh Maheshwary. 2024. Prompting with phonemes: Enhancing llms’ multilinguality for non-latin script languages. <i>arXiv preprint arXiv:2411.02398</i> .	Bharath Raj, Garvit Suri, Vikrant Dewangan, and Raghav Sonavane. 2024. When every token counts: Optimal segmentation for low-resource language models. <i>arXiv preprint arXiv:2412.06926</i> .	867
813			868
814			869
815			870
816			871
817	Thuat Nguyen, Chien Van Nguyen, Viet Dac Lai, Hieu Man, Nghia Trung Ngo, Franck Dernoncourt, Ryan A Rossi, and Thien Huu Nguyen. 2023. Culturax: A cleaned, enormous, and multilingual dataset for large language models in 167 languages. <i>arXiv preprint arXiv:2309.09400</i> .	Nived Rajaraman, Jiantao Jiao, and Kannan Ramchandran. 2024. Toward a theory of tokenization in llms. <i>arXiv preprint arXiv:2404.08335</i> .	872
818			873
819			874
820			875
821			876
822			877
823	Joakim Nivre, Marie-Catherine De Marneffe, Filip Ginter, Jan Hajič, Christopher D Manning, Sampo Pyysalo, Sebastian Schuster, Francis Tyers, and	Mike Schuster and Kaisuke Nakajima. 2012. Japanese and korean voice search. In <i>2012 IEEE international conference on acoustics, speech and signal processing (ICASSP)</i> , pages 5149–5152. IEEE.	878
824			879
825			

880	Rico Sennrich, Barry Haddow, and Alexandra Birch.	Megabyte: Predicting million-byte sequences with	936
881	2016. Neural machine translation of rare words with	multiscale transformers. <i>Advances in Neural Infor-</i>	937
882	subword units . In <i>Proceedings of the 54th Annual</i>	<i>mation Processing Systems</i> , 36:78808–78823.	938
883	<i>Meeting of the Association for Computational Lin-</i>		
884	<i>guistics (Volume 1: Long Papers)</i> , pages 1715–1725,	Crystina Zhang, Jing Lu, Vinh Q. Tran, Tal Schuster,	939
885	Berlin, Germany. Association for Computational Lin-	Donald Metzler, and Jimmy Lin. 2025. Tomato, tom-	940
886	guistics.	ahto, tomato: Do multilingual language models un-	941
		derstand based on subword-level semantic concepts?	942
887	Kevin Slagle. 2024. Spacebyte: Towards deleting	In <i>Findings of the Association for Computational</i>	943
888	tokenization from large language modeling. <i>Ad-</i>	<i>Linguistics: NAACL 2025</i> , pages 1821–1837, Albu-	944
889	<i>vances in Neural Information Processing Systems</i> ,	querque, New Mexico. Association for Computational	945
890	37:124925–124950.	Linguistics.	946
891	Jimin Sohn and David R Mortensen. 2025. Cross-	Xin Zhang, Dong Zhang, Shimin Li, Yaqian Zhou, and	947
892	lingual ipa contrastive learning for zero-shot ner.	Xipeng Qiu. 2023. Speechook: Unified speech	948
893	<i>arXiv preprint arXiv:2503.07214</i> .	tokenizer for speech large language models. <i>arXiv</i>	949
		<i>preprint arXiv:2308.16692</i> .	950
894	Pedro Javier Ortiz Suárez, Laurent Romary, and Benoît		
895	Sagot. 2020. A monolingual approach to contextu-	A International Phonetic Alphabet (IPA)	951
896	alized word embeddings for mid-resource languages.		
897	<i>arXiv preprint arXiv:2006.06202</i> .	The International Phonetic Alphabet (IPA) is a sys-	952
		tem of phonetic notation developed and maintained	953
898	Pedro Javier Ortiz Suárez, Benoît Sagot, and Laurent	by the International Phonetic Association to provide	954
899	Romary. 2019. Asynchronous pipeline for process-	a standardized, language-independent representa-	955
900	ing huge corpora on medium to low resource infras-	tion of the sounds of spoken language. Originally	956
901	tructures. In <i>7th Workshop on the Challenges in the</i>	created in the late 19th century, the IPA has under-	957
902	<i>Management of Large Corpora (CMLC-7)</i> . Leibniz-	gone multiple refinements over the years, with the	958
903	Institut für Deutsche Sprache.	most recent official chart standardized in 2020. ⁶	959
		The IPA uses a set of dedicated symbols ⁷ to	960
904	Yi Tay, Vinh Q Tran, Sebastian Ruder, Jai Gupta,	represent speech sounds, distinguishing four main	961
905	Hyung Won Chung, Dara Bahri, Zhen Qin, Si-	categories: consonants, vowels, diacritics, and	962
906	mon Baumgartner, Cong Yu, and Donald Metzler.	suprasegmentals. Consonant and vowel symbols	963
907	2021. Charformer: Fast character transformers via	correspond to specific phonemes and are organized	964
908	gradient-based subword tokenization. <i>arXiv preprint</i>	according to articulatory features such as place and	965
909	<i>arXiv:2106.12672</i> .	manner of articulation for consonants, and tongue	966
		height and backness for vowels. Diacritics are used	967
910	Linting Xue, Aditya Barua, Noah Constant, Rami Al-	to modify base symbols to indicate finer phonetic	968
911	Rfou, Sharan Narang, Mihir Kale, Adam Roberts,	details such as nasalization, aspiration, or devoicing.	969
912	and Colin Raffel. 2022. Byt5: Towards a token-free	Suprasegmentals, such as stress, length, intonation,	970
913	future with pre-trained byte-to-byte models. <i>Transac-</i>	and tone, are represented with additional symbols	971
914	<i>tions of the Association for Computational Linguis-</i>	and markings to capture prosodic aspects of speech.	972
915	<i>tics</i> , 10:291–306.	Notably, there are two common modes of IPA	973
		transcriptions—broad and narrow—which differ in	974
916	Linting Xue, Noah Constant, Adam Roberts, Mihir Kale,	the level of phonetic detail. Broad transcriptions	975
917	Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and	are <i>phonemic</i> , in that they only include contrastive	976
918	Colin Raffel. 2021. mt5: A massively multilin-	sounds needed to distinguish meaning, while nar-	977
919	gual pre-trained text-to-text transformer . <i>Preprint</i> ,	row transcriptions are <i>allophonic</i> , encoding more	978
920	<i>arXiv:2010.11934</i> .	fine-grained information about phonetic realization	979
		and variation. For example, a broad transcription	980
921	Yinfei Yang, Yuan Zhang, Chris Tar, and Jason	might not differentiate between the sound corre-	981
922	Baldrige. 2019a. PAWS-X: A cross-lingual adversarial	sponding to the letter “l” in “leaf” and “pool” (/lif/	982
923	dataset for paraphrase identification . In <i>Proceed-</i>		
924	<i>ings of the 2019 Conference on Empirical Methods</i>	⁶ https://www.internationalphoneticassociation.	
925	<i>in Natural Language Processing and the 9th Inter-</i>	org/content/full-ipa-chart	
926	<i>national Joint Conference on Natural Language Pro-</i>	⁷ The exact number varies by source, depending on IPA	
927	<i>cessing (EMNLP-IJCNLP)</i> , pages 3687–3692, Hong	variants and criteria for what counts as a single character, but	
928	Kong, China. Association for Computational Linguis-	typically ranges between 100 and 200 unique symbols.	
929	tics.		
930	Yinfei Yang, Yuan Zhang, Chris Tar, and Jason		
931	Baldrige. 2019b. Paws-x: A cross-lingual adver-		
932	sarial dataset for paraphrase identification . <i>arXiv</i>		
933	<i>preprint arXiv:1908.11828</i> .		
934	Lili Yu, Dániel Simig, Colin Flaherty, Armen Agha-		
935	janyan, Luke Zettlemoyer, and Mike Lewis. 2023.		

