TokSuite: Measuring the Impact of Tokenizer Choice on Language Model Behavior

Anonymous authors
Paper under double-blind review

ABSTRACT

Tokenizers provide the fundamental basis through which text is represented and processed by language models (LMs). Despite the importance of tokenization, its role in LM performance and behavior is poorly understood due to the challenge of measuring the impact of tokenization in isolation. To address this need, we present TokSuite, a collection of models and a benchmark that supports research into tokenization's influence on LMs. Specifically, we train fourteen models that use different tokenizers but are otherwise identical—using the same architecture, dataset, training budget, and initialization. Additionally, we curate and release a new benchmark that specifically measures model performance subject to real-world perturbations that are likely to influence tokenization. Put together, TokSuite allows robustly decoupling the influence of a model's tokenizer, supporting a series of novel findings that elucidate the respective benefits and shortcomings of a wide range of popular tokenizers.

1 Introduction

Language models (LMs) generally do not process "raw" text directly; instead, they operate on a sequence of "tokens" that represent words, subwords, or characters. As a result, tokenization fundamentally influences the representation learned by LMs and, consequently, affects downstream model capabilities (Mielke et al., 2021). For example, the tokenizer used in T5 (Raffel et al., 2020) cannot represent curly brace tokens, making the T5 models poorly suited for processing many coding languages (Wang et al., 2021c). The importance of tokenization naturally motivates not only understanding the impact of different tokenization strategies, but also the design of better tokenizers. However, tokenization is a relatively understudied aspect of language model development compared to, e.g., model architectures, training recipes, and dataset curation. In fact, the design of the tokenizer is often treated as an afterthought, with many open models simply using a preexisting tokenizer off the shelf. For instance, the GPT-2 tokenizer was directly reused for Meta's Open Pretrained Transformers (OPT) (Zhang et al., 2022), and EleutherAI's GPT-NeoX-20B tokenizer was directly used for the MPT-7B-8k model (Team, 2023) and Pythia models (Biderman et al., 2023).

We argue that one factor contributing to the paucity of research into the impact of tokenization is the relative difficulty—using existing artifacts—of decoupling the impact of the tokenizer with other possible variables (model architecture, training data, etc.). For example, it would be fraught to try to compare the Qwen 3 (Yang et al., 2025) and Llama 3 (Dubey et al., 2024)'s tokenizers by studying the respective models because differences in training data, training duration, and architectural details make it difficult to attribute performance differences specifically to tokenization. Understanding the downstream effects of tokenizer design choices is further complicated by the multifaceted nature of tokenization itself, involving various interrelated factors including the underlying segmentation algorithm (e.g., BPE Gage (1994); Sennrich et al. (2016), Unigram Kudo (2018), WordPiece Wu et al. (2016)), granularity level (e.g., byte-level Xue et al. (2022), character-level, word-level), vocabulary size constraints, and the composition of training data used to learn the vocabulary.

What would it take to reliably measure the impact of tokenization on model performance and behavior? We argue that reliable comparison can only be made through models that are completely identical apart from the tokenizer used, because otherwise differences in performance could be attributable to other factors. Given that no open collection of such models exists, we train and release ¹

¹URL redacted for anonymity

14 LMs with identical initialization, architecture, and training data composition, varying only in the tokenizer used.

Our suite of models covers a wide range of tokenizer types, selected among popular pretrained tokenizers as representatives of their main distinctive features, from the most granular byte-level tokenizers like ByT5, to subword-based approaches including BPE, SentencePiece, and WordPiece variants. This collection encompasses both English-only tokenizers trained on monolingual corpora and multilingual tokenizers designed to handle diverse language families and scripts. The tokenizers additionally exhibit varying approaches to out-of-vocabulary (OOV) handling, unicode normalization strategies, whitespace treatment protocols, continuation token markers for subword boundaries, and pretokenization splitting rules. These differences result in diverse vocabulary sizes ranging from compact, efficient lexicons to comprehensive multilingual vocabularies, each with distinct trade-offs between compression efficiency and linguistic coverage. Noting that different vocabularies might share tokens, we develop a novel vocabulary unification framework that creates bijective mappings between tokenizer-specific and unified token spaces. This allows us to use a unified parameter initialization where embeddings for shared tokens are initialized to the same value across models.

To test how tokenization choices affect model behavior, we introduce a novel benchmark with approximately 5,000 samples. Since the effect of different tokenizers can vary across languages (Ali et al., 2024; Dang et al., 2024b; Seo et al., 2025), our benchmark includes five orthographically and morphologically diverse languages: English, Turkish, Italian, Farsi, and Mandarin Chinese. Our benchmark includes 40 "canonical" multiple-choice text completion questions translated into all five languages. Each question has different perturbed versions manually curated by native speakers that reflect real-world changes users might make. For example, we test what happens when visually identical characters have different Unicode values (like replacing Latin "a" with Cyrillic "a"), when users type Turkish text with English keyboards (causing "ş" to become "s"), when Farsi text includes or omits optional accent marks, and when regular text uses special Unicode formatting like enclosed characters. We also add two specialized benchmarks: an elementary school math dataset and a science, technology, engineering, and mathematics (STEM) dataset, respectively with 20 and 44 "canonical" technical questions alongside targeted perturbations. This multi-domain approach allows us to assess tokenizer performance across general, mathematical, and scientific content.

