

MUTANT: A Recipe for Multilingual Tokenizer Design

Anonymous ACL submission

Abstract

Tokenizers play a crucial role in determining the performance, training efficiency, and the inference cost of Large Language Models (LLMs). Designing effective tokenizers for multilingual LLMs is particularly challenging due to diverse scripts and rich morphological variation. While subword methods like Byte Pair Encoding (BPE) are widely adopted, their effectiveness in multilingual settings remains underexplored. We present MUTANT, a recipe for building multilingual tokenizers, with careful vocabulary and training data design, language-aware pre-tokenization, and subword and multiword aware training. We also introduce MUTANT-Indic, a tokenizer for India-specific multilingual LLMs, that produces linguistically coherent tokens and achieves state-of-the-art performance. Evaluated across English, 22 Indian languages and code data, our tokenizer improves the average fertility score by 39.5% over LLaMA4 and by 18% over Sutra (the current best). This translates to 44% improvement in inference throughput over LLaMA4 while maintaining comparable performance on English and Indic benchmarks. We present detailed ablations across tokenizer training data size, vocabulary size, merging techniques, and pre-tokenization strategies, demonstrating the robustness of our design choices.

1 Introduction

Large Language Models (LLMs) (Touvron et al., 2023; Grattafiori et al., 2024; Abdin et al., 2025; Guo et al., 2025; Yang et al., 2025; Team et al., 2025) rely on the crucial step of tokenization, the process of converting raw text into discrete units called *tokens*. Among the many proposed approaches, subword tokenization schemes such as BPE (Sennrich et al., 2016a), Unigram (Kudo, 2018), WordPiece (Song et al., 2021), and their byte-level extensions have become the de facto choice. However, tokenization remains an understudied topic within the LLM literature (Dagan

et al., 2024; Mielke et al., 2021), especially in multilingual settings (Petrov et al., 2023), where, skewed fertility scores across languages, often lead to concerns around fairness, high inference latency, cost and context size. For instance, with 22 constitutionally recognized languages¹, these issues are especially pronounced for Indic languages comprising multiple scripts and a rich morphology. A key metric for evaluating tokenizers is the “fertility score” (or token-to-word ratio) (Ali et al., 2024) where, a lower fertility score is desirable due to more efficient (and hence cheaper) LLM training and inference. Our analysis suggests that tokenizers of popular multilingual LLMs (Team et al., 2025; OpenAI, 2025; AI, 2025b), largely designed for English, could produce fertility scores as high as 10.5 (LLaMA-4 tokenizer for Oriya; Table 3) for Indic languages, far worse than the near-ideal scores achieved for English. This leads to longer token sequences, higher compute overheads, and poor alignment with linguistic units.

Designing an efficient tokenizer involves making careful choices around the the tokenization algorithm, size of the vocabulary (of tokens) as well as tokenizer training data. In this work, we address five core questions and investigate how to design effective multilingual tokenizers by i) examining trade-offs between low- and high-resource languages, ii) joint versus language-specific training, iii) multilingual data balancing, iv) the role of pre-tokenization, and v) the benefits of incorporating multi-word expressions via curriculum training.

Supported by controlled experiments and systematic ablations, we present MULTilingual Tokenizer optimization ANd Training (MUTANT) as a practical and reproducible recipe (see Figure 1) for building equitable, efficient, and linguistically grounded multilingual tokenizers for LLMs. We

¹https://en.wikipedia.org/wiki/Languages_with_official_recognition_in_India

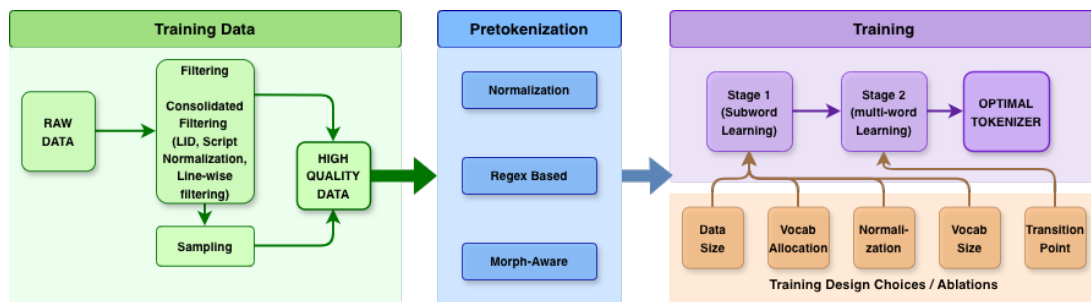


Figure 1: MUTANT: Multilingual tokenizer training recipe. Raw text is filtered and consolidated before a two-stage training process consisting of subword and multi-word learning with ablations to identify efficient design choices.

082 apply our recipe to train MUTANT-Indic, a state-
 083 of-the-art tokenizer for Indic LLMs, demonstrating
 084 the benefits of our methodology in a real-world
 085 multilingual setting. MUTANT-Indic combines lin-
 086 guistically grounded pre-tokenization with a two-
 087 stage subword–multi-word learning process (Liu
 088 et al., 2025a), yielding a compact and semantically
 089 faithful vocabulary. Figure 3 illustrates some ex-
 090 amples where our approach avoids fragmenting
 091 common words or idiomatic phrases into unnatural
 092 subunits across different languages. We make the
 093 following contributions:

- We propose MUTANT, a principled and general recipe for training optimal multilingual tokenizers. We systematically study the effect of vocabulary size, data distribution, language-specific pre-tokenization and multi-stage training on multilingual performance.
- We demonstrate the effectiveness by training MUTANT-Indic, that achieves state-of-the-art (SOTA) performance for Indic languages.
- To the best of our knowledge, we are the first to carry out a comprehensive benchmarking of tokenizers across multiple intrinsic and downstream LLM performance measures in both pretraining from scratch as well as continual pretraining settings. We will publicly release the evaluation framework to enable reproducibility and community benchmarking.

2 Related Work

112 **Multilingual Tokenizers.** Multilingual tokeniza-
 113 tion faces challenges from script diversity, morphol-
 114 ogy, and structural variation. Comparative studies
 115 show that vocabulary size and construction strate-
 116 gies strongly affect performance for morphologi-
 117 cally rich languages (Karthika et al., 2025a), while
 118 inefficiencies in underrepresented ones, such as

119 Ukrainian, translate to higher fertility and com-
 120 putational costs (Maksymenko and Turuta, 2025).
 121 Tokenization also influences how multilingual mod-
 122 els encode morphology, as demonstrated in mT5
 123 vs. ByT5 (Dang et al., 2024). For Indic languages,
 124 tailored resources (Kakwani et al., 2020) and In-
 125 dicBERT (AI4Bharat, 2022) highlight the value
 126 of domain-specific tokenization. Recent bench-
 127 marks further reveal economic implications, with
 128 BLOOM’s tokenizer achieving the best cost effi-
 129 ciency among popular multilingual LLMs (ADA
 130 Sci, 2024). Together, these studies show that cur-
 131 rent multilingual tokenizers fragment low-resource
 132 and morphologically rich languages, motivating
 133 approaches that combine normalization, language-
 134 tailored pre-tokenization, and multi-word learning
 135 to achieve better efficiency and fairness.

136 **Tokenization Algorithms.** Tokenization strate-
 137 gies differ in both theory and practice. While
 138 alternate sub-word tokenization algorithms have
 139 been explored in the past such as WordPiece (Song
 140 et al., 2021), Unigram LM (Kudo and Richard-
 141 son, 2018), Byte Pair Encoding (BPE) remains the
 142 most widely adopted. Originally developed for
 143 compression (Gage, 1994) and later adapted for
 144 neural MT (Sennrich et al., 2016b), BPE merges
 145 frequent character pairs to balance coverage with
 146 efficiency. Its variants aim to address inefficiencies:
 147 PickyBPE (Chizhov et al., 2024) discards uninfor-
 148 mative merges to improve vocabulary utility, while
 149 Scaffold-BPE (Lian et al., 2024) iteratively prunes
 150 low-frequency scaffold tokens to reduce imbalance
 151 and enhance downstream performance. Recent ex-
 152 tensions like SuperBPE (Liu et al., 2025b) expand
 153 beyond word boundaries, jointly learning subwords
 154 and multi-words yielding improved compression
 155 and inference efficiency in a 2-stage curriculum.
 156 BoundlessBPE (Schmidt et al., 2024), another con-
 157 temporary work, relaxes the Pre-tokenization word

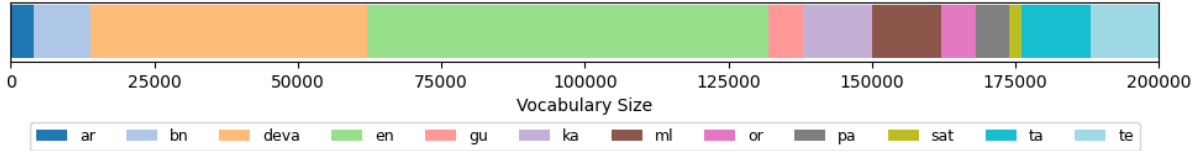


Figure 2: Vocabulary size distribution across language scripts. See Appendix A for script details.

boundary constraint in a single stage learning step. Our work compares these two recent approaches and shows that two-stage curriculum preserves subword coverage while capturing larger semantic units in morphologically rich Indian languages.

Pre-tokenization. Pre-tokenization plays a pivotal role in shaping token boundaries, directly influencing both compression efficiency and reasoning performance (Xue et al., 2024). SentencePiece (Kudo and Richardson, 2018) introduced a language-agnostic approach by treating input as raw streams, effective for languages without whitespace boundaries. More recent approaches like BoundlessBPE (Schmidt et al., 2024) relax pre-token constraints to improve frequency distributions, while regex-based designs continue to prove crucial for capturing script-specific structures.

3 MUTANT: Tokenizer Training Recipe

Language modeling involves estimating the probability distribution over text sequences, $P(S)$, where S may represent a sentence, paragraph, or document. To achieve this, the text is first converted into a sequence of discrete tokens through a tokenization function $g(S) = X = (X_1, X_2, \dots, X_n) \in V^n$, where V denotes the vocabulary and n the sequence length. Tokenizers can be open-vocabulary, ensuring any string can be represented (e.g., byte-level), or closed-vocabulary, where unseen text maps to an out-of-vocabulary symbol (e.g., word lists) (Rae et al., 2021). In our work, we adopt an open-vocabulary approach that combines byte-pair encoding (BPE) with a UTF-8 byte fallback, following Radford et al. (2018). In this section, we present a systematic recipe for designing multilingual tokenizers that jointly addresses data curation, pre-tokenization, and subword–multiword learning. Figure 1 provides an end-to-end overview of the pipeline, illustrating how raw multilingual text is filtered and sampled, transformed through a fixed pre-tokenization stage, and used to train a two-stage tokenizer with explicit design choices on vocabulary size and allocation, detailed below.