983 and /pul/), while a narrow transcription would ([li:f] 1025
984 and [pu:t])^{8,9} For more information about IPA, 1026
985 one can consult the detailed handbook provided by 1027
986 the [International Phonetic Association \(1999\)](#). 1028

987 B Tokenization 1030

988 Tokenization—the process of dividing a sequence 1031
989 of text¹⁰ into smaller, discrete units (i.e. *to-* 1032
990 *kens*)—has been an active area of research in recent 1033
991 years ([Hou et al., 2023](#); [Domingo et al., 2019](#); [Ra-](#) 1034
992 [jaraman et al., 2024](#), *inter alia*). It comes as the 1035
993 first component in the process of training an LLM, 1036
994 and has been shown to have a significant impact on 1037
995 model performance and generalization, especially 1038
996 in multilingual settings ([Chelombitko et al., 2024](#)). 1039

997 Tokenization methods can be systematically clas- 1040
998 sified along two axes: (1) the approach to deter- 1041
999 mining token boundaries—*rule-based* versus *data-* 1042
1000 *driven*; and (2) the granularity of the resulting to- 1043
1001 kens—*word-level*, *subword-level*, and *character-* 1044
1002 *or byte-level* ([Gastaldi et al., 2024](#)). Rule-based 1045
1003 tokenization relies on predefined, usually linguisti- 1046
1004 cally motivated heuristics, while data-driven ap- 1047
1005 proaches learn token boundaries directly from data, 1048
1006 typically optimizing for frequency, likelihood, or 1049
1007 other corpus-derived statistics. Regarding granu- 1050
1008 larity, in word-level tokenization, each token cor- 1051
1009 responds to a complete word form. Subword-level 1052
1010 tokenization segments text into units smaller than 1053
1011 or equal to full words. Character- and byte-level 1054
1012 methods treat each character or byte as a separate to- 1055
1013 ken, respectively. Table 2 categorizes some notable 1056
1014 tokenizers along these axes. 1057

1015 Early NLP systems relied predominantly on 1058
1016 rule-based, word-level tokenization. However, 1059
1017 this approach encounters significant limitations 1060
1018 such as rapid vocabulary growth, frequent out- 1061
1019 of-vocabulary (OOV) issues, and necessity for 1062
1020 language-specific hand-crafted rules that often re- 1063
1021 quire linguistic expertise ([Mielke et al., 2021](#)). Data- 1064
1022 driven, subword-level tokenizers emerged to ad- 1065
1023 dress these issues, offering several key advantages. 1066
1024 They significantly reduce vocabulary size, improve

⁸It would also not differentiate vowel length (indicated by “:”), unless this feature is contrastive (i.e. changes word meaning) in the particular language.

⁹Here we adopt the notation from the official IPA Handbook of enclosing phonemes within forward slashes, i.e. “/.../”, and phones within square brackets, i.e. “[...]”

¹⁰While we focus on text tokenization here, the concept applies more broadly to different types of sequential data e.g., audio stream in speech models ([Zhang et al., 2023](#)) or image pixels in computer vision models ([Qian et al., 2022](#)).

1025 generalization, and efficiently handle OOV terms 1026
1027 by supporting open-vocabulary modeling ([Gastaldi 1028
et al., 2024](#)). However, subword tokenization is not 1029
1030 without drawbacks. Most notably for our discus- 1031
1032 sion, these methods might neglect meaningful lin- 1033
1034 guistic structures ([Bostrom and Durrett, 2020](#); [Bein- 1035
born and Pinter, 2023](#)) and often unfairly allocate 1036
1037 vocabulary space, disproportionately favoring high- 1038
1039 resource languages in the training data ([Qin et al., 1040
2025](#); [Mielke et al., 2021](#); [Rust et al., 2020](#)). Re- 1041
1042 cently, fully character- and byte-level methods (also 1043
1044 collectively known as “tokenization-free” meth- 1045
1046 ods) gained attention as alternatives. They entirely 1047
1048 avoid vocabulary constraints, providing universal 1049
1050 coverage and script independence, alongside ro- 1051
1052 bustness to noise ([Xue et al., 2022](#)). Yet, these 1053
1054 methods also face certain limitations. Firstly, they 1055
1056 significantly increase computational cost due to 1057
1058 substantially longer input sequences ([Kallini et al., 1059
2024](#)). Secondly, because they operate on raw bytes, 1060
1061 these models inherit a non-linguistically motivated 1062
1063 Unicode¹¹-based scheme. In this scheme, every 1064
1065 Latin character is encoded in a single byte, whereas 1066

1067 While current research aims to resolve the chal- 1068
1069 lenges highlighted above, current state-of-the-art 1070
1071 LLMs and MLMs predominantly use data-driven 1072
1073 subword tokenizers—most notably, Byte-Pair En- 1074
1075 coding (used by, e.g., GPT-4 ([OpenAI Team, 2023](#)) 1076
1077 and LLaMA-3 ([Grattafiori et al., 2024](#))) and Uni- 1078
1079 gramLM (used by, e.g., mT5 ([Xue et al., 2021](#)) and 1079
1080 mBART ([Liu et al., 2020](#))). It is for that reason that 1080
1081 we explore these particular two tokenizer methods 1081
1082 in our experiments. 1082

1066 C Data and languages 1066

1067 C.1 Mapping of languages and scripts to their 1068 ISO codes 1068

1069 C.2 Data sampling strategies 1069

1070 We visualize the resulting language distributions of 1070
1071 the four different data sampling strategies in Fig. 7. 1071

¹¹<https://home.unicode.org/>

Word-level	Subword-level	Character-/ byte-level
Moses* (Koehn et al., 2007)	Morfessor† (Creutz and Lagus, 2005)	Charformer (Tay et al., 2021)
NLTK* (Bird et al., 2009)	WordPiece† (Schuster and Nakajima, 2012)	ByT5 (Xue et al., 2022)
spaCy* (Honnibal et al., 2020)	BPE† (Sennrich et al., 2016)	CANINE (Clark et al., 2022)
	UnigramLM† (Kudo, 2018)	MegaByte (Yu et al., 2023)
		SpaceByte (Slagle, 2024)
		BLT (Pagnoni et al., 2024)
		MrT5 (Kallini et al., 2024)

Table 2: Representative tokenizers grouped by granularity. Asterisk symbol (*) marks rule-based tokenizers, and dagger symbol (†) data-driven methods. Character-/byte-level column names tokenizer-free models that operate on raw characters/bytes using different approaches.

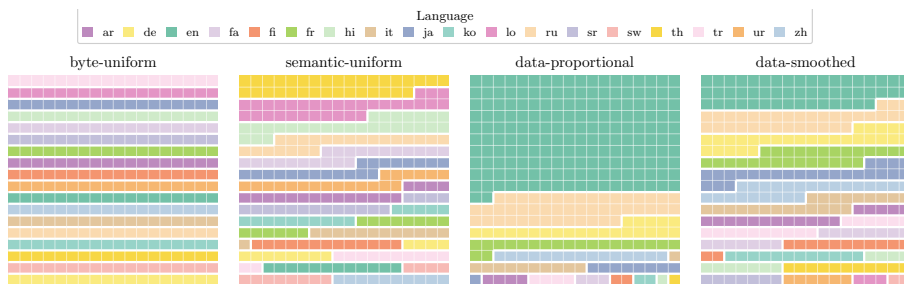


Figure 7: Resulting distribution of languages in the training data using four different sampling strategies. Each square represents $\approx 0.3\%$ of the training data.