By applying our new benchmark to our suite of models, we both uncover new findings and confirm existing beliefs relating tokenizer characteristics to model behaviors. For example, we find that perturbations tend to be more detrimental in non-English settings, even for tokenizers that were trained on non-English data. Additionally, we found that essentially all off-the-shelf tokenizers are sensitive to Unicode formatting and style perturbations. Further, we found that the two most unconventional tokenizers, ByT5 (Xue et al., 2022) and TokenMonster (Forsythe, 2025), tended to be more robust, suggesting that further investments should be made in the development of novel tokenizers. Put together, our models, dataset, and findings will support future research aiming to more deeply understand how tokenizer choices affect model behavior.

2 Background

Before focusing on how tokenization can affect downstream LM performance, we first explain how tokenizers are created and how design decisions can affect the final tokenizer.

Tokenizers Tokenization is the process of converting a sequence of input symbols into meaningful lexical tokens from some vocabulary \mathcal{V} . Each entry in the vocabulary corresponds to a particular string, and tokenizing an input string can be seen as segmenting it into strings from the vocabulary. When used as the input layer of an LM, the vocabulary is also used to map each token to an integer ID, $V:S\mapsto\{0,1,\ldots,|V|-1\}$. These IDs are then used to look up a vector representation of the token in an LM's embedding table, thus creating a real-valued vector input for each token in an input sequence. While $\mathcal V$ can be manually enumerated for languages with restrictive grammars, i.e., programming languages, the ambiguity and open-endedness of natural language necessitate estimating an optimal set of tokens from data.

Consequently, differences in tokenizers can result in different token sequences for the same string. These differences can affect both learnability and how information is processed in downstream models. For example, a tokenizer that maps the string "dogs" to two tokens "dog" and "s" allows the

model to "reuse" its understanding of the token for "dog", but requires composing with the meaning of the "s" token as pluralization. In contrast, a tokenizer that includes "dogs" as its own token packs both dog and its pluralization into a single token. These differences generally arise in the three main components involved in the training of tokenizers: the data, the learning algorithm, and preprocessing decisions.

Training Data In order to determine the collection of substrings in the vocabulary, tokenizers are generally trained on a text dataset. While the training process for different approaches to tokenization can vary (and are discussed in the following subsection), one straightforward effect of the training data is that if the training dataset does not include a given word or symbol then that word or symbol will not be included in the vocabulary.

Similarly, differences in tokenizer training datasets can result in different choices for tokens included in \mathcal{V} by different tokenizer learning algorithms. For example, if one tokenizer is trained on web data that includes many examples of the typo "teh", it is more likely to represent it as a single token in its vocabulary compared to a tokenizer that is only trained on highly edited text where this typo is rare.

The inclusion of multilingual data in the tokenizer training data can also have a large effect on the final vocabulary, especially when scripts that do not share an alphabet are included. Generally a much larger vocabulary is required—for example the increase from 32,000 to 256,000 when moving from T5 (Raffel et al., 2020) to mT5 (Xue et al., 2021).

Learning Algorithm When training a tokenizer, a learning algorithm produces a vocabulary \mathcal{V} that somehow "fits" the training data, with string inclusion primarily determined by frequency. Most tokenizers function as compressors (Lester et al., 2024), assigning common words to single tokens while splitting rarer ones.

Common algorithms include Byte-Pair Encoding (BPE) (Gage, 1994), which iteratively merges the most frequent symbol bigrams until reaching vocabulary size $|\mathcal{V}|$; WordPiece (Wu et al., 2016), which merges symbols by maximizing training data likelihood; and Unigram (Kudo, 2018), which starts with all possible segmentations and removes symbols causing minimal unigram loss increase. TokenMonster (Forsythe, 2025) uses an unusual approach, building a global vocabulary from all possible tokens and employing an "ungreedy" algorithm that revises tokenization by lookahead. Byte-level models like ByT5 (Xue et al., 2022) use predefined Unicode vocabularies rather than learned ones (Mielke et al., 2021).

Vocabulary size $|\mathcal{V}|$ significantly affects composition—larger vocabularies include more rare words as individual tokens. While most tokenizer training algorithms ensure that every string in the training set can be tokenized, "byte-fallback" forces \mathcal{V} to include the 256 bytes needed to represent any character in Unicode. This allows tokenization of symbols that do not appear in the training dataset, and is primarily important in cases where the training dataset is not large enough.

While most tokenizer training algorithms ensure that every string in the training set can be tokenized, "byte-fallback" forces \mathcal{V} to include the 256 bytes needed to represent any character in Unicode. This allows tokenization of symbols that do not appear in the training dataset, and is primarily important in cases where the training dataset is not large enough.

For a more in-depth discussion of various tokenization approaches, see Mielke et al. (2021).

Preprocessing Tokenization pipelines often use some form of pre-tokenization, which segments the input text into "intuitive" tokens, such as whitespace-separated words, before the learning algorithm is applied. This segmentation can limit which strings can be added to \mathcal{V} as the learning algorithms do not consider bigrams that cross pre-tokenization boundaries. This means that very common bigrams such as "New York" *cannot* be represented as a single token. While some work (Schmidt et al., 2025; Liu et al., 2025, *et alia*) explores methods that allow cross-boundary merges, most commonly used tokenizers do not.

As another example of pre-tokenization, the GPT-2 tokenizer (Radford et al., 2019) splits contractions—e.g., "we'll" \rightarrow "we", "'ll"—meaning that "we'll" cannot be a token in $\mathcal V$. In contrast, BLOOM's (Workshop et al., 2022) pre-tokenization process does not force contractions to a new token, thus allowing for "we'll" $\in \mathcal V$.

Similar differences exist in the handling of numbers. The pre-tokenization used in some models, like GPT-4 (Achiam et al., 2023), breaks contiguous digits into groups of three ("1337" \rightarrow "133", "7") while other models split numbers into their individual digits. There are also models that rely exclusively on the learning algorithm to decide how to segment numbers into digits. Each approach has trade-offs; for example, splitting numbers into thousands might be natural for math but is less natural for dates. Similar considerations exist for how repeated whitespace is handled, especially in domains like code where whitespace can be especially meaningful.