3.1 Training Data and Vocabulary

The composition of the training data and the vocabulary size are design choices that directly impact the quality of tokens and the tokenizer efficiency.

Training data volume: Effective tokenizer training relies on a sufficiently large yet carefully curated multilingual corpus; beyond moderate scale, further data scaling yields diminishing returns in tokenization quality (Section 5).

Data quality and diversity: We combine diverse, high-quality sources that is representative of natural language use (Hayase et al., 2024), such as, the Web, curated multilingual datasets, Wikipedia, math, code and others. The raw data is further subjected to filtering for language/script mixing, normalization, short documents (or lines) and controlled source-wise sampling (see Appendix B.1 for MUTANT-Indic).

Vocabulary size and multilingual allocation: In multilingual settings, vocabulary allocation must be explicitly language (or script)-aware to preserve subword and multi-word granularity across both high- and low-resource languages. We determine script-wise vocabulary budgets through ablations that vary allocation across scripts, identifying distributions that balance fragmentation and coverage (Figure 2). To realize these target allocations in practice, we evaluate two strategies: *i) explicit merging*, where we train script-specific tokenizers and concatenates their vocabularies via rule stacking; and *ii) corpus-driven alignment*, where we train a single tokenizer on a multilingual corpus whose data proportions are aligned with the desired script-wise vocabulary budgets.

Explicit merging introduces distributional interference and inconsistent segmentation across scripts (Table 10). In contrast, corpus-driven alignment allows the learned vocabulary to naturally converge toward the target script-wise proportions, better mirroring corpus composition (Table 9) and achieving consistently lower fertility across scripts (Table 3), outperforming both explicit merging and public baselines.

	as	bn	brx	code	doi	eng	gom	gu	hi	kas	kn	mai	ml	mni	mr	nep	or	pa	san	sat	snd	ta	te	urd
Size (MB)	56	562	13	4	1	148	67	83	422	30	252	60	311	34	152	82	46	144	110	0.47	38	502	486	38
# Lines (K)	65	681	19	118	2	449	135	91	545	21	273	129	337	64	257	126	59	198	139	1	69	658	773	27
Avg W/Line	51	47	41	3	43	56	37	59	59	145	41	39	35	43	33	40	45	57	36	26	70	32	32	185

Table 1: Evaluation corpus statistics across 22 Indic languages, English, and code in Section 3.4 including size (in MB), number of lines (in K) and average words per line. We report standard ISO codes here (Appendix Section A).

3.2 Tokenizer Training

To operationalize our design principles for optimal multilingual tokenizers, we adopt a two-stage curriculum following the framework of Liu et al. (2025b). This procedure forms the training strategy used in MUTANT and its Indic instantiation, MUTANT-Indic.

Stage 1 (Subword Learning): Training begins with standard byte-pair encoding (BPE) applied after whitespace pre-tokenization. This ensures that merges occur only within word boundaries, allowing the tokenizer to learn fine-grained *subword units* such as roots, affixes, and common morphemes. Stage 1 continues until the vocabulary reaches a pre-defined *transition point* t ($< |V|$).

Stage 2 (Multi-word Learning): After reaching t , training resumes without whitespace constraints, allowing BPE to merge across word boundaries. This enables the formation of *multi-words*, frequent multiword expressions or collocations (e.g., “in the morning” in Figure 3) improving compression and reducing token counts for common phrases.

Language	Tokenizer	Tokens
English	MI	I wake up early in the morning and get ready for school. My mother makes tea and puts
	Sutra	I wake up early in the morning and get ready for school. My mother makes tea and puts
	Gemma	I wake up early in the morning and get ready for school. My mother makes tea and puts
Hindi	MI	मैं सुबह जल्दी उठ जाता हूँ और तैयार हो जाता हूँ। मैं चाय बनाती हूँ और नारंगी लगा देती हूँ। नारंगी करने के बाद
	Sutra	मैं सुबह जल्दी उठ जाता हूँ और तैयार हो जाता हूँ। मैं चाय बनाती हूँ और नारंगी लगा देती हूँ। नारंगी करने के बाद
	Gemma	मैं सुबह जल्दी उठ जाता हूँ और तैयार हो जाता हूँ। मैं चाय बनाती हूँ और नारंगी लगा देती हूँ। नारंगी करने के बाद मैं
Bengali	MI	এই প্রায় শিকার জন্য কাঁজ করা হয়, এবং সেখানে সাক্ষরতার হার কম বয়সী শিশুদের মাথা দেখা যায়। জুল
	Sutra	এই প্রায় শিকার জন্য কাঁজ করা হয়, এবং সেখানে সাক্ষরতার হার কম বয়সী শিশুদের মাথা দেখা যায়। জুল
	Gemma	এই প্রায় শিকার জন্য কাঁজ করা হয়, এবং সেখানে সাক্ষরতার হার কম বয়সী শিশুদের মাথা দেখা যায়। জুল
Tamil	MI	இந்த கிராமத்தில் ஒரு பள்ளி உள்ளது என்று கூற மக்கள் சொல்லுகிறார்கள். அந்த பள்ளியில்
	Sutra	இந்த கிராமத்தில் ஒரு பள்ளி உள்ளது என்று கூற மக்கள் சொல்லுகிறார்கள். அந்த பள்ளி
	Gemma	இந்த கிராமத்தில் ஒரு பள்ளி உள்ளது என்று கூற மக்கள் சொல்லுகிறார்கள். அந்த பள்ளியில்

Figure 3: MUTANT-Indic (MI) captures multi-words (e.g. “wake up”, “in the morning”) and avoids fragmenting Indic words (see for e.g. Bengali, Tamil).

This two-stage tokenizer training is particularly effective for scripts with complex variations where, meaningful subwords are first anchored and then composed into frequent multiword units.

3.3 Pre-tokenization

As illustrated in Figure 1, pre-tokenization is a fixed, upstream component of our tokenizer design. It segments raw text before subword learning

to improve token consistency and efficiency. For MUTANT-Indic, we explore regex-based, Unicode normalization, and morphology-aware strategies. Unicode-aware regex separates punctuation, handles numeric groups, and aligns tokens with semantic units, while NFKC normalization standardizes visually identical characters to reduce orthographic sparsity (Table 22). To improve robustness across Indic scripts, we explore GPT-2 pre-tokenization rules against LLaMA-4 regex in Stage 1, where the latter yields substantially lower token fragmentation and improving token-to-word ratios by 38–40% on Indic languages (Table 2).

Script-agnostic segmentation (as in LLaMA-4 regex rules) across Indic as well as non-Indic scripts helps reduce token fragmentation, while relaxing whitespace constraints in Stage 2 enables the learning of frequent multi-word expressions. However, unconstrained merging can destabilize generation by creating sentence-crossing tokens. Enforcing additional sentence-level boundary constraints in Stage 2 regex prevents this while preserving multi-word modeling benefits.

Additionally, morphology-aware segmentation decomposes words into roots and affixes to capture recurring morphemes. Although we explore morphology-aware segmentation, integrating it into the tokenization pipeline without incurring latency overhead is challenging (Appendix C.2).

3.4 Evaluation Framework

We develop a modular evaluation framework supporting HuggingFace², SentencePiece³, and TikToken⁴ tokenizers along with a comprehensive set of intrinsic metrics (Appendix D), including Fertility score, Normalized Sequence Length (NSL), Rényi entropy and efficiency, and Bytes Per Token (BPT). All metrics are computed at the line level and aggregated to the language level. We will release both the evaluation dataset and the framework for reproducible benchmarking and fair comparison of multilingual tokenizers. To assess tokenizer

²<https://github.com/huggingface/tokenizers>

³<https://github.com/google/sentencepiece>

⁴<https://github.com/openai/tiktoken>

Regex	as	bn	brx	code	doi	eng	gom	gu	hi	kas	kn	mai	ml	mni	mr	nep	or	pa	san	sat	snd	ta	te	urd
GPT-2	4.36	4.72	4.67	1.57	2.88	1.32	3.95	4.12	3.47	2.47	5.95	3.17	7.08	3.30	4.86	4.37	4.44	3.28	5.97	2.71	1.30	6.53	5.61	1.29
LLaMA-4	1.83	1.74	1.99	1.54	1.56	1.33	2.17	1.83	1.36	1.36	2.15	1.56	2.24	2.27	1.61	1.59	1.65	1.47	2.51	3.60	1.45	2.07	1.83	1.47

Table 2: Fertility scores (\downarrow) showing LLaMA-4 regex outperforms GPT-2 in stage-1 tokenizer training.

Tokenizer (\downarrow)	as	bn	brx	code	doi	eng	gom	gu	hi	kas	kn	mai	ml	mni	mr	nep	or	pa	san	sat	snd	ta	te	urd
Gemma-3	2.65	1.69	2.84	1.79	1.69	1.39	2.60	2.50	1.47	1.48	3.34	1.91	3.45	2.07	2.03	2.03	4.42	2.83	3.37	5.16	2.03	2.50	2.94	1.44
GPT-OSS	2.66	2.41	3.17	1.51	1.89	1.33	2.73	2.37	1.72	1.58	3.34	2.01	3.51	2.41	2.61	2.10	6.26	2.71	3.89	13.01	1.76	3.18	3.13	1.51
LLaMA-4	4.40	2.93	3.34	1.46	2.00	1.34	2.84	3.37	1.83	1.72	4.23	2.28	4.95	2.73	2.79	2.46	10.51	3.23	4.12	9.04	2.13	5.87	4.53	1.76
Sarvam	4.24	1.91	2.92	2.14	1.85	1.66	3.01	2.11	1.53	1.91	2.53	2.11	3.19	4.60	1.94	2.35	2.43	1.67	3.78	13.07	7.62	2.49	2.63	7.93
Sutra	2.12	2.07	3.06	2.12	1.78	1.17	2.68	2.15	1.62	1.48	2.71	2.08	3.10	2.40	2.18	2.01	2.24	1.50	3.76	2.03	2.23	2.58	2.77	1.55
MUTANT-Indic	1.85	1.74	2.04	1.47	1.45	1.12	2.17	1.77	1.23	1.21	2.19	1.58	2.30	2.28	1.63	1.62	1.65	1.39	2.59	3.72	1.45	2.12	1.88	1.44

Table 3: Fertility score (\downarrow) comparison for Indic focused and good Indic support tokenizers across languages. MUTANT-Indic performs best in 20 of 24 languages. An extended version in Table 21 (Appendix).

behavior in Indic use cases, we construct an evaluation set spanning 22 Indic languages, English, and code. Table 1 reports dataset statistics: text volume, number of lines, and average words per line per language.