1072 C.3 WikiPron details

1073 Table 4 shows WikiPron dataset details by language
1074 on the sample that we use in this work.

1075 D Epitran usage details

1076 A typical usage involves instantiating an Epitran
1077 object with a language–script code (e.g.,
1078 `epitran.Epitran('eng-Latn')` for English)
1079 and then calling its `transliterate()` method,
1080 which returns the IPA representation. Internally,
1081 Epitran relies on language-specific mapping files
1082 and optional pre- and post-processing rules. An al-
1083 ternative utility is the `Backoff` class, which allows
1084 cascading multiple language–script mappings (for
1085 example, to handle mixed-script inputs by falling
1086 back from one transliterator to another). However,
1087 the `Backoff` class currently does not support the
1088 full range of text pre- and post-processing routines
1089 applied by the standard `transliterate()`, lead-
1090 ing us adopt a different transliteration approach. In
1091 this section, we describe this approach, as well as
1092 other modifications that we made to this tool.

1093 **Bug fixes and mapping extensions.** We cor-
1094 rected identified bugs in the Serbian mapping and
1095 added several missing mappings for Urdu, Arabic,
1096 Persian, Hindi, Finnish, and Turkish. These addi-

tions were based on empirically observed character
sequences occurring above 1% frequency in our
training data and consulted via well-documented
Wikipedia pages of the respective languages. The
exact modifications we made per aforementioned
language will be accessible in our codebase, once
published.

Efficient processing for Chinese and Japanese.

1104 While most languages were processed by applying
1105 `transliterate()` *once* per full document, this ap-
1106 proach proved inefficient for Chinese and Japanese
1107 due their special treatment by Epitran under the
1108 hood, which resulted in a non-linear runtime in-
1109 creases with increased input length (see Figure 8).
1110 To address this, we split documents into smaller
1111 segments along Chinese/Japanese specific punctua-
1112 tion, ensuring word boundaries were preserved.
1113 This chunking preserved performance while main-
1114 taining accuracy. In contrast, this segmentation
1115 approach notable overhead—for other languages
1116 (see Figure 9), so it was only applied to Chinese
1117 and Japanese. 1118

1119 **Sequential mapping for multiple variants.** For
1120 languages supported by multiple mapping files we
1121 employed a specific strategy: we applied the first
1122 mapping to the full text, then identified characters

Language	ISO 639-1	Script	ISO 15924
Arabic	ar	Arabic	Arab
Amharic	am	Ethiopic	Ethi
Burmese	my	Myanmar	Mymr
Chinese	zh	Han	Hani
English	en	Latin	Latn
Finnish	fi	Latin	Latn
French	fr	Latin	Latn
German	de	Latin	Latn
Greek [†]	el	Greek	Grek
Hebrew [†]	he	Hebrew	Hebr
Hindi	hi	Devanagari	Deva
Italian	it	Latin	Latn
Japanese*	ja	Japanese	Jpan
Korean	ko	Hangul	Hang
Lao	lo	Lao	Laoo
Persian	fa	Arabic	Arab
Polish	pl	Latin	Latn
Russian	ru	Cyrillic	Cyrl
Serbian*	sr	Latin/Cyrillic	Latn/Cyrl
Spanish	es	Latin	Latn
Swahili	sw	Latin	Latn
Thai	th	Thai	Thai
Turkish	tr	Latin	Latn
Urdu	ur	Arabic	Arab

Table 3: Languages used in this work. Each language is followed by its two-letter ISO 639-1 code, its script, and the script’s four-letter ISO 15924 code. Languages marked with an asterisk (*) use multiple scripts, while languages not supported by Epitran are marked with a dagger symbol (†).

still untransliterated (i.e., non-IPA), and selectively applied subsequent mappings only to those residual characters. This procedure ensured each script variant was handled thoroughly, while avoiding unnecessary repeated processing.

Post-conversion cleanup. After IPA conversion, any remaining characters not recognized or mapped by Epitran (primarily external script tokens included in CulturaX data for the particular language, but also, at times, characters not supported by Epitran) were removed. Across languages, this accounted for an average of only 0.3% of tokens per language, a removal rate negligible in volume yet important for ensuring consistency and purity of the phonological representation space.

E Experimental setup

E.1 Tokenizer training hyper-parameters

F Intrinsic tokenization metrics

F.1 Compression

Let W_ℓ be the set of unique word types for language ℓ , and let $\tau(\cdot)$ be a tokenizer mapping an input string to a token sequence. Denote the num-

ber of produced tokens by $n_\tau(x) = |\tau(x)|$ and the character length of a string by $|x|$ (in the tokenizer’s input representation).

Word Fertility (WF). Average number of subword tokens per word type (lower is better):

$$\text{WF}_\ell = \frac{1}{|W_\ell|} \sum_{w \in W_\ell} n_\tau(w). \quad (1)$$

Proportion of Continued Words (PCW). Fraction of word types split into multiple tokens (lower is better):

$$\text{PCW}_\ell = \frac{1}{|W_\ell|} \sum_{w \in W_\ell} \mathbb{I}[n_\tau(w) > 1]. \quad (2)$$

Average Token Length (ATL). Average length of produced tokens, measured in characters of the tokenizer input representation (**Text** or **IPA**); higher indicates longer token segments. Let $\mathcal{T}_\ell = \bigsqcup_{w \in W_\ell} \tau(w)$ be the multiset of all tokens produced for W_ℓ , and let $|t|$ be the character length of token t (excluding any boundary marker, e.g. leading “_”):

$$\text{ATL}_\ell = \frac{1}{|\mathcal{T}_\ell|} \sum_{t \in \mathcal{T}_\ell} |t|. \quad (3)$$

Compression Rate (CR). Average number of raw input characters per token at the sentence level (higher indicates stronger compression). For a sentence set S_ℓ :

$$\text{CR}_\ell = \frac{1}{|S_\ell|} \sum_{s \in S_\ell} \frac{|s|}{n_\tau(s)}. \quad (4)$$

F.2 Token frequency distribution shape

Let $c(v)$ be the count of token type v on evaluation data for language ℓ , with total token count $N = \sum_v c(v)$ and empirical probabilities $p(v) = c(v)/N$.

Rényi Entropy (RE). Measures how evenly token usage is distributed across the vocabulary:

$$H_\alpha = \frac{1}{1 - \alpha} \log \sum_v p(v)^\alpha. \quad (5)$$

We report $\alpha = 1$ (Shannon entropy), $\alpha = 2$ (collision entropy), and $\alpha = \infty$ (min-entropy). Higher values indicate a more uniform distribution.