3 THE TOKSUITE MODELS

3.1 TOKENIZER SELECTION AND CHARACTERISTICS

To systematically investigate how different tokenization design choices affect model performance and robustness, we began by selecting a diverse set of 14 preexisting tokenizers, specifically ByT5 (Xue et al., 2022), TokenMonster (Forsythe, 2025), Phi-3 (Abdin et al., 2024), GPT-2 (Radford et al., 2019), Comma (Kandpal et al., 2025), mBERT (Devlin et al., 2019), Llama-3.2 (Dubey et al., 2024), Tekken (AI, 2024), Qwen-3 (Yang et al., 2025), GPT-40 (Hurst et al., 2024), BLOOM (Workshop et al., 2022), Aya (Dang et al., 2024a), Gemma-2 (Team et al., 2024), and XGLM (Lin et al., 2021). Our selection provides comprehensive coverage across vocabulary sizes (ranging from 259 tokens in byte-level tokenizers like ByT5 to over 256,000 tokens in models such as Aya or XGLM), tokenization algorithms (BPE, WordPiece, Unigram, TokenMonster, and byte-level approaches). This diversity enables systematic analysis of how different tokenizers handle out-of-vocabulary words, morphological variations, and adversarial inputs. The selected tokenizers also encompass notable variation in preprocessing strategies that affect robustness, including different approaches to numerical content handling (digit splitting vs. grouping), contraction processing (rule-based vs. learned), Unicode normalization schemes, and multilingual support ranging from monolingual to 100+ languages. Additionally, the tokenizers vary in their out-of-vocabulary handling mechanisms, with some incorporating byte-fallback and others relying on unknown tokens, providing insight into how these design choices propagate to model robustness under various challenges. Detailed technical specifications for each tokenizer are provided in Table 2 and Table 3 in the Appendix.

3.2 Cross-Tokenizer Vocabulary Alignment

To align vocabularies across tokenizers, we first create a unified "super vocabulary". For each tokenizer i we extract its individual vocabulary \mathcal{V}_i , accounting for tokenizer-specific quirks (like WordPiece's "##" prefixes or Unigram's "_" whitespace markers) in this conversion. We also unify the strings that denote the beginning of a sequence— $\langle s \rangle$, $\langle | \text{beginoftext} | \rangle$, etc. Then, we create a super vocabulary, \mathcal{SV} , by taking the union of all vocabularies $\mathcal{SV} = \bigcup_i \mathcal{V}_i$. Note that this unification is based on the UTF-8 byte representation of each element in the vocabularies; see Appendix A for an example of how multiple tokens that are visually the same can appear in \mathcal{SV} .

Finally, for each tokenizer, we create a mapping, $SV:V(X)\mapsto SV(X)$ that translates a tokenizer's original token IDs to the corresponding positions in the unified super vocabulary. This causes a given token string to always map to the same index—regardless of which tokenizer was used—that is, $\forall i, jSV(V_i(S)) = SV(V_j(S))$, assuming $S \in \mathcal{V}_i \land \in \mathcal{V}_j$. The use of the super vocabulary allows us to use the same initialization for the embeddings for shared tokens across models. This shared starting point removes one factor of variation across models, allowing more rigorous attribution of downstream performance to tokenizer characteristics.

3.3 MODEL ARCHITECTURE AND TRAINING CONFIGURATION

We trained fourteen LMs (one for each tokenizer) using the lingua framework (Videau et al., 2024). Our model architecture and training hyper-parameters follow lingua's Llama-1B configuration with approximately one billion non-embedding layers, which follows conventions from the Llama model family (Dubey et al., 2024). All models use a shared initialization based on the super vocabulary. See Appendix B.1 for more information. All models were trained for 100,000 steps with batches of of 256 length-4096 sequences. We use the AdamW optimizer (Loshchilov & Hutter, 2019) with a weight decay of 0.1 and a peak learning rate of 0.001 with cosine annealing and 2000 warm-up steps.

We train all models on a multilingual corpus totaling approximately 100 billion tokens. For English content, we use FineWeb-Edu (Penedo et al., 2024; Lozhkov et al., 2024), which provides high-quality content filtered from Common Crawl data. For the multilingual components, we use the Chinese, Turkish, Italian, and Farsi subsets of the FineWeb-2 HQ Dataset (Messmer et al., 2025), which is a pre-training dataset curated from FineWeb-2 (Penedo et al., 2025) to select high-quality across languages. The final corpus composition consists of 40B English tokens and 60B multilingual tokens equally distributed across the four target languages (15B each).

We acknowledge that training models with different tokenizers under the same token budget means that each model has seen a different collection of text. For example, 100B tokens correspond to approximately 100GB (ByT5), 413GB (Comma), and 698GB (Gemma-2) of text. However, we consider the alternative—training each model on the same text, but for a different number of training steps—to be more problematic, because training duration heavily influences model performance and some models would be relatively under- or over-trained. Additionally, a tokenizer's efficiency in compressing the training data is a relevant factor in tokenizer selection.

As an initial sanity check that our trained models perform reasonably well, we evaluated their performance on standard benchmarks used to evaluate base (i.e., non-post-trained) LMs: HellaSWAG (Zellers et al., 2019), ARC (Clark et al., 2018), PIQA (Bisk et al., 2020), and XNLI (Conneau et al., 2018). Results are shown in Fig. 1. Overall, we find that our models attain reasonable performance for their parameter and training budget. However, we do find notable differences in performance across different models. Since our models are otherwise equivalent, this performance difference can be attributed directly to tokenization, which we discuss further in Section 5.