4 Experiments and Results

As part of our analysis, we compare MUTANT-Indic designed using our recipe against 9 tokenizers, comprising: *i) Indic-focused tokenizers*: tokenizers designed primarily for Indian languages: Sutra (Tamang and Bora, 2024) and Sarvam-2B (Team, 2024b) (referred as Sarvam). *ii) Good Indic support tokenizers*: multilingual tokenizers with demonstrated capabilities for Indic languages: Gemma-3-27B-it (Team et al., 2025) (referred as Gemma-3), GPT-OSS (OpenAI, 2025) and LLaMA-4 (AI, 2025b). *iii) General tokenizers*: tokenizers of widely-used general-purpose LLMs: Qwen3-32B (Team, 2024a) (referred as Qwen-3), LLaMA-3.2-1B (Dubey et al., 2024), Mistral-Nemo (AI, 2024) and DeepSeek-R1 (AI, 2025a).

4.1 Intrinsic Evaluation of Tokenizers

We evaluate tokenization quality using four complementary intrinsic metrics that capture efficiency, granularity, and information utilization: (i) **Fertility score** (Rust et al., 2021; Scao et al., 2022), measuring subword fragmentation; (ii) **Normalized Sequence Length (NSL)** (Dagan et al., 2024), quantifying relative compression against a base tokenizer; (iii) **Rényi entropy and efficiency** (Zouhar et al., 2023), assessing information density and vocabulary utilization; and (iv) **Bytes per token** (Kocetkov et al., 2022), reflecting memory and storage efficiency. All metrics are reported as micro-averages per line at the language level. Formal definitions are provided in Appendix D.

As shown in Table 3, MUTANT-Indic achieves state-of-the-art fertility scores across all evaluated languages, consistently yielding the lowest token-to-word ratios among nine tokenizers with strong Indic support (see Appendix Table 21 for extended comparisons). This indicates substantially reduced fragmentation, particularly for morphologically rich scripts. Bytes-per-token results (Table 4) show that MUTANT-Indic achieves higher values across languages, reflecting more information-dense tokens and improved sequence compactness. Normalized Sequence Length scores (Table 5) further confirms that MUTANT-Indic produces shorter tokenized sequences relative to baseline tokenizers, indicating superior compression efficiency. Finally, Rényi entropy and efficiency results (Table 7) demonstrate that MUTANT-Indic achieves consistently higher efficiency across languages, reflecting effective and balanced utilization of the vocabulary. Collectively, these intrinsic results establish MUTANT-Indic as a robust and efficient tokenizer across diverse scripts and languages.

4.2 Extrinsic Evaluation

We also evaluate the downstream model performance (see Table 6) by pretraining LLaMA-3.2 1B models using two tokenizers: i) MUTANT-Indic, our proposed tokenizer optimized for morphologically meaningful segmentation in Indic and multilingual settings, and (ii) LLaMA-4 tokenizer, chosen for comparable vocabulary size and widespread use. Both models were trained on the same dataset in iso-compute setting to ensure a fair comparison. More details in the Appendix B. We find that our tokenizer shows competitive performance across the English and Indic benchmarks. We additionally trained a model using the Stage-1 tokenizer, which also attains strong downstream performance. Table 23 in Appendix shows that the Stage-1 tokenizer

Tokenizer (\uparrow)	as	bn	brx	code	doi	eng	gom	gu	hi	kas	kn	mai	ml	mni	mr	nep	or	pa	san	sat	snd	ta	te	urd
Gemma-3	6.37	10.45	5.87	2.33	6.75	4.36	5.29	6.31	9.16	7.01	6.73	6.42	7.57	6.23	8.90	8.31	3.76	4.62	6.66	2.59	3.87	9.60	6.82	5.55
GPT-oss	6.36	7.35	5.27	2.77	6.04	4.55	5.02	6.68	7.83	6.54	6.74	6.11	7.43	5.34	6.94	8.04	2.65	4.83	5.79	1.03	4.46	7.56	6.41	5.28
LLaMA-4	3.84	6.05	4.99	2.85	5.70	4.53	4.84	4.69	7.37	6.03	5.33	5.39	5.26	4.71	6.49	6.84	1.58	4.05	5.45	1.48	3.69	4.10	4.43	4.54
Sarvam-2B	3.92	9.42	5.70	1.95	6.16	3.65	4.55	7.62	8.92	5.29	9.07	5.83	8.63	2.81	9.46	7.20	7.17	7.95	6.03	1.02	1.02	9.74	8.46	1.00
Sutra	8.04	8.50	5.44	1.97	6.39	5.15	4.98	7.36	8.33	7.00	8.38	5.88	8.75	5.36	8.35	8.45	7.73	8.76	6.04	6.59	3.49	9.38	8.04	5.15
MUTANT-Indic	9.12	10.15	8.18	2.84	7.86	5.44	6.29	8.95	11.01	8.59	10.30	7.80	11.33	5.67	11.11	10.39	10.07	9.40	8.70	3.60	5.40	11.32	10.70	5.55

Table 4: Bytes-per-token score (\uparrow) comparison for Indic focused and good support tokenizers across languages here. MUTANT-Indic performs best in 22 of 24 languages.

Tokenizer (\downarrow)	as	bn	brx	code	doi	eng	gom	gu	hi	kas	kn	mai	ml	mni	mr	nep	or	pa	san	sat	snd	ta	te	urd
Gemma-3	0.63	0.59	0.87	1.31	0.91	1.06	0.94	0.76	0.83	0.93	0.81	0.89	0.73	0.81	0.76	0.83	0.44	0.89	0.84	0.59	0.99	0.45	0.67	0.85
GPT-oss	0.63	0.83	0.95	1.03	0.96	1.00	0.96	0.71	0.94	0.95	0.79	0.90	0.72	0.89	0.94	0.85	0.60	0.85	0.94	1.43	0.83	0.56	0.71	0.88
Sutra	0.55	0.74	0.93	2.09	0.92	0.89	0.96	0.68	0.92	0.91	0.67	0.94	0.65	0.92	0.84	0.82	0.24	0.51	0.91	0.26	1.10	0.47	0.59	0.90
Sarvam	0.99	0.66	0.91	1.50	1.00	1.27	1.13	0.64	0.85	1.19	0.62	0.99	0.65	2.19	0.72	0.96	0.24	0.54	0.93	1.45	3.63	0.45	0.56	4.25
MUTANT-Indic	0.45	0.60	0.65	0.94	0.78	0.85	0.82	0.54	0.68	0.80	0.53	0.76	0.50	0.91	0.61	0.67	0.18	0.45	0.66	0.45	0.72	0.38	0.44	0.86

Table 5: NSL score (\downarrow) comparison for Indic focused and Good support tokenizers across languages here. MUTANT-Indic performs best in 23 of 24 languages. An extended version in Table 20 (Appendix).

Dataset	LLaMA-4	MUTANT-Indic
English Benchmarks		
HellaSwag	0.353	0.357
CommonsenseQA	0.206	0.204
OpenBookQA	0.216	0.218
Winogrande	0.504	0.510
GSM8K	0.016	0.018
ARC Easy	0.623	0.630
ARC Challenge	0.291	0.292
MMLU	0.252	0.249
DROP	0.048	0.036
Average	0.279	0.279
Indic Benchmarks		
Indic COPA	0.544	0.556
Indic Sentiment	0.524	0.551
Indic XNLI	0.347	0.346
Indic Paraphrase	0.534	0.539
MILU (Indic Multi-turn LU)	0.261	0.258
ARC Challenge (Indic)	0.236	0.244
TriviaQA (Indic)	0.268	0.262
Average	0.388	0.394

Table 6: Extrinsic evaluation of MUTANT-Indic vs LLaMA-4 on English and Indic benchmarks.

itself constitutes a strong and competitive baseline.

The pretraining corpus (Table 17 in Appendix) balances coverage and domain diversity. It combines web-scale sources (Nemotron CC) for general context with structured data including MegaMath, StackV2, synthetic generations, and books. Indic-language content constitutes roughly 20% of the corpus, drawn from Indic CC, Wikipedia, and Sangraha Verified, providing sufficient signal to evaluate cross-lingual and morphologically rich representation quality.

	Gemma-3	GPT-OSS	LLaMA-4	Sarvam	Sutra	MUTANT-Indic
Entropy \downarrow	20.70	20.81	21.09	20.71	20.62	20.42
Efficiency \uparrow	0.22	0.19	0.14	0.21	0.23	0.28

Table 7: Rényi’s Entropy and Efficiency across top Indic tokenizers. Higher efficiency indicates better balance between vocabulary capacity and token usage.

Model	TTFT (ms) \downarrow	OTPT (tokens/s) \uparrow
LLaMA-4	19.17 \pm 0.15	117.99
MUTANT-Indic	18.98 \pm 0.36	169.42

Table 8: Inference latency comparison of 1B models trained with LLaMA-4 and MUTANT-Indic tokenizers.

Metric	ar	bn	deva	en	gu	ka	ml	pa	ta	te
Data size (MB)	106	396	2200	3590	124	644	580	307	616	617
Percentage	1.12	4.18	23.25	37.94	1.31	6.81	6.13	3.24	6.51	6.52
Vocab %	2.69	6.32	20.89	32.92	2.38	7.82	6.76	4.68	7.04	8.50

Table 9: Script-specific training data size (Total corpus size 9.4 GB) and resulting vocabulary % distribution.