Seen Languages (18)						
Language	ISO	Original	After Removal	Scope	Filtered	
Arabic	ar	2 252	1 596	Broad	No	
German	de	47 779	42 914	Broad	Yes	
English	en	81 576	69 639	Broad	Yes	
Persian	fa	34 033	7 357	Narrow	No	
Finnish	fi	158 880	156 923	Broad	No	
French	fr	80 690	71 229	Broad	Yes	
Hindi	hi	24 640	15 867	Broad	Yes	
Italian	it	79 865	73 136	Broad	Yes	
Japanese*	ja	32 749	29 092	Narrow	Yes	
Korean	ko	22 072	21 383	Narrow	Yes	
Lao	lo	4 180	2 085	Narrow	No	
Russian	ru	411 651	395 845	Narrow	No	
Serbian*	sr	46 991	45 553	Broad	Yes	
Swahili	sw	110	106	Broad	No	
Thai	th	16 689	11 589	Broad	No	
Turkish	tr	7 266	6 933	Broad	No	
Urdu	ur	4 493	3 637	Broad	No	
Chinese	zh	158 873	126 804	Broad	No	
Total	—	1 214 798	1 081 688			

Unseen Languages (6)						
Language	ISO	Original	After Removal	Scope	Filtered	
Amharic	am	378	371	Broad	No	
Greek	el	14 825	14 825	Broad	Yes	
Spanish	es	99 043	98 787	Broad	Yes	
Hebrew	he	1 957	1 945	Broad	No	
Burmese	my	6 062	4 486	Broad	Yes	
Polish	pl	132 558	115 865	Broad	No	
Total	—	256 608	238 057			

Table 4: WikiPron details by language. “Original” represents the total number of words in the original dataset; “After Removal” shows the number of remaining words after deduplication. “Scope” and “Filtered” columns direct to the specific WikiPron file that was used per language (we always selected *broad* and *filtered* versions, unless unavailable). For languages marked with an asterisk (*), we include different scripts (Katakana and Hiragana for Japanese; Latin and Cyrillic for Serbian).

Zipf Deviation (ZipfD). Following Lotz et al. (2025), we quantify how closely the empirical token rank–frequency distribution follows a Zipfian power law by fitting a linear function to the log–log curve and reporting its mean absolute deviation. Concretely, let tokens be sorted by descending frequency, with rank $r \in \{1, \dots, R\}$ and frequency f_r . Define $x_r = \log r$ and $y_r = \log f_r$, and estimate parameters β_0, β_1 via least squares for $f(x) = \beta_0 + \beta_1 x$. Zipf deviation is then

$$\text{ZipfD} = \frac{1}{R} \sum_{r=1}^R |\beta_0 + \beta_1 x_r - y_r|. \quad (6)$$

Lower ZipfD indicates closer agreement with Zipf’s law.

F.3 Vocabulary usage

Using the same notation above, let V be the learned tokenizer vocabulary with size $|V|$, and let $U =$

$\{v \in V : c(v) > 0\}$ be the set of observed token types.

Vocabulary Utilization (VU). Share of the learned vocabulary used at least once on evaluation data (higher is better):

$$\text{VU} = \frac{|U|}{|V|}. \quad (7)$$

Type–Token Ratio (TTR). Distinct token types relative to the total number of produced tokens (higher indicates more diverse usage):

$$\text{TTR} = \frac{|U|}{N}. \quad (8)$$

F.4 Cross-lingual equity

Let S_ℓ^{\parallel} be a parallel sentence set paired with English, containing pairs $(s_{\text{en}}^{(i)}, s_\ell^{(i)})$.

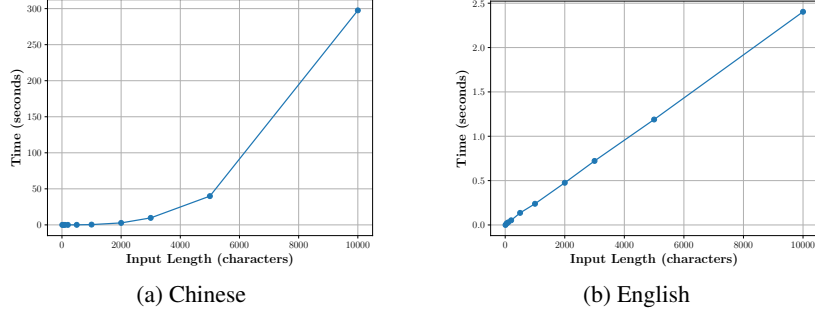


Figure 8: Runtime vs. Input length efficiency trade-off for (a) Chinese and (b) English.

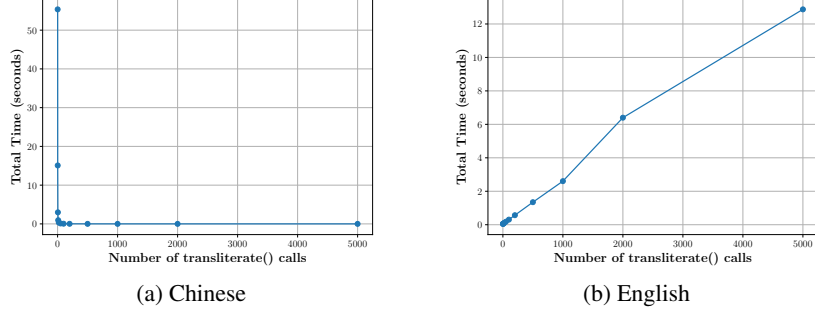


Figure 9: Runtime vs. Number of `transliterate()` calls efficiency trade-off for (a) Chinese and (b) English.

Tokenization Parity (TP) Following the discussion of Petrov et al. (2023), we compare tokenization length relative to English on parallel sentences:

$$TP_\ell = \frac{1}{|S_\ell^{\parallel}|} \sum_i \frac{n_\tau(s_\ell^{(i)})}{n_\tau(s_{en}^{(i)})}. \quad (9)$$

Values closer to 1 indicate more similar tokenization length relative to English.

Tokenization Fairness Gini (TFG). Following Meister (2025), we measure inequality in tokenization cost across languages using the (equal-weight) Gini coefficient. Let $\mathcal{L} = \{1, \dots, n\}$ be the set of languages. For each language $\ell \in \mathcal{L}$, define the token cost as

$$c_\ell = \frac{\sum_{s \in S_\ell} n_\tau(s)}{\sum_{s \in S_\ell} |s|_{\text{bytes}}}, \quad (10)$$

where $n_\tau(s)$ is the number of tokens produced for sentence s and $|s|_{\text{bytes}}$ is its UTF-8 byte length. The mean cost is

$$\mu = \frac{1}{n} \sum_{\ell=1}^n c_\ell. \quad (11)$$

TFG is then

$$TFG = \frac{\sum_{i=1}^n \sum_{j=1}^n |c_i - c_j|}{2n^2\mu}. \quad (12)$$

Lower TFG indicates more equal byte-normalized token costs across languages.

G Tokenizer selection

Training autoregressive LMs for all tokenizer configurations is computationally prohibitive, so we use intrinsic evaluation to select a small set of tokenizers for GPT-2 pre-training. This appendix formalizes the ranking procedure and reports the full ranking heatmaps (Fig. ??–Fig. ??).

Set-up. Let $R \in \{\text{TEXT}, \text{IPA}\}$ denote the input representation, and let \mathcal{C}_R be the set of tokenizer configurations evaluated under R (here $|\mathcal{C}_R| = 32$). Each configuration $c \in \mathcal{C}_R$ is defined by the tuple

$$c = (\text{algorithm}, \text{vocab size}, \text{sampling}).$$

Let \mathcal{L} be the set of evaluation languages and let \mathcal{M} be the set of intrinsic metrics (ten metrics in our study). For each language $\ell \in \mathcal{L}$, metric $m \in \mathcal{M}$, and configuration c , we obtain a per-language score $s_m(c, \ell)$ (as defined in §4.4).