4 THE TOKSUITE BENCHMARKS

To systematically study tokenization's impact on model performance, we develop a new benchmark that captures the type of input variations models may encounter in real-world deployment. Unlike existing evaluations that focus on clean, canonical text, our benchmark specifically targets naturally occurring perturbations that expose tokenization-dependent issues across our target languages—Chinese (ZH), English (EN), Farsi (FA), Italian (IT), and Turkish (TR)—and domains including general knowledge, basic arithmetic, and STEM. Since the goal of the benchmark is to measure robustness against changing tokenization schemes, we specifically select straightforward canonical questions that establish a strong baseline performance across all models. The selection of canonical questions follows a model-in-the-loop process in which we iteratively test question candidates across our model suite to ensure high baseline accuracy, allowing us to cleanly measure performance degradation when perturbations are applied.

4.1 Multi-lingual Parallel Dataset

We begin by selecting a seed set of 40 *canonical* questions in multiple-choice text completion format in English that almost all of the fourteen models answer correctly, such as "The capital of France is," "The chemical formula for water is," and "The number of continents on Earth is". We aim for canonical questions that our base models get correct so that we can study cases where perturbations flip the answer incorrect. Each canonical question is then translated into FA, IT, TR, and ZH by native speakers. Subsequently, each example undergoes targeted *perturbations* designed to reflect the morphological and orthographic characteristics of each language. Canonical questions in English are provided in Appendix D.1.2, and further examples to each category with detailed case studies on tokenization differences are presented in Appendix E.

Orthographic Perturbations encompass variations in writing systems, diacritics, script-specific features, input medium challenges, and orthographic errors. Writing System Variations include script variations such as traditional vs. simplified Chinese characters, and romanization—writing text in Latin script like Pinyin for Chinese or Finglish for Farsi. Input medium challenges capture typing scenarios where users employ non-native keyboards, leading to systematic character substitutions. This category also includes spacing irregularities with zero-width characters, and homoglyphs—visually similar characters with different Unicode values. Diacritics perturbations include presence of optional diacritics, where text remains valid with or without marks—fata for /a/, kasra for /e/

²https://anonymous.4open.science/r/toksuite-3CEA/

in FA—and common accent errors (è \rightarrow é). Orthographic errors represent spelling mistakes and character-level variations commonly encountered in real-world text, including vowel substitutions, consonant errors, phonetic spelling variants, common misspellings, and punctuation errors. Register & Style captures variations in linguistic register and stylistic conventions across different contexts. This includes web search query formatting with shortened keyword expressions, standard and domain-specific abbreviations, and word reordering that reflect old orthographic conventions. This category encompasses informal digital communication patterns such as colloquial language, emoji or character substitution, and letter repetition for emphasis.

Morphological challenges cover contractions, compound words, inflectional variations, case marking, and derivations that may fragment or alter token boundaries. These challenges are particularly pronounced in agglutinative languages such as Turkish.

Noise perturbations introduce realistic types of textual noise encountered in practice, including typos, character or space deletion, character permutation, and formatting inconsistencies arising from sources such as OCR or other data processing pipelines. These variations test the robustness of the tokenizer under imperfect input conditions that the models must handle.

Grammatical errors represent typical mistakes made by non-expert speakers. Examples in the benchmark include subject-verb agreement errors, article omission or misuse, wrong preposition, incorrect verb tenses and structural errors.

Linguistic variety covers variations in expressing the same semantic content across different linguistic contexts. It includes equivalent expressions with different syntactic structures, codeswitching, similar words, historical spelling variations, and dialects representing regional language varieties with different vocabulary and spelling conventions.

Structural text elements includes Unicode-based formatting (see Fig. 4) and stylistic variations that preserve semantic content while altering visual presentation.

4.2 MATH & STEM DATASETS

Beyond testing simple world knowledge, a subset of our benchmark tests basic arithmetic and STEM, which allows TokSuite to include additional domain-specific perturbations.

LaTeX **Formatting** variations include straightforward such and examples as \$6\$ and \$N_2\$, as well as more complex formatted expressions like $\frac{kg} \cdot \frac{m}^2} {\text{s}^2}$. We also include ASCIIbased structural representations such as molecular diagrams, tree structures, and flowcharts.

Multilingual basic arithmetic is tested by translating the canonical questions into ZH, FA, TR, and IT.

4.3 THE TOKSUITE EVALUATION FRAMEWORK

Robustness We evaluated each model using lm-eval (Gao et al., 2024) library, with byte-length normalized log-likelihood. For fair comparison between models with different baseline capabilities, we report relative accuracy drop for each model against its canonical performance within each category, computed as $\frac{Acc_{can}-Acc_{pert}}{Acc_{can}}$, where lower values indicate greater robustness.

Intrinsic Tokenization Efficiency We evaluate tokenizers' efficiency in compressing text from the five target languages using 10,000 parallel Flores200 (Team et al., 2022) samples with three metrics: 1) *Subword fertility (SF)*: mean tokens per word, where lower values indicate less segmentation; (2) *Parity*: cross-lingual fairness measured as $\frac{T(s_A)}{T(sS_B)} \approx 1$ for parallel sentences (Ali et al., 2024); (3) *Proportion of continued words (PCW)*: fraction of words requiring multiple tokens (Rust et al., 2020). See Appendix C for detailed results.