4.3 Impact on latency and throughput

Next, we evaluate how tokenization impacts end-to-end model efficiency. We train two 1B-parameter models under identical conditions: one with our tokenizer and one with the LLaMA tokenizer of similar vocabulary size. We then evaluate inference efficiency over 200 samples spanning Indic languages and English, with varying input lengths. Latency⁵ was measured using standard metrics, including Time-To-First-Token (TTFT), Output Throughput (OTPT), and Input Sequence Length (ISL), across 200 instances (See Appendix C.4 for details) with 5 warm-up requests and results averaged over 10 runs. Experiments were done on 8 H100 GPUs using Triton Inference Server as backend, with a maximum generation limit of 256 new tokens. Our tokenizer yields clear efficiency gains (Table 8) which stem from improved compression: shorter token sequences encode more information per token, thereby lowering per-request computation without compromising expressivity. Overall, this demonstrates that tokenizer design directly shapes not only pretraining efficiency but also real-world deployment latency, making it a critical factor for

⁵<https://tinyurl.com/4e7nh7c8>

practical model performance.

4.4 Vocabulary Allocation: Explicit vs. Corpus-Driven

We compare explicit vocabulary merging with corpus-driven joint training under script-aware budget constraints. Explicitly merging script-specific tokenizers leads to fragmented segmentation and higher fertility due to cross-script interference (Table 10). In contrast, corpus-driven joint training naturally aligns vocabulary allocation with data distribution (Table 9), achieving consistently lower fertility across scripts and outperforming merged and public baselines (Table 3). These results indicate that joint, corpus-driven training is a more effective and scalable strategy.

Tokenizer	as	bn	hi	mai	mr	san	te
Individual	2.05	2.13	1.21	1.35	1.75	2.49	1.40
Merged	2.32	2.14	1.55	1.57	1.73	2.79	1.95
MUTANT-Indic	1.85	1.74	1.23	1.58	1.63	2.59	1.88

Table 10: Fertility comparison (\downarrow) between individual script tokenizers and the merged tokenizer across selected Indic languages. Lower values are better.

Dataset	w/ Original	w/ MUTANT-Indic
English Benchmarks		
Winogrande	0.60	0.61
GSM8K	0.05	0.05
ARC Challenge	0.40	0.39
MMLU	0.32	0.29
Average	0.34	0.34
Indic Benchmarks		
Indic COPA	0.58	0.56
Indic Sentiment	0.82	0.85
Indic XNLI	0.35	0.34
Indic Paraphrase	0.57	0.53
Average	0.58	0.57

Table 11: Performance comparison on English and Indic benchmarks in Continual Pretraining (CPT) setting of LLaMA-3.2 1B model (best scores in bold).

4.5 Quality Analysis: Undertrained Tokens

We analyze under-trained ‘‘Glitch’’ tokens in our tied-embedding LLaMA-3.2-1B models trained with both the MUTANT-Indic tokenizer and a comparable BPE tokenizer of similar vocabulary size trained on the same corpus. Both tokenizers share the first 90% of the vocabulary. The MUTANT-Indic tokenizer switches to multi-word training for the last 10% whereas the base BPE tokenizer continues standard subword training. Following Land

and Bartolo (2024) to construct a reference for unused embeddings, we introduced a small set of dummy tokens into the vocabulary that have zero occurrences in the training data. Their embeddings were averaged to obtain a mean reference vector. We then retrieve the top- K nearest neighbors (cosine distance), which represent potential ‘‘glitch’’ tokens (Geiping et al., 2024). As shown in Figure 6 (in the Appendix C.3), the MUTANT-Indic tokenizer produces far fewer such glitch tokens than the base BPE tokenizer. These results suggest that incorporating multi-words promotes more efficient utilization of the vocabulary, while purely subword-based tokenizers overfit in the long tail, yielding a higher proportion of under-trained tokens.

4.6 Can we replace Opensource model tokenizer with MUTANT-Indic?

Following ReTok (Gu et al., 2024), we replace the tokenizer of a pre-trained LLaMA-3.2-1B model (denoted LLaMA-3.2-ORIG) (Grattafiori et al., 2024) with MUTANT-Indic (referred as LLaMA-3.2-MUTANT-Indic). Let V_{orig} and $V_{\text{MUTANT-Indic}}$ denote their corresponding vocabularies. For a token $t \in V_{\text{MUTANT-Indic}}$, we initialize its embedding $E_{\text{init}}(t)$ as: if $t \in V_{\text{orig}} \cap V_{\text{MUTANT-Indic}}$, then $E_{\text{init}}(t) = E_{\text{orig}}(t)$, its embedding from the pre-trained model, otherwise, if $t \in V_{\text{MUTANT-Indic}} \setminus V_{\text{orig}}$ and decomposes under the original tokenizer into (t_1, \dots, t_k) , then $E_{\text{init}}(t) = \frac{1}{k} \sum_{i=1}^k E_{\text{orig}}(t_i)$.

We then continually pretrained the LLaMA-3.2-MUTANT-Indic model, keeping just the embedding and LM head layers trainable, on a 40B-token corpus comprising English, Indic, code, and mathematics (see Appendix for details). As seen in Table 11, the LLaMA-3.2-MUTANT-Indic model performs competitively with the original LLaMA-3.2-ORIG. This suggests that, in addition to pretraining-from-scratch settings, an optimized multilingual tokenizer, such as MUTANT-Indic, could also be leveraged in opensource models through CPT (Continual Pretraining (Chen et al., 2024)) leading to significant throughput gains (see Table 8).

5 Ablation studies

Two-Stage vs. One-Stage: Controlling Vocabulary Recently, BoundlessBPE (Schmidt et al., 2024) also explored a one-stage training paradigm in which pre-tokenization is governed by a fixed regular expression, enabling the direct learning of multiword units in a single pass for English lan-

Tokenizer	as	bn	brx	code	doi	eng	gom	gu	hi	kas	kn	mai	ml	mni	mr	nep	or	pa	san	sat	snd	ta	te	urd
Single-stage (200K)	1.86	1.76	2.05	1.75	1.62	1.37	2.20	1.86	1.39	1.39	2.19	1.61	2.29	2.30	1.66	1.67	1.69	1.49	2.68	3.61	1.56	2.12	1.88	1.54
Two-stage (180K/200K)	1.85	1.74	2.04	1.47	1.45	1.12	2.17	1.77	1.23	1.21	2.19	1.58	2.30	2.28	1.63	1.62	1.65	1.39	2.59	3.72	1.45	2.12	1.88	1.44

Table 12: Fertility score (\downarrow) comparison between one-stage and two-stage MUTANT-Indic tokenizers.

Parameter	as	bn	brx	code	doi	eng	gom	gu	hi	kas	kn	mai	ml	mni	mr	nep	or	pa	san	sat	snd	ta	te	urd	Average
Data size																									
1G	3.02	2.32	2.71	1.62	1.64	1.33	1.97	1.62	1.50	1.43	2.16	1.85	2.83	2.62	1.72	2.13	1.68	1.50	2.46	13.02	1.43	1.92	1.82	1.91	2.42
5G	1.71	1.93	2.58	1.63	1.58	1.33	2.18	1.72	1.40	1.36	2.04	1.57	2.43	2.28	1.68	1.48	1.61	1.57	2.48	4.74	1.30	2.02	1.87	1.43	1.91
10G	1.83	1.74	1.99	1.54	1.56	1.33	2.17	1.83	1.36	1.36	2.15	1.56	2.24	2.27	1.61	1.59	1.65	1.47	2.51	3.60	1.45	2.08	1.83	1.47	1.80
25G	1.75	1.84	2.56	1.62	1.57	1.33	2.15	1.78	1.39	1.36	2.04	1.56	2.32	2.23	1.67	1.47	1.63	1.55	2.45	3.92	1.31	2.01	1.86	1.34	1.86
30G	1.76	1.84	2.32	1.62	1.57	1.33	2.13	1.78	1.39	1.36	2.03	1.57	2.31	2.24	1.67	1.47	1.63	1.54	2.45	4.02	1.31	2.00	1.87	1.35	1.86
50G	1.72	1.82	2.25	1.60	1.57	1.34	2.14	1.82	1.39	1.36	2.03	1.58	2.28	2.22	1.69	1.49	1.64	1.52	2.44	4.54	1.31	2.01	1.87	1.34	1.87
Transition point																									
60	1.91	1.80	2.05	1.39	1.38	1.04	2.16	1.77	1.16	1.15	2.17	1.53	2.30	2.29	1.56	1.58	1.68	1.39	2.48	3.89	1.43	2.11	1.86	1.45	1.81
75	1.91	1.79	2.05	1.41	1.38	1.04	2.16	1.77	1.16	1.15	2.16	1.53	2.30	2.28	1.56	1.58	1.68	1.39	2.47	3.91	1.43	2.10	1.86	1.45	1.81
80	1.89	1.78	2.03	1.41	1.38	1.05	2.15	1.77	1.16	1.16	2.14	1.53	2.28	2.26	1.56	1.57	1.67	1.39	2.46	3.83	1.42	2.08	1.83	1.44	1.80
85	1.87	1.77	2.01	1.43	1.39	1.06	2.13	1.76	1.17	1.16	2.13	1.53	2.26	2.25	1.56	1.56	1.66	1.39	2.46	3.78	1.42	2.07	1.82	1.44	1.80
90	1.85	1.74	2.04	1.47	1.45	1.12	2.17	1.77	1.23	1.21	2.19	1.58	2.30	2.28	1.63	1.62	1.65	1.39	2.59	3.72	1.45	2.12	1.88	1.44	1.83
95	1.85	1.75	1.98	1.47	1.42	1.10	2.13	1.74	1.21	1.20	2.12	1.53	2.23	2.24	1.56	1.56	1.66	1.41	2.46	3.68	1.43	2.06	1.81	1.44	1.79
Vocab size																									
162K/180K	1.89	1.78	2.08	1.48	1.47	1.13	2.21	1.80	1.24	1.22	2.22	1.60	2.35	2.27	1.65	1.65	1.68	1.42	2.62	3.84	1.48	2.16	1.91	1.47	1.86
180K/200K	1.85	1.74	2.04	1.47	1.45	1.12	2.17	1.77	1.23	1.21	2.19	1.58	2.30	2.27	1.63	1.62	1.65	1.39	2.59	3.72	1.45	2.12	1.88	1.44	1.82
202K/225K	1.81	1.70	1.99	1.44	1.43	1.10	2.14	1.72	1.20	1.19	2.14	1.55	2.24	2.21	1.59	1.59	1.60	1.36	2.55	3.59	1.41	2.08	1.82	1.41	1.79
225K/250K	1.78	1.67	1.95	1.42	1.42	1.09	2.11	1.69	1.19	1.17	2.10	1.53	2.20	2.17	1.57	1.57	1.57	1.34	2.52	3.45	1.38	2.04	1.77	1.38	1.75

Table 13: Ablation of tokenizer training data size (in GB), transition point (as a % of vocab size 200K), vocab size ($t = 90\%$) and its impact on fertility score (\downarrow).