Macro-averaged metric scores. To obtain a single score per metric and configuration, we macro-average over languages:

$$\bar{s}_m(c) = \frac{1}{|\mathcal{L}|} \sum_{\ell \in \mathcal{L}} s_m(c, \ell). \quad (13)$$

For metrics where “lower is better” (e.g., WF, PCW, ZipfD, TFG), we convert them to a consistent

Mode	Parameter	Value
Text	model_type	[bpe, unigram]
	vocab_size	[40 000, 80 000, 100 000, 200 000]
	character_coverage	0.9995
	split_by_unicode_script	True
	byte_fallback	True
	normalization_rule_name	nmt_nfkc
IPA	model_type	[bpe, unigram]
	vocab_size	[40 000, 80 000, 100 000, 200 000]
	character_coverage	1.0
	split_by_unicode_script	False
	byte_fallback	False
	normalization_rule_name	identity

Table 5: Tokenizer training hyper-parameter settings used for Text and IPA representations. We show only the relevant parameters, for all others we used the default SentencePiece values for both modes. Values enclosed in square brackets signify that multiple configurations were considered.

“higher is better” orientation by defining

$$\tilde{s}_m(c) = \begin{cases} \bar{s}_m(c) & \text{if higher is better for } m, \\ -\bar{s}_m(c) & \text{if lower is better for } m. \end{cases} \quad (14)$$

For TP (where values closer to 1 are preferred), we use

$$\tilde{s}_{\text{TP}}(c) = -|\bar{s}_{\text{TP}}(c) - 1|. \quad (15)$$

Per-metric ranks and mean rank aggregation.

For each metric m and representation R , we rank configurations by their oriented macro score $\tilde{s}_m(c)$:

$$r_m(c) = 1 + \sum_{c' \in \mathcal{C}_R} \mathbb{I}[\tilde{s}_m(c') > \tilde{s}_m(c)], \quad c \in \mathcal{C}_R, \quad (16)$$

where $r_m(c) = 1$ denotes the best configuration for metric m (ties can be handled by average ranks; in practice ties are rare given continuous scores). We aggregate metric-wise ranks into a single scalar score by averaging across metrics:

$$\text{MR}(c) = \frac{1}{|\mathcal{M}|} \sum_{m \in \mathcal{M}} r_m(c). \quad (17)$$

Lower $\text{MR}(c)$ indicates a better overall intrinsic ranking. The full rank matrices $\{r_m(c)\}$ and mean-rank scores $\text{MR}(c)$ are visualized in Fig. 12a–Fig. 12b.

Selecting Opt and Subopt tokenizers. For each representation R , we select the intrinsically best tokenizer as

$$c_R^* = \arg \min_{c \in \mathcal{C}_R} \text{MR}(c). \quad (18)$$

We refer to the resulting tokenizers as *Text Opt* and *IPA Opt*. In our experiments, *Text Opt* corresponds to [BPE, 200k, data-proportional],

whereas *IPA Opt* corresponds to [UnigramLM, 200k, byte-uniform]. Because these optima arise under different hyperparameter settings, they differ along both representation and tokenizer configuration axes. To control for this, we also select cross-swapped “Subopt” tokenizers: a Text tokenizer trained with the IPA Opt configuration (*Text Subopt*) and an IPA tokenizer trained with the Text Opt configuration (*IPA Subopt*). We then pre-train four GPT-2 models in total: Text Opt, Text Subopt, IPA Opt, and IPA Subopt.

H GPT-2 training and fine-tuning hyper-parameters

I General benefits of IPA as an input representation

I.1 General benefits of IPA as an input representation

In the introduction section, we have outlined several beneficial properties offered by the IPA. Here we show what implications those properties have in practice.

IPA improves cross-lingual sharing of character inventories. In Figure 10 we show the percentage of shared characters across each pair of our 18 languages on the CulturaX sample we used for pre-training on the original text and on the converted IPA text. In Text heatmap, we see biggest sharing between languages that use the same script, but that percentage drops for language pairs with different scripts. We see that languages with complex orthographies (Chinese, Japanese, and Korean) have almost zero overlap with any other language. Other languages do show some percentage of over-

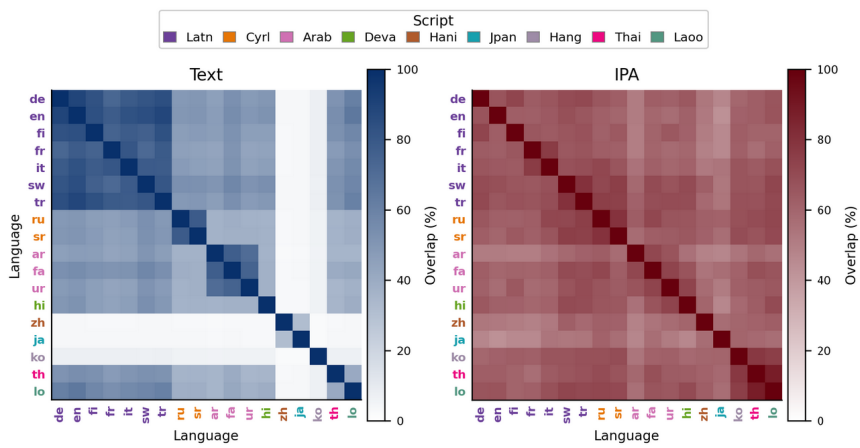


Figure 10: Pairwise overlap of character inventories (%) between the 18 pre-training languages on the CulturaX byte-uniform sample, computed on the original **Text** input (left) and its **IPA** transcription (right). Each cell reports the percentage of unique characters shared by a language pair (diagonal = 100%). In **Text**, overlap is largely confined to languages that share a script, with near-zero sharing for languages with distinct orthographies (e.g., ZH, JA, KO); in **IPA**, overlap becomes higher across languages and more uniform across scripts due to the shared phonemic symbol inventory.

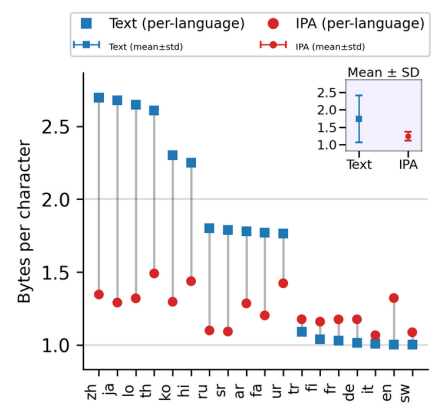


Figure 11: Average bytes-per-character (BPC) by language on the CulturaX byte-uniform sample, for the original **Text** input and its **IPA** transcription. **Text** exhibits large cross-lingual disparities, with higher BPC for languages whose scripts occupy higher Unicode ranges (e.g., ZH, JA, LO, TH), while Latin-script languages cluster near 1 BPC. **IPA** concentrates most symbols in low-byte code points, reducing mean BPC and making storage cost more uniform across languages; the inset (top-right) summarizes the mean \pm standard deviation across languages.

lap even though they differ in scripts, partly because many languages in our sample do contain inputs from other languages, most notably English. However, the IPA heatmap reveals a much more evenly distributed and overall larger character overlap across languages. This is a direct consequence of IPA’s shared input representation for all of the languages. Overall, this increased sharing is beneficial for multilingual tokenizers, as cross-lingual

token overlap has been shown to improve model performance (Zhang et al., 2025; Kallini et al., 2025).

IPA reduces storage overhead across languages on average. Figure 11 compares the average number of bytes needed to encode a single character per language for Text vs IPA, estimated from our pre-training CulturaX sample. Standard text representation shows notable disparities: languages with complex orthographies (Chinese, Japanese, Lao, and Thai) require much larger amount (> 2.5) of bytes-per-character (BPC), while languages using Latin script are all very close to 1 byte-per-character. This difference is due to the position of different characters in the Unicode range. The symbols used by the IPA mostly use ASCII characters (1 BPC), together with a set of special IPA symbols (2 BPC). Some additional, but less frequent IPA markers require 3 BPC. This fact, together with IPA’s shared language representation space, reduces the average BPC across our set of languages, requiring slightly larger numbers for Latin-script languages than Text, but much lower BPC for all other ones. In addition, the required storage memory is much more balanced across languages.