Table 1: Tokenization robustness under multilingual text perturbations. Values represent relative performance drop $(\frac{Acc_{can} - Acc_{pert}}{Acc_{can}})$; lower values indicate greater robustness. Perturbation types: Input: non-native keyboard/romanization; Diacr.: optional diacritics; Orth. Errors: orthographic errors; Morph.: derivations/inflections/contractions; Noise: homoglyphs/OCR/typos/spacing; LaTeX: LaTeX-style math formatting; STEM: scientific diagrams and notations; Unic.: Unicode styling characters. NEN=non-English. Break-down of each category and detailed case studies are presented in Appendix E. **Green** and red entries indicate notable rosbustness and fragility, respectively.

Model	Input NEN	Diacr. NEN	Orth EN	. Errors NEN	M EN	orph. NEN	I EN	Noise NEN	LaTeX EN	STEM EN	Unic EN	Avg
TokenMonster	0.23	0.33	0.09	0.02	0.23	-0.05	0.11	0.19	0.23	0.11	0.52	0.18
XGLM	0.35	0.49	0.10	0.12	0.25	0.07	0.12	0.22	0.30	0.29	0.12	0.22
BLOOM	0.31	0.35	0.13	0.08	0.18	0.11	0.18	0.19	0.25	0.11	0.57	0.22
Comma	0.29	0.43	0.05	0.07	0.18	0.00	0.11	0.21	0.23	0.29	0.61	0.23
ByT5	0.30	0.44	0.04	0.06	0.27	0.05	0.14	0.18	0.18	0.29	0.53	0.23
mBERT	0.33	0.44	0.11	0.11	0.23	0.06	0.18	0.22	0.15	0.22	0.62	0.24
GPT-4o	0.30	0.52	0.08	0.05	0.21	0.06	0.16	0.20	0.25	0.33	0.55	0.25
GPT-2	0.35	0.46	0.07	0.10	0.25	0.06	0.15	0.21	0.25	0.35	0.53	0.25
Phi-3	0.33	0.46	0.16	0.09	0.27	0.08	0.17	0.21	0.25	0.22	0.55	0.25
Qwen-3	0.36	0.42	0.14	0.11	0.25	0.06	0.16	0.23	0.26	0.29	0.57	0.26
Gemma-2	0.32	0.43	0.14	0.15	0.24	0.03	0.16	0.25	0.22	0.37	0.57	0.26
Llama-3.2	0.33	0.55	0.11	0.10	0.25	0.08	0.15	0.24	0.18	0.29	0.59	0.26
Aya	0.31	0.47	0.15	0.10	0.22	0.03	0.19	0.25	0.23	0.38	0.58	0.26
Tekken	0.34	0.48	0.18	0.03	0.31	0.10	0.14	0.21	0.27	0.44	0.55	0.28
Avg	0.32	0.45	0.11	0.08	0.24	0.05	0.15	0.22	0.23	0.29	0.53	0.24

5 FINDINGS

How does tokenization algorithm design impact robustness across diverse multilingual settings? While orthographic and morphological diversities present universal difficulties across tokenizers, TokenMonster's performance is particularly striking given its architectural constraints. Despite having a 32,000-token vocabulary trained exclusively on English text—roughly one-tenth the size of multilingual competitors like Aya or XGLM—it achieves the best average robustness score across all multilingual perturbations, with the lowest average relative performance drop of 0.18 (see Table 1). This effectiveness stems not from its vocabulary, but from its unique "ungreedy" tokenization algorithm that allows it to revise the token sequence by looking ahead.

ByT5 also demonstrates exceptional multilingual robustness, on average outperforming 9 models (see Table 1) despite using only a 259-token vocabulary. Its byte-level "token-free" design achieves minimal performance degradation across diverse perturbations: 0.04/0.06 drops for English/non-English orthographic errors (see Table 1), 0.00 drop for English grammatical errors (see Table 9), and top average 0.18 drop for multilingual noise (e.g., typos, OCR errors, etc.) (see Table 14). The model shows particular strength in Turkish and Chinese scenarios, including romanized Pinyin handling and even performance improvements (-0.11) with zero-width characters (see Table 7). However, this robustness comes at an efficiency cost, with the highest subword fertility and PCW scores across all languages (see Appendix C), reflecting the robustness-efficiency trade-off. These findings demonstrate that tokenization algorithm design and segmentation consistency can be critical factors for multilingual performance, often more so than massive training data or vocabulary size.

How does multilingual noise amplify tokenization vulnerabilities? Noise-based perturbations create systematic degradation across all tokenizers, but the average performance drop due to noise is markedly more severe for non-English languages (0.22) compared to English (0.15) (see Table 1). This degradation can stem from the core mechanics of subword tokenization: when noise corrupts a familiar word, the tokenizer fragments it into unfamiliar or non-sensical subword units. This effect is particularly damaging in morphologically complex languages. For instance, a simple spacing error in the Turkish phrase "gün sayısı" (day count) causes it to be re-tokenized into chaotic and less meaningful sequences like gün, ##s, ay, ##1s1 by mBERT or gü, ns, ay, 1s1 by Llama-3.2. In contrast, the byte-level tokenizer ByT5 proves more resilient, as character-level errors result in a predictably altered sequence of known bytes rather than catastrophic fragmentation. This suggests that the reliance on a fixed vocabulary in subword models creates an inherent brittleness

that is significantly exacerbated by noise in multilingual contexts. See Section E.3 for a detailed case study of this fragmentation phenomenon.