498 guage. Our approach instead introduces a two-
499 stage procedure. We replicate the one-stage setup
500 of BoundlessBPE using its released regex and com-
501 pare against our strategy. Our method consistently
502 achieves lower fertility across the top 10 Indic lan-
503 guages and English (Table 12). Overall, the com-
504 parison highlights a clear trade-off: while one-stage
505 methods capture surface-level patterns indiscrim-
506 inately, our two-stage design balances efficiency
507 and linguistic integrity by decoupling subword and
508 multiword learning.

509 **Dataset Size** Similar to Reddy et al. (2025), we
510 study the effects of scaling training data in Stage
511 1 of our training. Table 13 shows that the perfor-
512 mance plateaus after 10G of data.

513 **Transition Point** We ablate the transition point
514 t (Section 3.2) at which training shifts from sub-
515 word to cross-word merges. Varying t reveals a
516 clear trade-off: early transitions favor frequent mul-
517 tiword expressions but weaken morphological cov-
518 erage, while late transitions preserve subwords at
519 the cost of longer sequences. Across Indic and non-
520 Indic languages, intermediate values of 85-90% t
521 yields the best balance, improving token efficiency
522 and cross-lingual consistency (Table 13).

523 **Vocabulary Size** Vocabulary size strongly influ-
524 ences tokenization-model efficiency with trade-offs.
525 Smaller vocabularies yield finer subword units that
526 generalize well to unseen words but with increased
527 compute costs. Larger vocabularies shorten se-
528 quences by encoding frequent forms as single to-
529 kens, but waste capacity on rare items, inflate em-

530 beddings and softmax layers (Shazeer et al., 2017),
531 and have bias toward high-resource languages, hurt-
532 ing multilingual balance. With the same transition
533 point at 90%, we found no significant impact on
534 fertility scores beyond 200K (Table 13).

535 **Effect of Normalization** Prior works show that
536 unicode normalization is crucial for multilingual
537 settings (Karthika et al., 2025b), which reduces
538 token fragmentation and inflated vocabulary size.
539 Table 22 (in Appendix) shows that NFKC yielded
540 marginal but consistent gains by unifying character
541 forms. Accordingly, we adopt NFKC to reduce
542 variability and improve tokenizer robustness.

543 6 Conclusion

544 In this work, we revisit tokenization as a central
545 design choice for multilingual LLMs and introduce
546 MUTANT, a principled and general recipe for mul-
547 tilingual tokenizer training. We demonstrate the
548 effectiveness of this framework with MUTANT-
549 Indic, which achieves SOTA performance across
550 Indic languages through systematic analyses of vo-
551 cabulary size, data distribution, language-specific
552 pre-tokenization, and multi-stage learning. Exten-
553 sive evaluation across intrinsic efficiency metrics,
554 downstream tasks, ablations, and inference latency,
555 in both pretraining-from-scratch and continual pre-
556 training settings, highlights consistent gains in effi-
557 ciency and deployment cost. We will additionally
558 publicly release our evaluation framework to sup-
559 port reproducibility and future research.

560 Limitations

561 While MUTANT-Indic achieves strong intrinsic ef-
562 ficiency gains and practical throughput improve-
563 ments, our study has several limitations. First,
564 although we evaluate across 22 Indic languages,
565 English, and code, the analysis is restricted to these
566 language families, and extending the framework to
567 other low-resource or non-Indic languages remains
568 future work. Second, downstream evaluations are
569 limited to models up to 1B parameters due to com-
570 pute constraints; while prior scaling-law results
571 suggest these gains extrapolate to larger models,
572 we do not validate this empirically. Third, although
573 we analyze morphology-aware pre-tokenization,
574 we do not integrate it into the final tokenizer due
575 to its high inference latency, leaving efficient im-
576 plementations to future work. Finally, while our
577 evaluation framework standardizes intrinsic met-
578 rics, it does not capture all linguistic dimensions
579 relevant to downstream tasks.

580 Ethical Considerations

581 **Ethics Statement** This work focuses on the re-
582 sponsible development of multilingual tokenization
583 methods for Indian languages. We did not collect
584 or utilize any sensitive or Personally Identifiable
585 Information (PII). AI tools (like ChatGPT) were
586 used to refine text writing and section phrasing
587 in the paper. All external datasets, libraries, and
588 tools employed in this work are appropriately ac-
589 knowledged through citations. Since the study did
590 not involve personal, medical, or otherwise sen-
591 sitive information, formal IRB approval was not
592 required. Throughout the process, we aimed to min-
593 imize biases that could disadvantage low-resource
594 languages. We provide our exhaustive study to ad-
595 vance the development of inclusive and efficient
596 multilingual language models.

597 **Reproducibility Statement** To promote trans-
598 parency and reproducibility, we will release the
599 artifacts publicly to benchmark performance of In-
600 dian tokenizers, along with detailed documentation.
601 Detailed records of experimental setups, hyperpa-
602 rameters, and evaluation protocols are maintained
603 to allow replication of our results with the imple-
604 mentation details in the Appendix. In addition,
605 we provide ablation studies to facilitate fair bench-
606 marking and enable future research on Indian and
607 multilingual tokenization.

References

- 608
609 Marah Abdin, Sahaj Agarwal, Ahmed Awadallah,
610 Vidhisha Balachandran, Harkirat Behl, Lingjiao
611 Chen, Gustavo de Rosa, Suriya Gunasekar, Mo-
612 jan Javaheripi, Neel Joshi, and 1 others. 2025.
613 Phi-4-reasoning technical report. *arXiv preprint*
614 *arXiv:2504.21318*.
- 615 ADA Sci. 2024. Multilingual tokenization efficiency in
616 large language models: A study on Indian languages.
617 *Lattice*, 5(2).
- 618 DeepSeek AI. 2025a. Deepseek-r1: Incentiviz-
619 ing reasoning capability in llms via reinforcement
620 learning. [https://github.com/deepseek-ai/](https://github.com/deepseek-ai/DeepSeek-R1)
621 *DeepSeek-R1*.
- 622 Meta AI. 2025b. The llama 4 herd: The
623 beginning of a new era of natively multi-
624 modal intelligence. [https://ai.meta.com/blog/](https://ai.meta.com/blog/llama-4-multimodal-intelligence/)
625 *llama-4-multimodal-intelligence/*.
- 626 Mistral AI. 2024. Mistral nemo. [https://mistral.](https://mistral.ai/news/mistral-nemo)
627 *ai/news/mistral-nemo*.
- 628 AI4Bharat. 2022. IndicBERT: A multilingual AL-
629 BERT model for Indian languages. [https://](https://huggingface.co/ai4bharat/indic-bert)
630 *huggingface.co/ai4bharat/indic-bert*.
- 631 Mehdi Ali, Michael Fromm, Klaudia Thellmann,
632 Richard Rutmann, Max Lübbering, Johannes Lev-
633 eling, Katrin Klug, Jan Ebert, Niclas Doll, Jasper
634 Buschhoff, and 1 others. 2024. Tokenizer choice for
635 llm training: Negligible or crucial? In *Findings of the*
636 *Association for Computational Linguistics: NAACL*
637 *2024*, pages 3907–3924.
- 638 Maharaj Brahma, NJ Karthika, Atul Singh, Devaraj
639 Adiga, Smruti Bhate, Ganesh Ramakrishnan, Rohit
640 Saluja, and Maunendra Sankar Desarkar. 2025. Mor-
641 phtok: Morphologically grounded tokenization for
642 indian languages. *arXiv preprint arXiv:2504.10335*.
- 643 Jie Chen, Zhipeng Chen, Jiapeng Wang, Kun Zhou, Yu-
644 tao Zhu, Jinhao Jiang, Yingqian Min, Wayne Xin
645 Zhao, Zhicheng Dou, Jiaxin Mao, and 1 others.
646 2024. Towards effective and efficient continual pre-
647 training of large language models. *arXiv preprint*
648 *arXiv:2407.18743*.
- 649 Pavel Chizhov, Catherine Arnett, Elizaveta Korotkova,
650 and Ivan P. Yamshchikov. 2024. BPE gets picky: Effi-
651 cient vocabulary refinement during tokenizer training.
652 In *Proceedings of the 2024 Conference on Empiri-
653 cal Methods in Natural Language Processing*, pages
654 16587–16604. Association for Computational Lin-
655 guistics.
- 656 Gautier Dagan, Gabriel Synnaeve, and Baptiste Roziere.
657 2024. Getting the most out of your tokenizer for
658 pre-training and domain adaptation. *arXiv preprint*
659 *arXiv:2402.01035*.
- 660 Anh Dang, Limor Raviv, and Lukas Galke. 2024. Tok-
661 enization and morphology in multilingual language
662 models: A comparative analysis of mT5 and ByT5.
663 *arXiv preprint arXiv:2410.11627*.