J Per-language results for suboptimal tokenizers

We provide a detailed per-language intrinsic evaluation of the suboptimal tokenizers in Fig. 13.

1320
1321

1322
1323
1324
1325
1326
1327
1328
1329

1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343

1344
1345

1346
1347

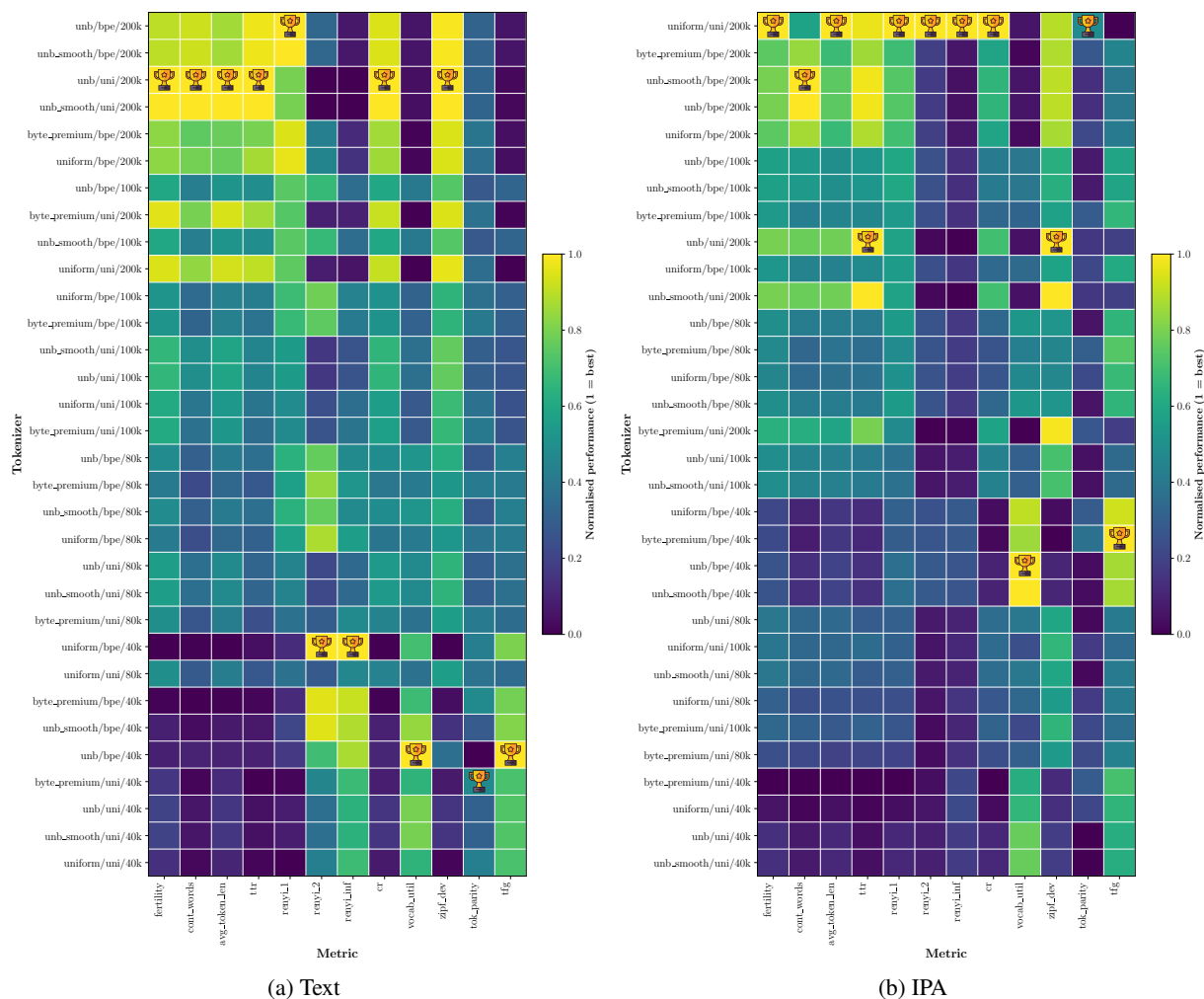


Figure 12: Tokenizer selection. We rank Text and IPA tokenizers to select the optimal settings for GPT-2 pre-training. Optimal setting for Text is [BPE, 200k, data-smoothed] and for IPA it is [UnigramLM, 200k, byte-uniform], according to the mean rank across metrics and languages.

K Detailed per-script daignostic

We provide a detailed per-script intrinsic evaluation for both the optimal (Fig. 15) and suboptimal tokenizers ().

L Interpretation of tokenizer configuration effects

Compression metrics improve with vocab size.

This is consistent across algorithms, sampling strategies, and modes (text/ipa). Consequently, vocab size of 200k performs best in terms of compression (WF, PCW, ATL, CR). This is expected, as larger vocabs will have more space for storing tokens, which means fewer tokens needed to represent the same input.

Renyi H_2 and H_∞ best at smallest vocab; H_1 best at 200k. This pattern is consistent with the different sensitivities of R’enyi orders: H_1 (Shannon

entropy) reflects the full distribution and increases when the vocabulary grows to include many rare, low-probability types (e.g., longer tokens approaching full words). In contrast, H_2 and especially H_∞ place increasing weight on the head of the distribution and thus tend to favor settings where probability mass is less concentrated in a small set of very frequent tokens—often achieved by smaller vocabularies that encourage greater sharing of subword units. Because the entropy scale depends strongly on vocabulary size, we treat these trends primarily as descriptive of configuration effects and focus comparisons on the Text vs. IPA axis under matched vocabulary sizes. One exception is [IPA, UnigramLM, byte-uniform], where all three entropies are maximized at 200k; we leave a more detailed analysis of this case to future work.

Component	Hyperparameter	Value
Data & packing	Context length (block size)	2048
	Vocabulary size	200,000
	BOS / EOS token id	2 / 2
Model (GPT-2)	Layers / heads / hidden size	12 / 12 / 768
	Attention implementation	SDPA
Optimization	Training steps	10,000
	Per-device batch size	8
	Gradient accumulation steps	8
	Effective batch size (seq./device)	64
	Tokens per optimizer step	131,072
	Learning rate	5×10^{-4}
	LR scheduler	cosine
	Warmup ratio	0.03
	Weight decay	0.1
Precision	bfloat16	
Eval & logging	Evaluation cadence	every 500 steps
	Logging cadence	every 50 steps
	Checkpointing	none (final model only)
Reproducibility	Random seed	42

Table 6: GPT-2 pretraining hyperparameters shared across Text Opt/Subopt and IPA Opt/Subopt.

Hyperparameter	XNLI	PAWS-X
Training language	English only	English only
Number of labels	3	2
Max steps	6128	3860
Train batch size (per device)	32	32
Gradient accumulation steps	8	8
Effective batch size (seq./device)	256	256
Learning rate	1×10^{-4}	3×10^{-5}
LR scheduler	cosine	cosine
Warmup ratio	0.06	0.06
Weight decay	0.0	0.0
Max grad norm	1.0	1.0
Precision	bfloat16	bfloat16
Eval cadence (steps)	300	50
Logging cadence (steps)	50	50
Best model criterion	validation loss	validation loss

Table 7: Fine-tuning hyperparameters for XNLI and PAWS-X, shared across Text Opt/Subopt and IPA Opt/Subopt within each task.