How do mathematical and STEM content dependencies reveal fundamental limitations in input structure processing for tokenizers? Technical content presents unique tokenization challenges that extend beyond vocabulary coverage. Analysis of mathematical and STEM content reveals critical tokenizer dependencies, with models showing moderate but significant performance degradation (average drops of 0.23 for LaTeX and 0.29 for STEM content, (see Table 1). Even in simplified text completion format with mild technical notation, models exhibit vulnerability to descriptive STEM content. The clearest example of destructive tokenization is XGLM, with the highest LaTeX performance drop (0.30) and notable performance drop for STEM (0.29). This is likely due to XGLM's tokenizer employing an aggressive normalization strategy that creates a stark performance trade-off. It excels at ignoring superficial text styling but fails significantly on technical domains like STEM and LaTeX, where its "lossy" pre-processing destroys the essential structural and spatial information required for comprehension. These domains rely heavily on precise whitespace treatment, symbol placement, and structural conventions—parallel to challenges in coding tasks where spacing and formatting carry semantic meaning. See Appendix E.4 for a detailed case study.

Are there any universal challenges across tokenizers? Formatting presents a universal challenge across tokenizers. Unicode styling and character transformations degrade performance consistently across nearly all models, with an average drop of 0.53—among one of the highest drops observed in our study (see Tables 1, 16, 17). XGLM shows strong robustness to these perturbations thanks to its NFKC normalization during preprocessing. While it mitigates performance degradation from styled characters, it also means that the tokenizer cannot faithfully represent or generate the diverse Unicode formatting present in real-world text.

Similar patterns emerge at scale Evaluation of the original models (including instruction-tuned and larger variants) from which we sourced our tokenizers, which tend to be larger and/or trained for longer than the models in TokSuite, shows modest improvements to robustness (see Table 18) yet consistent patterns persist despite their exposure to orders of magnitude more training data, indicating that tokenization design can influence these robustness characteristics much more heavily than simply training for longer.

6 RELATED WORK

While tokenization is relatively understudied compared to other aspects of LM development, some past work has also studied how tokenization design choices influence model performance and computational cost.

Tokenization Design Factors: Ali et al. (2024) shows that using English-centric tokenizers in a multilingual setting leads to severe downstream degradation and up to 68% additional training cost owing to inefficient token coverage for non-English languages. Rust et al. (2020) found that monolingual tokenizers play an equally important role for pretraining data size in downstream performance. Islam et al. (2022) showed vocabulary-free neural tokenizers yielded substantial improvements for low-resource languages in multilingual natural language inference.

On algorithmic choice, ByT5 notably shows that a byte-level tokenizer can match or outperform subword-level tokenizers on generative tasks. A comparative work compares mT5 (Xue et al., 2021) and ByT5, which share architecture and data but differ in tokenization, and finds that while their overall performance is comparable, the ByT5 model requires more layers to encode morphological information and performs differently across languages (Dang et al., 2024b). Hou et al. (2023) showed that morphological segmentation consistently outperformed BPE across morphologically rich languages, achieving lower perplexity and more efficient training convergence while enabling smaller models to match larger BPE-trained counterparts. Richburg et al. (2020) provided controlled evidence that Unigram language models perform translation more effectively and exhibit superior recall for rare words compared to BPE, particularly in morphologically rich languages like Swahili and Turkish for neural machine translation (NMT). The original SentencePiece work (Kudo & Richardson, 2018) reported processing speeds up to 380 times faster than subword-based NMT in this setting, while achieving comparable or improved performance in machine translation. In

another thread, Huang et al. (2025) argues for decoupling input and output vocabularies and shows a log-linear benefit from scaling the input vocabulary, i.e., larger token sets often reduce loss and improve performance. Schmidt et al. (2024) explores how vocabulary sizes over a specific range perform similarly across a moderate range for English tasks, suggesting diminishing returns from very large vocabularies in that regime. Tao et al. (2024) demonstrated that most current LLMs use insufficient vocabulary sizes, with their analysis suggesting Llama2-70B's optimal vocabulary size should be 216K tokens, 7 times larger than its actual vocabulary size with 32K tokens.

Tokenization Robustness and Vulnerabilities: Like our work, Chai et al. (2024) study LM's sensitivity to typographical errors and ambiguities caused by the internal token structure; while scaling model parameters mitigates this sensitivity it doesn't eliminate it. Wang et al. (2024) developed an adversarial dataset for tokenizer (ADT) framework, successfully degrading the performance of state-of-the-art LM's through vocabulary-based adversarial examples that exploit tokenization vulnerabilities. They created "trap words" where concatenating two vocabulary tokens forms a different existing vocabulary token, causing LLMs to incorrectly tokenize inputs and produce completely wrong responses, with particularly high effectiveness in Chinese due to tokenization complexity. Geh et al. (2025) demonstrated "adversarial tokenization" using non-canonical segmentations that preserve semantic meaning while evading safety alignment. Their approach successfully bypassed existing defense mechanisms, including LlamaGuard and ShieldGemma, revealing fundamental flaws in current LLM safety training pipelines. Several other previous works (Dhole et al., 2021; Wang et al., 2021a;b) have also evaluated LM's vulnerability to noise.

Limitations in the Background Work: Despite recent advances, tokenization research suffers from critical gaps: lack of open-source model collections differing solely in tokenization, limited robustness benchmarks for tokenizer evaluation, and narrow coverage of languages and tokenizer types. To address these limitations, we trained and open-sourced 14 models with different tokenizers using identical architectures, developed a multilingual robustness benchmark, and evaluated models across diverse input variations to isolate tokenization's impact on performance and stability.