664	Sumanth Doddapaneni, Rahul Aralikkatte, Gowtham Ramesh, Shreya Goyal, Mitesh M. Khapra, Anoop Kunchukuttan, and Pratyush Kumar. 2023. Towards leaving no indic language behind: Building monolingual corpora, benchmark and models for indic languages . <i>Preprint</i> , arXiv:2212.05409.	
670	Abhimanyu Dubey and 1 others. 2024. The llama 3 herd of models. <i>arXiv preprint arXiv:2407.21783</i> .	
672	Philip Gage. 1994. A new algorithm for data compression. <i>The C Users Journal</i> , 12(2):23–38.	
674	Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, and 5 others. 2024. The language model evaluation harness .	
682	Jonas Geiping, Alex Stein, Manli Shu, Khalid Saifullah, Yuxin Wen, and Tom Goldstein. 2024. Coercing llms to do and reveal (almost) anything . <i>arXiv preprint arXiv:2402.14020</i> .	
686	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> .	
691	Shuhao Gu, Mengdi Zhao, Bowen Zhang, Liangdong Wang, Jijie Li, and Guang Liu. 2024. Retok: Replacing tokenizer to enhance representation efficiency in large language model . <i>Preprint</i> , arXiv:2410.04335.	
695	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. <i>arXiv preprint arXiv:2501.12948</i> .	
701	Jonathan Hayase, Alisa Liu, Yejin Choi, Sewoong Oh, and Noah A Smith. 2024. Data mixture inference attack: Bpe tokenizers reveal training data compositions. <i>Advances in Neural Information Processing Systems</i> , 37:8956–8983.	
706	Shashank Mohan Jain. 2022. Hugging face . In <i>Introduction to transformers for NLP: With the hugging face library and models to solve problems</i> , pages 51–67. Springer.	
710	Divyanshu Kakwani and 1 others. 2020. Monolingual corpora, evaluation benchmarks and pre-trained models for 11 major Indian languages. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 4948–4961.	
715	N J Karthika, Maharaj Brahma, Rohit Saluja, Ganesh Ramakrishnan, and Maunendra Sankar Desarkar. 2025a. Multilingual tokenization through the lens of indian languages: Challenges and insights . <i>arXiv preprint arXiv:2506.17789</i> .	
	N J Karthika, Maharaj Brahma, Rohit Saluja, Ganesh Ramakrishnan, and Maunendra Sankar Desarkar. 2025b. Multilingual tokenization through the lens of indian languages: Challenges and insights . <i>Preprint</i> , arXiv:2506.17789.	720 721 722 723 724
	Mohammed Safi Ur Rahman Khan, Priyam Mehta, Ananth Sankar, Umashankar Kumaravelan, Sumanth Doddapaneni, Suriyaprasaad G, Varun Balan G, Sparsh Jain, Anoop Kunchukuttan, Pratyush Kumar, Raj Dabre, and Mitesh M. Khapra. 2024. Indicllmsuite: A blueprint for creating pre-training and fine-tuning datasets for indian languages . <i>arXiv preprint arXiv: 2403.06350</i> .	725 726 727 728 729 730 731 732
	Denis Kocetkov, Raymond Li, Loubna Ben Allal, Jia Li, Chenghao Mou, Carlos Muñoz Ferrandis, Yacine Jernite, Margaret Mitchell, Sean Hughes, Thomas Wolf, and 1 others. 2022. The stack: 3 tb of permissively licensed source code. <i>Transactions on Machine Learning Research</i> .	733 734 735 736 737 738
	Taku Kudo. 2018. Subword regularization: Improving neural network translation models with multiple subword candidates . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 66–75, Melbourne, Australia. Association for Computational Linguistics.	739 740 741 742 743 744 745
	Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 66–71.	746 747 748 749 750 751
	Anoop Kunchukuttan. 2020. The IndicNLP Library. https://github.com/anoopkunchukuttan/indic_nlp_library/blob/master/docs/indicnlp.pdf .	752 753 754 755
	Sander Land and Max Bartolo. 2024. Fishing for magikarp: Automatically detecting under-trained tokens in large language models . <i>Preprint</i> , arXiv:2405.05417.	756 757 758 759
	Haoran Lian, Yizhe Xiong, Jianwei Niu, Shasha Mo, Zhenpeng Su, Zijia Lin, Hui Chen, Peng Liu, Jungong Han, and Guiguang Ding. 2024. Scaffold-BPE: Enhancing byte pair encoding for large language models with simple and effective scaffold token removal . <i>arXiv preprint arXiv:2404.17808</i> .	760 761 762 763 764 765
	Alisa Liu, Jonathan Hayase, Valentin Hofmann, Sewoong Oh, Noah A. Smith, and Yejin Choi. 2025a. Superbpe: Space travel for language models . <i>Preprint</i> , arXiv:2503.13423.	766 767 768 769
	Alisa Liu, Jonathan Hayase, Valentin Hofmann, Sewoong Oh, Noah A. Smith, and Yejin Choi. 2025b. SuperBPE: Space travel for language models . <i>arXiv preprint arXiv:2503.13423</i> .	770 771 772 773
	Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi,	774 775

776	Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, and 1 others. 2024. Starcoder 2 and the stack v2: The next generation. <i>arXiv preprint arXiv:2402.19173</i> .	830
777		831
778		832
779	Daniil Maksymenko and Oleksii Turuta. 2025. Tokenization efficiency of current foundational large language models for the Ukrainian language. <i>PMC Digital Health</i> , 12380774.	833
780		834
781		835
782		836
783	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> .	837
784		838
785		839
786		840
787		841
788		842
789	Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Taffjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, and 21 others. 2025. <i>2 olmo 2 furious</i> . <i>Preprint</i> , arXiv:2501.00656.	843
790		844
791		845
792		846
793		847
794		848
795		849
796	OpenAI. 2025. <i>gpt-oss-120b & gpt-oss-20b model card</i> . <i>Preprint</i> , arXiv:2508.10925.	850
797		851
798	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.	852
799		853
800		854
801		855
802		856
803	Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, and 1 others. 2018. Improving language understanding by generative pre-training. <i>TODD</i> .	857
804		858
805		859
806	Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, and 1 others. 2021. Scaling language models: Methods, analysis & insights from training gopher. <i>arXiv preprint arXiv:2112.11446</i> .	860
807		861
808		862
809		863
810		864
811		865
812	Varshini Reddy, Craig W Schmidt, Yuval Pinter, and Chris Tanner. 2025. How much is enough? the diminishing returns of tokenization training data. <i>arXiv preprint arXiv:2502.20273</i> .	866
813		867
814		868
815		869
816	Phillip Rust, Jonas Pfeiffer, Ivan Vulić, Sebastian Ruder, and Iryna Gurevych. 2021. How good is your tokenizer? on the monolingual performance of multilingual language models. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 3118–3135.	870
817		871
818		872
819		873
820		874
821		875
822		876
823		877
824	Sarvam AI. 2025. Arc-challenge-indic. https://huggingface.co/datasets/sarvamai/arc-challenge-indic .	878
825		879
826		880
827	Teven Le Scao, Angela Fan, Christopher Akiki, Elie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, and 1 others. 2022. Bloom: A 176b-parameter open-access multilingual language model. <i>arXiv preprint arXiv:2211.05100</i> .	881
828		882
829		883
		884
	Craig W. Schmidt, Varshini Reddy, Chris Tanner, and Yuval Pinter. 2024. Boundless byte pair encoding: Breaking the pre-tokenization barrier. <i>arXiv preprint arXiv:2504.00178</i> .	882
		883
		884
	Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016a. Neural machine translation of rare words with subword units. In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.	885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

885	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,	Appendix	896
886	Binyuan Hui, Bo Zheng, Bowen Yu, Chang	A Language Details	897
887	Gao, Chengen Huang, Chenxu Lv, and 1 others.	We provide more details about the 22 languages	898
888	2025. Qwen3 technical report. <i>arXiv preprint</i>	and corresponding scripts considered in our work	899
889	<i>arXiv:2505.09388</i> .	in Table 14. Additionally, ISO code mapping is	900
890	Vilém Zouhar, Clara Meister, Juan Luis Gastaldi, Li Du,	also provided in Table 15.	901
891	Mrinmaya Sachan, and Ryan Cotterell. 2023. Tok-	B Implementation	902
892	enization and the noiseless channel. In <i>Proceedings</i>	B.1 Tokenizer implementation	903
893	<i>of the 61st Annual Meeting of the Association for</i>	Our training code for the tokenizer is based on	904
894	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	the open implementation of SuperBPE ⁶ using Hug-	905
895	pages 5184–5207.	gingFace library (Jain, 2022). We first curate a	906
		multilingual corpus of ~50G from OLMo (OLMo	907
		et al., 2025), Wikipedia ⁷ , books, PDFs, Common	908
		Crawl, and the Sangraha dataset (Khan et al., 2024)	909
		where we sample our tokenizer training data from	910
		varied sources to get enough representation of the	911
		different languages individually. English data was	912
		mostly sampled from OLMo and code data from	913
		Stackv2 (Lozhkov et al., 2024) We applied filter-	914
		ing rules like LID, Script normalization as well as	915
		line-wise filtering to get high quality data. The final	916
		MUTANT-Indic tokenizer uses a shared vocabulary	917
		of 200K tokens, distributed across language scripts	918
		and is trained on 10GB of multilingual high quality	919
		data.	920
		In our work, we also explored merging tokeniz-	921
		ers based on the default priority based BPE in Sen-	922
		tencePiece ⁸ . While we explored implementing the	923
		multi-word two stage curriculum in the Sentence-	924
		Piece, we found that it was not trivial. On the	925
		other hand, HuggingFace showed issues with the	926
		merging strategy. We thus relied on different im-	927
		plementations for different approaches.	928
		B.2 Training details	929
		We provide more details about our training setup as	930
		discussed in Section 4.2. Each model was trained	931
		for 50B tokens under matched hyperparameters	932
		(learning rate, batch size, training steps), aligning	933
		FLOPs to isolate tokenizer effects. The evalua-	934
		tion was performed using lm-eval-harness (Gao	935
		et al., 2024) across standard English benchmarks	936
		(MMLU, GSM8K, Winogrande, TriviaQA, Hel-	937
		laSwag, ARC, OpenBookQA, CommonsenseQA,	938

⁶<https://github.com/PythonNut/superbpe/tree/main>

⁷<https://en.wikipedia.org/wiki/>

⁸<https://github.com/google/sentencepiece>

Family	Script	Languages
Indo-Aryan	Devanagari	Hindi, Marathi, Maithili, Dogri, Konkani, Sanskrit, Nepali, Kashmiri
	Bengali (bn)	Assamese, Bengali
	Gurmukhi (pa)	Punjabi
	Arabic (ar)	Urdu, Sindhi
Dravidian	Kannada (kn)	Kannada
	Malayalam (ml)	Malayalam
	Tamil (ta)	Tamil
	Telugu (te)	Telugu
Tibeto-Burman	Devanagari	Bodo
	Meitei Mayek	Manipuri (Meitei Mayek script)
Austroasiatic	Oi Chiki (sat)	Santali

Table 14: Linguistic composition of the 22 scheduled Indian languages analyzed in this work, with their corresponding scripts.

Code	Language	Code	Language	Code	Language
as	Assamese	bn	Bengali	brx	Bodo
doi	Dogri	gu	Gujarati	hi	Hindi
kn	Kannada	ks	Kashmiri	gom	Konkani
mai	Maithili	ml	Malayalam	mni	Manipuri
mr	Marathi	ne	Nepali	or	Odia
pa	Punjabi	san	Sanskrit	sat	Santali
snd	Sindhi	ta	Tamil	te	Telugu
ur	Urdu				

Table 15: Mapping of ISO codes to corresponding 22 Indic languages.