Zipf Deviation decreases with vocabulary size.

Across all settings, Zipf Deviation (ZipfD) is lowest at the largest vocabulary (200k). In Fig. 2, ZipfD appears as higher-is-better because we invert and normalize metrics for readability; in the original definition, lower ZipfD indicates that the empirical rank-frequency curve is closer to an ideal Zipfian trend. The observed monotonic improvement with vocabulary size is intuitive: as vocabularies grow, token inventories increasingly contain longer units (often close to whole-word tokens), and the resulting token frequency distribution more closely mirrors the underlying word-frequency distribution, which is approximately Zipfian in natural language. At the same time, ZipfD is strongly affected by vocabu-

lary granularity, so we interpret the vocabulary-size trend as a configuration effect and emphasize comparisons along the Text vs. IPA axis under matched vocabulary sizes.

Vocabulary utilization always best for 40k, while TTR always best for 200k. Although both metrics quantify vocabulary usage, they capture different aspects of the token distribution. VU measures the fraction of the learned vocabulary that is observed at least once on evaluation data; smaller vocabularies encourage heavier reuse of a limited set of subword units, making it more likely that most entries are encountered. In contrast, TTR measures the number of distinct token types relative to the

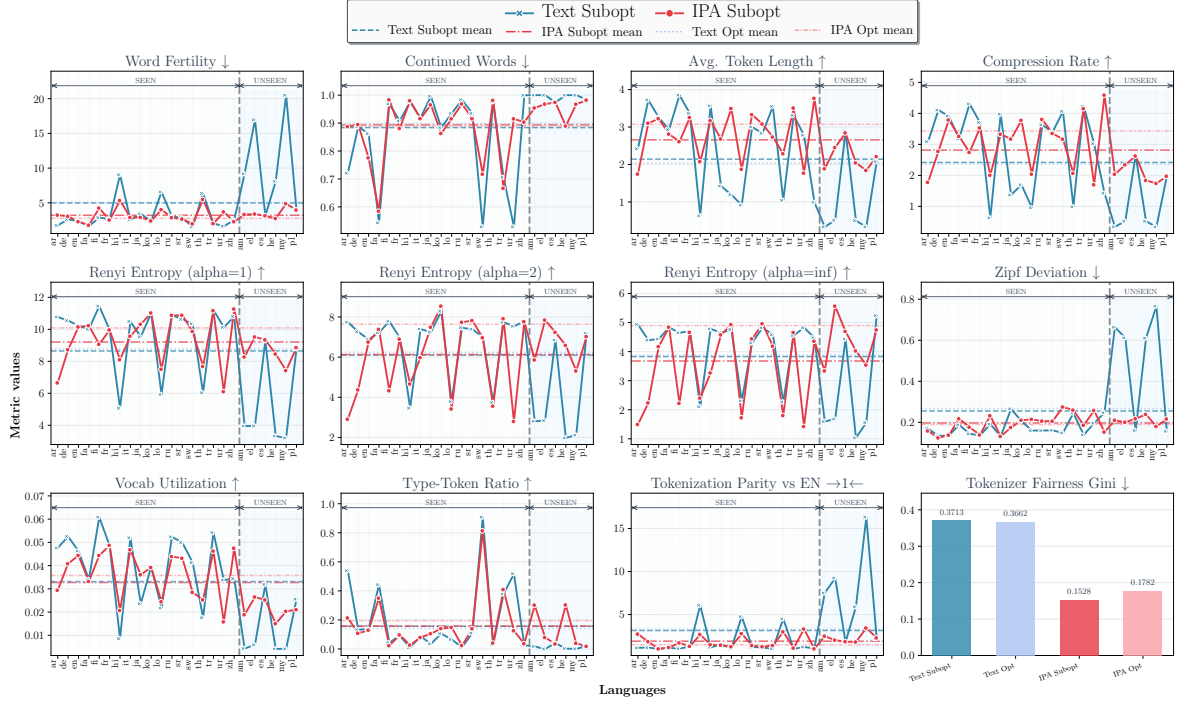


Figure 13: Detailed per-language performance of the suboptimal tokenizers.

total number of produced tokens; larger vocabularies introduce more specific (and often rarer) token types, increasing the diversity of observed types even if many are used only infrequently.

M Additional statistics: win-rate, paired effect size, and mean z -score

We summarize intrinsic differences between IPA and TEXT using (i) a per-language win-rate, (ii) a paired effect size d_z per metric, and (iii) an overall per-language score based on mean z -scored improvements. All computations are paired on the intersection of languages (or scripts) available for both representations. We exclude TFG from per-language/script computations because it is a global metric.

Direction-corrected improvement. For metric m and language/script i , let $x_{m,i}^{\text{IPA}}$ and $x_{m,i}^{\text{TEXT}}$ be the corresponding metric values. We convert each metric to an “IPA-better” improvement value $\Delta_{m,i}$:

$$\Delta_{m,i} = \begin{cases} x_{m,i}^{\text{IPA}} - x_{m,i}^{\text{TEXT}} & (\uparrow) \\ x_{m,i}^{\text{TEXT}} - x_{m,i}^{\text{IPA}} & (\downarrow) \\ |x_{m,i}^{\text{TEXT}} - 1| - |x_{m,i}^{\text{IPA}} - 1| & (\rightarrow 1 \leftarrow) \end{cases} \quad (19)$$

Here (\uparrow) denotes metrics where larger values are better, (\downarrow) denotes metrics where smaller values are better, and $(\rightarrow 1 \leftarrow)$ denotes metrics where values

closer to 1 are better. Thus, $\Delta_{m,i} > 0$ means IPA is better than TEXT on metric m for item i .

Win-rate (per language/script). Let \mathcal{M}_i be the set of metrics available for item i . The win-rate is the fraction of metrics improved by IPA:

$$\text{WinRate}(i) = \frac{1}{|\mathcal{M}_i|} \sum_{m \in \mathcal{M}_i} \mathbb{I}[\Delta_{m,i} > 0]. \quad (20)$$

Paired effect size d_z (per metric). For each metric m , we compute the mean and standard deviation of $\Delta_{m,i}$ across items $i \in \mathcal{I}_m$ (the items with non-missing paired values):

$$\bar{\Delta}_m = \frac{1}{|\mathcal{I}_m|} \sum_{i \in \mathcal{I}_m} \Delta_{m,i}, \quad s_m = \text{std}(\{\Delta_{m,i}\}_{i \in \mathcal{I}_m}). \quad (21)$$