7 Future Work & Limitations

TokSuite models are trained exclusively on five languages with higher mixing rates than massively multilingual models (for example, the highest mixing rate across *all* languages in mT5 (Xue et al., 2021)'s training was less than 5%). This setup may underestimate multilingual interference effects present in more realistic settings, where cross-lingual interference could degrade performance. While additional training data may alleviate some vulnerabilities, tokenizers provide a cost-free inductive bias that fundamentally shapes robustness and efficiency. Critically, intrinsic properties like compression rates directly constrain information processing within token budgets, forcing inefficient tokenizers to underconsume or learn subpar representations for certain languages. While coding tasks could present interesting challenges related to non-natural text and whitespace handling, we excluded them from our benchmark due to inconsistent model performance at the scale we considered. Future research should expand to include these domains and broader linguistic coverage, and investigate whether tokenization vulnerabilities persist at larger model scales.

8 Conclusion

Despite tokenization's fundamental role in language model behavior, practitioners commonly adopt off-the-shelf tokenizers without systematic understanding of their impact. To address this, we introduced TokSuite: 14 identical language models differing only in their tokenizer, plus a benchmark curated by native speakers probing natural variations that capture orthographic and morphological challenges across 5 languages and technical domains. Our results show that tokenizer design can matter more than vocabulary size—for example, an English-only tokenizer (TokenMonster) outperformed larger multilingual ones on certain perturbations, while byte-level models proved more robust to multilingual noise and subword fragmentation. Technical content analysis revealed critical vulnerabilities where trivial formatting differences caused catastrophic performance degradation. Our work provides clear evidence that tokenizer choice directly impacts model robustness and capability across diverse contexts and will support future work on understanding the impact of tokenization on LM performance.

REFERENCES

Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Hassan Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Singh Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Caio C'esar Teodoro Mendes, Weizhu Chen, Vishrav Chaudhary, Parul Chopra, Allison Del Giorno, Gustavo de Rosa, Matthew Dixon, Ronen Eldan, Dan Iter, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos Karampatziakis, Dongwoo Kim, Young Jin Kim, Mahoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Chen Liang, Weishung Liu, Eric Lin, Zeqi Lin, Piyush Madan, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Corby Rosset, Sambudha Roy, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Xianmin Song, Olatunji Ruwase, Praneetha Vaddamanu, Xin Wang, Rachel Ward, Guanhua Wang, Philipp Andre Witte, Michael Wyatt, Can Xu, Jiahang Xu, Sonali Yadav, Fan Yang, Ziyi Yang, Donghan Yu, Cheng-Yuan Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yunan Zhang, Xiren Zhou, and Yifan Yang. Phi-3 technical report: A highly capable language model locally on your phone. ArXiv, abs/2404.14219, 2024. URL https://api.semanticscholar.org/CorpusID:269293048.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.

Mistral AI. mistral-common. https://github.com/mistralai/mistral-common, 2024.

Mehdi Ali, Michael Fromm, Klaudia Thellmann, Richard Rutmann, Max Lübbering, Johannes Leveling, Katrin Klug, Jan Ebert, Niclas Doll, Jasper Buschhoff, et al. Tokenizer choice for llm training: Negligible or crucial? In *Findings of the Association for Computational Linguistics:* NAACL 2024, pp. 3907–3924, 2024.

Lucas Bandarkar, Davis Liang, Benjamin Muller, Mikel Artetxe, Satya Narayan Shukla, Donald Husa, Naman Goyal, Abhinandan Krishnan, Luke Zettlemoyer, and Madian Khabsa. The belebele benchmark: a parallel reading comprehension dataset in 122 language variants. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pp. 749–775. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.ACL-LONG.44. URL https://doi.org/10.18653/v1/2024.acl-long.44.

Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023.

Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. PIQA: reasoning about physical commonsense in natural language. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pp. 7432–7439. AAAI Press, 2020. doi: 10.1609/AAAI.V34I05.6239. URL https://doi.org/10.1609/aaai.v34i05.6239.

Yekun Chai, Yewei Fang, Qiwei Peng, and Xuhong Li. Tokenization falling short: On subword robustness in large language models. *arXiv preprint arXiv:2406.11687*, 2024.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the AI2 reasoning challenge. *CoRR*, abs/1803.05457, 2018. URL http://arxiv.org/abs/1803.05457.

541

543

544

546

547

548

549

550

551

552

553 554

555

556

558

559

561

563

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583 584

585

586

588

592

Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel R. Bowman, Holger Schwenk, and Veselin Stoyanov. XNLI: evaluating cross-lingual sentence representations. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii (eds.), *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pp. 2475–2485. Association for Computational Linguistics, 2018. doi: 10.18653/V1/D18-1269. URL https://doi.org/10.18653/v1/d18-1269.

John Dang, Shivalika Singh, Daniel D'souza, Arash Ahmadian, Alejandro Salamanca, Madeline Smith, Aidan Peppin, Sungjin Hong, Manoj Govindassamy, Terrence Zhao, et al. Aya expanse: Combining research breakthroughs for a new multilingual frontier. *arXiv preprint arXiv:2412.04261*, 2024a.

Thao Anh Dang, Limor Raviv, and Lukas Galke. Tokenization and morphology in multilingual language models: A comparative analysis of mt5 and byt5. *arXiv preprint arXiv:2410.11627*, 2024b.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pp. 4171–4186, 2019.