DROP) and Indic benchmarks (IndicCOPA, Indic-Sentiment, IndicXParaphrase, IndicXNLI (Dodapaneni et al., 2023), ARC Challenge Indic (Sarvam AI, 2025), and MILU (Verma et al., 2024)). We report EM for GSM8K and TriviaQA, F1 for DROP, and Accuracy for other benchmarks. Shot settings were fixed per task: 25-shot for ARC/ARC Challenge Indic, 10-shot for HellaSwag, 5-shot for MMLU, GSM8K, and TriviaQA, and zero-shot for the remainder. This setup allows a direct assessment of how tokenizer design influences pre-training efficiency, semantic representation, and generalization across English and Indic tasks.

C Additional Discussion

C.1 Mismatch Between Loss and Task Performance

Although tokenizers, incorporating multi-word often show slightly higher loss (Liu et al., 2025b) during training compared to models using traditional atomic tokenizers like SentencePiece/BPE, this does not necessarily translate to worse down-

stream performance. We hypothesize that this is due to two complementary factors. First, the introduction of longer or multi-word tokens such as “to the” or “as well as” increases the number of semantically overlapping candidates, making the model’s prediction space less sharply peaked. This means the model may distribute probability across several plausible completions (e.g., “to”, “to the”, “to be”), thereby lowering the maximum assigned probability to the correct token and inflating the cross-entropy loss. In contrast, other BPE tokenizers often yield only one atomic candidate for such function words, allowing sharper predictions with lower loss. Second, MUTANT-Indic tokenizes text into fewer, more meaningful units, so when computing the average loss per token, each mistake contributes more heavily to the total. As a result, although the model learns more compact and generalizable representations, its token-level loss appears higher. This creates a divergence between model loss and real-world task accuracy, indicating that traditional loss curves may underrepresent the representational efficiency and practical utility of compositional tokenizers like MUTANT-Indic.

C.2 Morphologically grounded token splitting

We investigate the impact of incorporating morphological information into tokenization for Indic languages (Brahma et al., 2025). The approach involves pre-processing text with a morphology analyzer to segment words into morphemes prior to training. This experiment focuses on languages in the Devanagari script.

We compare two variants: Tokenizer A, trained on raw text, and Tokenizer B, trained on morpho-

Tokenizer	Architecture	Parameters	Data Size (B)	Learning Rate	Train Steps	Context Length	Batch Size	Vocab Size
LLaMA-4	LLaMA-3.2	1B	53.24	5×10^{-5}	68000	4096	192	201134
MUTANT-Indic	LLaMA-3.2	1B	53.18	5×10^{-5}	68000	4096	192	200008

Table 16: Pretraining configuration for different tokenizers.

Category	Sources	Percentage (%)	Token Count (B)
Web	Nemotron CC	30	15
Math	MegaMath	15	7.5
Code	StackV2	15	7.5
Synthetic	New Generations	10	5
Books	Archive	10	5
Indic	Indic CC	8	4
Indic	Indic Wiki	4	2
Indic	Sangraha Verified	8	4
Total		100	50

Table 17: Pretraining corpus distribution across domains and token count. Indic content is emphasized to reflect multilingual objectives.

logically segmented text using morphology analyzer (Kunchukuttan, 2020). At inference time, Tokenizer B requires the same pre-processing for consistency. Tokenizer B exhibits more semantically coherent multi-words, reflecting meaningful morpheme combinations (Figure 4, 5). This promotes better generalization across related forms and reduces the raw token-to-word ratio, as morpheme-based units are more compressible. Sample outputs (Figures 4, 5) illustrate the contrast between surface-level splits and linguistically aligned segmentations.

Despite these gains, we do not adopt this approach in our final tokenizer. The primary limitation is latency, as the pipeline requires both language identification and morphological analysis. For completeness, we evaluated a Hindi morphology-aware tokenizer augmented with a morphological analyzer (Kunchukuttan, 2020) combined with language identification (LID)⁹. We performed inference on approximately 4-5 MB of Hindi text and measured throughput over 10 runs (with 5 warm-up runs), comparing against our MUTANT-Indic tokenizer. Our tokenizer achieved 194K tokens/sec, whereas the morphology-aware tokenizer achieved 90K tokens/sec, representing a 53.28% reduction in throughput. This slowdown arises entirely from the additional LID and morphology-analysis stages, underscoring the efficiency advantages of our approach even when compared to linguistically informed baselines. Extending robust analyzers across all Indic languages also introduces engineering overhead and brittle

⁹<https://pypi.org/project/langdetect/>

dependencies. Nevertheless, morphology-aware tokenization remains a promising direction if fast, reliable analyzers become widely available.

C.3 More on Glitch tokens

For each tokenizer, we vary $K \in \{10, 50, 100, 150, \dots, 400\}$ to select the top- K embeddings closest to a reference vector derived from artificially unused tokens in the vocabulary (Land and Bartolo, 2024; Geiping et al., 2024). For the MUTANT-Indic tokenizer, we count the number of multi-word tokens within the top- K . For the BPE variant, we count tokens with IDs $> 180,000$, which corresponds to the upper 20K of the vocabulary. Both tokenizers share the first 180K IDs; the difference lies in how the final 20K IDs are utilized: MUTANT-Indic allocates this space for frequent multi-word tokens, while the BPE tokenizer continues learning subwords. This design choice allows MUTANT-Indic to more effectively utilize the tail of the vocabulary for meaningful units, whereas the BPE tokenizer exhibits overfitting in low-frequency subwords. The trend of these counts across different top- K values is visualized in Figure 6. As K increases, the fraction of multi-word tokens in MUTANT-Indic remains low but stable, while the BPE variant consistently shows a higher fraction of under-trained subwords, indicating overfitting in the residual vocabulary space.

C.4 More on Latency and Throughput Evaluation

To obtain reliable latency and throughput measurements, we constructed a 200-example multilingual inference set intended to approximate realistic LLM workloads. The set contains diverse sentence-completion style prompts representative of common generation patterns. We include 20 inputs per language across English and nine major Indic languages, ensuring balanced coverage of script diversity, lexical variation, and syntactic complexity. Table 18 presents the token-length distribution of these examples under both the LLaMA-4 tokenizer and our MUTANT-Indic, allowing a controlled comparison of inference efficiency across tokenization schemes.

Tokenized Output:

Words: 79 Tokens: 100

शिक्षा विदों द्वारा पाठ्य पुस्तकों का पुनर्लेखन शिक्षा क्षेत्र के व्यापक आधुनिकीकरण की प्रक्रिया का हिस्सा है। इस नवप्रवर्तन शील प्रयास में भाषाविज्ञानियों, समाजशास्त्रियों तथा मनोवैज्ञानिकों का सक्रिय सहयोग अनिवार्य होता है। विद्यार्थियों की बहुभाषिकता और विविध संवेदनशीलता को ध्यान में रखते हुए शिक्षण विधियों का पुनरावलोकन किया जा रहा है ताकि ज्ञानार्जन की संप्रेषणीयता में निरंतरता बनी रहे।

Figure 4: Tokenized output of morph-aware tokenizer

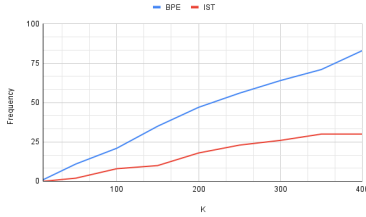


Figure 6: Trend of potential glitch tokens in upper 20K of vocabulary for different K.

Table 18: Token-length statistics for the 200-example inference set. We report min, p75, p90, p99, maximum, and average token lengths.

Tokenizer	min	p75	p90	p99	max	avg
Llama-4	288	805	1157	2583	2869	784
MUTANT-Indic	178	440	541	654	676	379

C.5 Details About Baseline Tokenizers

Tokenizer fertility is shaped by multiple factors including training data distribution, vocabulary construction, and underlying algorithmic choices, yet publicly available documentation on these aspects is often limited. Table 19 summarizes the vocabulary sizes, training methodologies, and any disclosed data distributions for all baseline tokenizers considered in our study.

Tokenizer	Vocab Size	Training Algorithm / Framework	Data Distribution
DeepSeek-R1	128K	BPE (undisclosed variant)	Not publicly disclosed
Gemma-3	262K	SentencePiece	140+ languages
GPT-OSS	200K	o200k_harmony (TikToken variant)	Not publicly disclosed
LLaMA-3.2-1B	128K	BPE / SentencePiece-based	Not publicly disclosed
LLaMA-4	200K	BPE	Not fully disclosed
Mistral-Nemo	131K	Tekken tokenizer (TikToken-based)	100+ languages; multilingual + code
Qwen-3	151K	Byte-level BPE	Not publicly disclosed
Sarvam	68K	Not publicly disclosed	Not publicly disclosed
Sutra	256K	SentencePiece (unigram/BPE hybrid)	Balanced multilingual; uniform sampling

Table 19: Summary of baseline tokenizers and publicly available training details.

D Metrics Definitions

Here, we discuss the different intrinsic metrics used in our evaluation framework.

Tokenized Output:

Words: 53 Tokens: 110

शिक्षाविदों द्वारा पाठ्यपुस्तकों का पुनर्लेखन शिक्षा क्षेत्र के व्यापक आधुनिकीकरण की प्रक्रिया का हिस्सा है। इस नवप्रवर्तनशील प्रयास में भाषाविज्ञानियों, समाजशास्त्रियों तथा मनोवैज्ञानिकों का सक्रिय सहयोग अनिवार्य होता है। विद्यार्थियों की बहुभाषिकता और विविधसंवेदनशीलता को ध्यान में रखते हुए शिक्षणविधियों का पुनरावलोकन किया जा रहा है ताकि ज्ञानार्जन की संप्रेषणीयता में निरंतरता बनी रहे।

Figure 5: Tokenized output of non morph-aware tokenizer

D.1 Token-to-Word Ratio

The Token-to-word ratio measures the average number of tokens required to represent a single word. It captures the degree of segmentation induced by a tokenizer and is particularly informative for morphologically rich languages where excessive fragmentation increases sequence length. We report this metric to evaluate whether tokenizers balance compact representations with sufficient linguistic coverage.