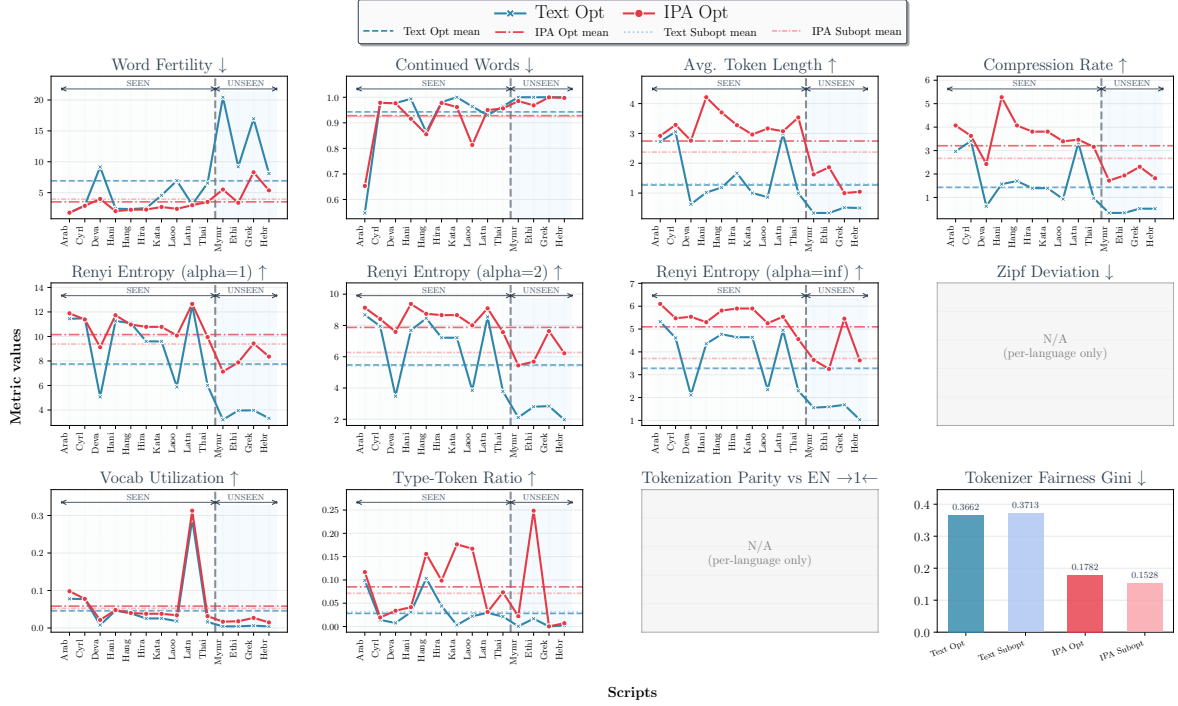
We then report Cohen’s paired standardized mean difference:

$$d_z(m) = \frac{\bar{\Delta}_m}{s_m + \varepsilon}, \quad (22)$$

where ε is a small constant for numerical stability. Positive d_z indicates that IPA tends to outperform TEXT on metric m .

Mean z -score (overall per language/script). To combine metrics with different scales, we z -score improvements within each metric across items:

$$z_{m,i} = \frac{\Delta_{m,i} - \mu_m}{\sigma_m}, \quad (23)$$



Scripts

Figure 14: Per-script intrinsic tokenization metrics: **Text Opt** vs. **IPA Opt**. Solid lines show *Text Opt* and *IPA Opt* per language; dashed lines are means. Arrows: \uparrow = higher is better, \downarrow = lower is better, $\rightarrow 1 \leftarrow$ = closer-to-1 is better; TFG is a single global tokenizer metric (not per-language), hence shown as bars. Zipf Deviation and Tokenization Parity are per-language only, so they are omitted here. SEEN/UNSEEN mark languages included in vs. absent from tokenizer training data. **Takeaway:** *IPA* improves tokenization quality for non-Latin scripts, while maintaining the same quality on Latin script languages as *Text* tokenizers. Improvements are most drastic for scripts unseen during tokenizer training.

where μ_m and σ_m are the mean and standard deviation of $\Delta_{m,i}$ over $i \in \mathcal{I}_m$ (if $\sigma_m = 0$, we set $z_{m,i} = 0$). The overall score for item i is the mean z across its available metrics:

$$\text{MeanZ}(i) = \frac{1}{|\mathcal{M}_i|} \sum_{m \in \mathcal{M}_i} z_{m,i}. \quad (24)$$

Higher $\text{MeanZ}(i)$ indicates that *IPA* improves more metrics, and by a larger margin, relative to *TEXT*.

N CulturaX post-processing

We apply a lightweight, line-level cleaning step to CulturaX before writing the sampled text to disk. For each streamed document, we first replace newline characters with spaces and then apply the following operations in order: we remove URLs matching `http(s)://...` or `www....`; remove any remaining tag-like spans of the form `<...>`; remove emoji characters; filter the text to retain only Unicode letters (L), numbers (N), punctuation (P), and standard spaces (Zs), deleting all other characters; and finally normalize whitespace by collapsing consecutive whitespace to a single space and stripping

leading/trailing spaces. After this pipeline, we discard empty lines and keep non-empty cleaned lines for sampling.

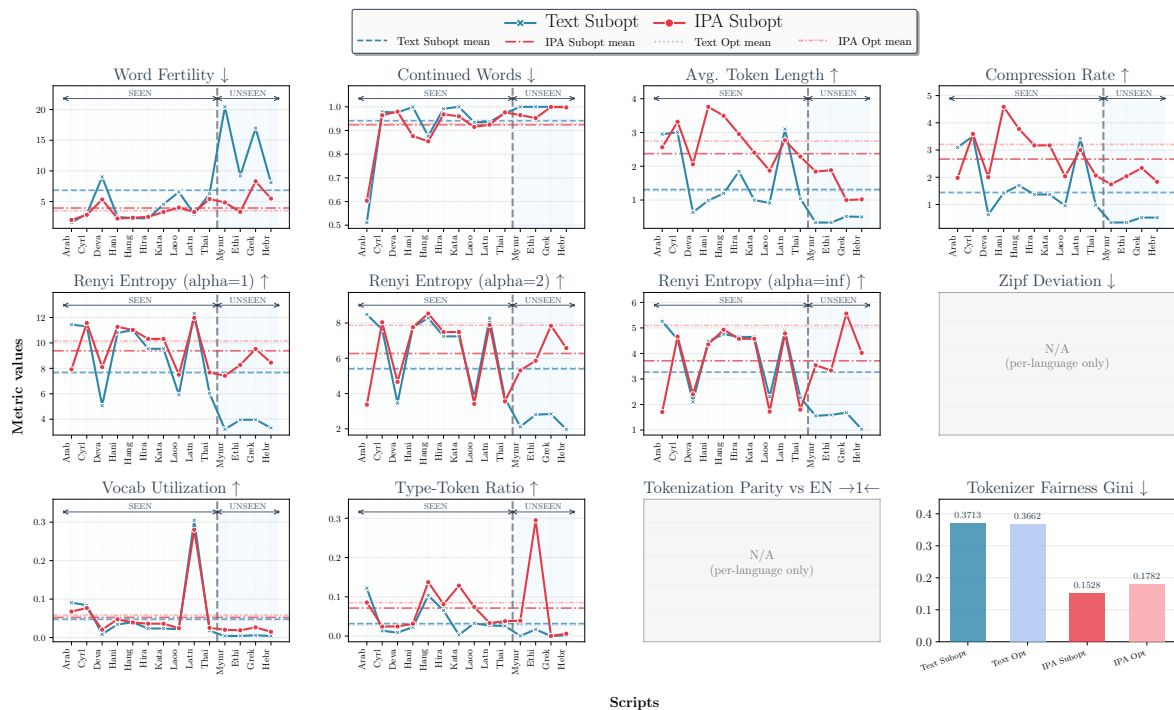


Figure 15: Per-script intrinsic tokenization metrics: **Text Subopt** vs. **IPA Subopt**. Solid lines show *Text Subopt* and *IPA Subopt* per language; dashed lines are means. Arrows: ↑ = higher is better, ↓ = lower is better, → 1 ← = closer-to-1 is better; TFG is a single global tokenizer metric (not per-language), hence shown as bars. Zipf Deviation and Tokenization Parity are per-language only, so they are omitted here. SEEN/UNSEEN mark languages included in vs. absent from tokenizer training data.