D.2 Bytes-per-token

Bytes-per-token quantifies the average number of raw text bytes contained in a token. Since scripts differ substantially in character set size and encoding, this metric provides a language-agnostic measure of efficiency. Higher values indicate that tokens encode more information per unit, which reduces sequence length. We include this metric to enable direct comparison of tokenizers across writing systems.

D.3 Normalized Sequence Length

Normalized sequence length measures the average length of tokenized sequences relative to a chosen base tokenizer. Instead of reporting absolute sequence lengths, this metric highlights how much longer or shorter sequences become when compared to an established reference. It enables fairer cross-tokenizer comparisons since raw lengths can vary significantly across languages and corpora. A normalized value greater than one indicates that the tokenizer produces longer sequences than the baseline, while a value less than one reflects more compact tokenization. We include this metric to directly assess relative efficiency in sequence compression.

1119 **D.4 Rényi’s Efficiency**

1120 Rényi’s entropy measures the uncertainty of to-
1121 ken distributions induced by a tokenizer, extending
1122 Shannon entropy by allowing different orders to
1123 emphasize frequent or rare tokens. A tokenizer
1124 with a very large vocabulary may contain many
1125 infrequent tokens that are poorly utilized, while a
1126 very small vocabulary forces overuse of common
1127 tokens. Entropy therefore reflects how effectively
1128 the vocabulary is allocated. To complement this,
1129 Rényi’s efficiency normalizes entropy with respect
1130 to vocabulary size, providing a scale-invariant view
1131 of how well the vocabulary capacity is utilized.
1132 Together, these metrics characterize both the dis-
1133 tributional balance of tokens and the comparative
1134 efficiency of different vocabulary scales.

1135 **E Extended Results**

1136 In the main paper, due to space constraints, we lim-
1137 ited the number of tokenizers presented. Here, we
1138 provide an extended list including all of our base-
1139 line tokenizers. Table 20 and 21 provides the NSL
1140 scores and fertility scores of the other tokenizers
1141 considered in our work. Table 22 shows the fertility
1142 scores for different normalization techniques in our
1143 MUTANT-Indic tokenizer. Additionally, we also
1144 compare the performance of two models trained
1145 with only Stage-1 and the 2-stage curriculum tok-
1146 enizer in Table 23.

Table 20: Comparison of NSL scores (Base LLaMA-4) for different tokenizers across all languages.

Tokenizer (↓)	as	bn	brx	code	doi	eng	gom	gu	hi	kas	kn	mai	ml	mni	mr	nep	or	pa	san	sat	snd	ta	te	urd
DeepSeek-R1	0.83	0.97	1.25	1.03	1.29	0.98	1.28	1.48	1.59	1.29	1.41	1.34	1.52	0.99	1.49	1.61	0.67	1.41	1.19	0.69	1.34	0.82	1.34	1.21
Gemma-3	0.63	0.59	0.87	1.31	0.91	1.06	0.94	0.76	0.83	0.93	0.81	0.89	0.73	0.81	0.76	0.83	0.44	0.89	0.84	0.59	0.99	0.45	0.67	0.85
GPT-OSS	0.63	0.83	0.95	1.03	0.96	1.00	0.96	0.71	0.94	0.95	0.79	0.90	0.72	0.89	0.94	0.85	0.60	0.85	0.94	1.43	0.83	0.56	0.71	0.88
LLaMA-3.2-1B	1.90	2.71	1.08	1.02	1.36	0.99	1.22	2.91	1.47	1.36	3.30	1.16	3.25	1.92	1.41	1.44	1.48	2.45	1.19	1.34	1.33	2.11	3.01	1.58
LLaMA-4	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
Mistral-Nemo	1.00	0.95	1.06	1.15	1.07	1.06	1.09	1.09	1.12	1.08	0.91	1.08	0.95	0.95	1.13	1.21	1.57	0.98	1.04	1.34	1.20	0.63	0.64	0.95
Qwen-3	1.68	2.37	1.78	1.11	1.85	1.03	1.72	2.59	2.65	2.16	2.69	1.97	2.57	1.72	2.35	2.47	1.19	2.37	1.92	0.96	1.37	1.63	2.45	1.63
Sutra	0.55	0.74	0.93	2.09	0.92	0.89	0.96	0.68	0.92	0.91	0.67	0.94	0.65	0.92	0.84	0.82	0.24	0.51	0.91	0.26	1.10	0.47	0.59	0.90
Sarvam	0.99	0.66	0.91	1.50	1.00	1.27	1.13	0.64	0.85	1.19	0.62	0.99	0.65	2.19	0.72	0.96	0.24	0.54	0.93	1.45	3.63	0.45	0.56	4.25
MUTANT-Indic-BR	0.45	0.61	0.66	1.28	0.89	1.04	0.84	0.57	0.77	0.91	0.54	0.80	0.50	0.94	0.63	0.69	0.18	0.48	0.70	0.45	0.78	0.38	0.44	0.92
MUTANT-Indic	0.45	0.60	0.65	0.94	0.78	0.85	0.82	0.54	0.68	0.80	0.53	0.76	0.50	0.91	0.61	0.67	0.18	0.45	0.66	0.45	0.72	0.38	0.44	0.86

Table 21: Fertility scores across tokenizers and languages. Lower is better.

Tokenizer (↓)	as	bn	brx	code	doi	eng	gom	gu	hi	kas	kn	mai	ml	mni	mr	nep	or	pa	san	sat	snd	ta	te	urd
DeepSeek-R1	3.54	2.88	4.23	1.53	2.66	1.34	3.68	4.92	3.02	2.49	6.01	3.21	7.95	2.67	4.17	3.97	7.13	4.48	5.07	6.12	2.82	4.92	6.13	2.17
Gemma3	2.65	1.69	2.84	1.79	1.69	1.39	2.60	2.50	1.47	1.48	3.34	1.91	3.45	2.07	2.03	2.03	4.42	2.83	3.37	5.16	2.03	2.50	2.94	1.44
GPT-OSS	2.66	2.41	3.17	1.51	1.89	1.33	2.73	2.37	1.72	1.58	3.34	2.01	3.51	2.41	2.61	2.10	6.26	2.71	3.89	13.01	1.76	3.18	3.13	1.51
Llama-3.2-1B	8.44	8.08	3.64	1.51	2.92	1.35	3.46	9.95	2.74	2.70	14.44	2.79	16.26	5.31	3.90	3.52	15.68	7.88	4.86	12.15	2.85	12.25	13.68	2.73
LLaMA-4	4.40	2.93	3.34	1.46	2.00	1.34	2.84	3.37	1.83	1.72	4.23	2.28	4.95	2.73	2.79	2.46	10.51	3.23	4.12	9.04	2.13	5.87	4.53	1.76
Mistral-Nemo	4.28	2.82	3.52	1.75	2.12	1.41	3.08	3.63	2.05	1.82	3.84	2.48	4.82	2.67	3.10	2.97	16.92	3.04	4.34	12.16	2.51	3.67	3.71	1.65
Qwen3-32B	7.47	7.11	6.10	1.68	4.05	1.41	5.08	8.87	4.86	3.70	11.48	4.53	12.77	4.76	6.56	6.10	12.37	7.60	8.04	8.81	2.95	9.69	11.10	2.90
Sarvam-2B	4.24	1.91	2.92	2.14	1.85	1.66	3.01	2.11	1.53	1.91	2.53	2.11	3.19	4.60	1.94	2.35	2.43	1.67	3.78	13.07	7.62	2.49	2.63	7.93
Sutra	2.12	2.07	3.06	2.12	1.78	1.17	2.68	2.15	1.62	1.48	2.71	2.08	3.10	2.40	2.18	2.01	2.24	1.50	3.76	2.03	2.23	2.58	2.77	1.55
MUTANT-Indic-BR	1.86	1.76	2.05	1.75	1.62	1.37	2.20	1.86	1.39	1.39	2.19	1.61	2.29	2.30	1.66	1.67	1.69	1.49	2.68	3.61	1.56	2.12	1.88	1.54
MUTANT-Indic	1.85	1.74	2.04	1.47	1.45	1.12	2.17	1.77	1.23	1.21	2.19	1.58	2.30	2.28	1.63	1.62	1.65	1.39	2.59	3.72	1.45	2.12	1.88	1.44

Tokenizer	as	bn	brx	code	doi	eng	gom	gu	hi	kas	kn	mai	ml	mni	mr	nep	or	pa	san	sat	snd	ta	te	urd
NFC	1.8520	1.7449	2.0412	1.4658	1.4520	1.1167	2.1741	1.7664	1.2250	1.2042	2.1845	1.5761	2.3025	2.2421	1.6273	1.6241	1.6464	1.3915	2.5859	3.7170	1.4515	2.1226	1.8754	1.4371
NFD	1.8518	1.7454	2.0413	1.4665	1.4521	1.1168	2.1661	1.7667	1.2252	1.2044	2.1905	1.5765	2.3019	2.2487	1.6274	1.6246	1.6465	1.3917	2.5864	3.7170	1.4523	2.1227	1.8757	1.4377
NFKC	1.8512	1.7430	2.0409	1.4647	1.4520	1.1155	2.1738	1.7644	1.2239	1.2041	2.1812	1.5762	2.2991	2.2327	1.6258	1.6234	1.6420	1.3884	2.5855	3.7172	1.4505	2.1200	1.8724	1.4369

Table 22: Fertility scores with NFC, NFD, NFKC normalization for all languages.

Table 23: Comparison of downstream performance between MUTANT-Indic (Stage-1) and MUTANT-Indic (Stage-2).

English Benchmarks			Indic Benchmarks		
Dataset	MUTANT-Indic-Stage-1	MUTANT-Indic-Stage-2	Dataset	MUTANT-Indic-Stage-1	MUTANT-Indic-Stage-1
HellaSwag		0.348	Indic COPA	0.556	0.556
CommonsenseQA	0.193	0.204	Indic Sentiment	0.557	0.551
OpenBookQA	0.214	0.218	Indic XNLI	0.366	0.346
Winogrande	0.515	0.510	Indic Paraphrase	0.562	0.539
GSM8K	0.021	0.018	MILU (Indic Multi-turn LU)	0.265	0.258
ARC Easy	0.625	0.630	ARC Challenge (Indic)	0.247	0.244
ARC Challenge	0.279	0.292	TriviaQA (Indic)	0.268	0.262
MMLU	0.255	0.249			
DROP	0.042	0.036			
Average	0.277	0.279	Average	0.403	0.394