Kaustubh D. Dhole, Varun Gangal, Sebastian Gehrmann, Aadesh Gupta, Zhenhao Li, Saad Mahamood, Abinaya Mahendiran, Simon Mille, Ashish Srivastava, Samson Tan, Tongshuang Wu, Jascha Sohl-Dickstein, Jinho D. Choi, Eduard H. Hovy, Ondrej Dusek, Sebastian Ruder, Sajant Anand, Nagender Aneja, Rabin Banjade, Lisa Barthe, Hanna Behnke, Ian Berlot-Attwell, Connor Boyle, Caroline Brun, Marco Antonio Sobrevilla Cabezudo, Samuel Cahyawijaya, Emile Chapuis, Wanxiang Che, Mukund Choudhary, Christian Clauss, Pierre Colombo, Filip Cornell, Gautier Dagan, Mayukh Das, Tanay Dixit, Thomas Dopierre, Paul-Alexis Dray, Suchitra Dubey, Tatiana Ekeinhor, Marco Di Giovanni, Tanya Goyal, Rishabh Gupta, Louanes Hamla, Sang Han, Fabrice Harel-Canada, Antoine Honore, Ishan Jindal, Przemyslaw K. Joniak, Denis Kleyko, Venelin Kovatchev, Kalpesh Krishna, Ashutosh Kumar, Stefan Langer, Seungjae Ryan Lee, Corey James Levinson, Hualou Liang, Kaizhao Liang, Zhexiong Liu, Andrey Lukyanenko, Vukosi Marivate, Gerard de Melo, Simon Meoni, Maxime Meyer, Afnan Mir, Nafise Sadat Moosavi, Niklas Muennighoff, Timothy Sum Hon Mun, Kenton Murray, Marcin Namysl, Maria Obedkova, Priti Oli, Nivranshu Pasricha, Jan Pfister, Richard Plant, Vinay Prabhu, Vasile Pais, Libo Qin, Shahab Raji, Pawan Kumar Rajpoot, Vikas Raunak, Roy Rinberg, Nicholas Roberts, Juan Diego Rodriguez, Claude Roux, Paulo Henrique Santos Vasconcellos, Ananya B. Sai, Robin M. Schmidt, Thomas Scialom, Tshephisho Sefara, Saqib Shamsi, Xudong Shen, Yiwen Shi, Haoyue Shi, Anna Shvets, Nick Siegel, Damien Sileo, Jamie Simon, Chandan Singh, Roman Sitelew, Priyank Soni, Taylor Sorensen, William Soto, Aman Srivastava, K. V. Aditya Srivasta, Tony Sun, Mukund Varma T., A. Tabassum, Fiona Anting Tan, Ryan Teehan, Mo Tiwari, Marie Tolkiehn, Athena Wang, Zijian Wang, Zijie J. Wang, Gloria Wang, Fuxuan Wei, Bryan Wilie, Genta Indra Winata, Xinyi Wu, Witold Wydmanski, Tianbao Xie, Usama Yaseen, Michael A. Yee, Jing Zhang, and Yue Zhang. Nl-augmenter: A framework for task-sensitive natural language augmentation. CoRR, abs/2112.02721, 2021. URL https://arxiv.org/abs/2112.02721.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pp. arXiv–2407, 2024.

Alasdair Forsythe. tokenmonster: Ungreedy subword tokenizer and vocabulary trainer for python, go & javascript, 2025. URL https://github.com/alasdairforsythe/tokenmonster.

Philip Gage. A new algorithm for data compression. *C Users Journal*, 12(2):23–38, 1994. URL http://www.pennelynn.com/Documents/CUJ/HTML/94HTML/19940045.HTM.

Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The language model evaluation harness, 07 2024. URL https://zenodo.org/records/12608602.

- Renato Lui Geh, Zilei Shao, and Guy Van den Broeck. Adversarial tokenization. *arXiv preprint arXiv*:2503.02174, 2025.
 - Jue Hou, Anisia Katinskaia, Anh-Duc Vu, and Roman Yangarber. Effects of sub-word segmentation on performance of transformer language models. *arXiv preprint arXiv:2305.05480*, 2023.
 - Hongzhi Huang, Defa Zhu, Banggu Wu, Yutao Zeng, Ya Wang, Qiyang Min, and Xun Zhou. Overtokenized transformer: Vocabulary is generally worth scaling. *arXiv preprint arXiv:2501.16975*, 2025.
 - Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
 - Md Mofijul Islam, Gustavo Aguilar, Pragaash Ponnusamy, Clint Solomon Mathialagan, Chengyuan Ma, and Chenlei Guo. A vocabulary-free multilingual neural tokenizer for end-to-end task learning. arXiv preprint arXiv:2204.10815, 2022.
 - Nikhil Kandpal, Brian Lester, Colin Raffel, Sebastian Majstorovic, Stella Biderman, Baber Abbasi, Luca Soldaini, Enrico Shippole, A Feder Cooper, Aviya Skowron, et al. The common pile v0. 1: An 8tb dataset of public domain and openly licensed text. *arXiv preprint arXiv:2506.05209*, 2025.
 - Taku Kudo. Subword regularization: Improving neural network translation models with multiple subword candidates. In Iryna Gurevych and Yusuke Miyao (eds.), *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 66–75, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1007. URL https://aclanthology.org/P18-1007/.
 - Taku Kudo and John Richardson. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv* preprint arXiv:1808.06226, 2018.
 - Brian Lester, Jaehoon Lee, Alexander A Alemi, Jeffrey Pennington, Adam Roberts, Jascha Sohl-Dickstein, and Noah Constant. Training LLMs over neurally compressed text. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL https://openreview.net/forum?id=pRvhMSV48t. Featured Certification.
 - Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, et al. Few-shot learning with multilingual language models. *arXiv preprint arXiv:2112.10668*, 2021.
 - Alisa Liu, Jonathan Hayase, Valentin Hofmann, Sewoong Oh, Noah A. Smith, and Yejin Choi. Superbpe: Space travel for language models, 2025. URL https://arxiv.org/abs/2503.13423.
 - Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=Bkg6RiCqY7.
 - Anton Lozhkov, Loubna Ben Allal, Leandro von Werra, and Thomas Wolf. Fineweb-edu: the finest collection of educational content, 2024. URL https://huggingface.co/datasets/HuggingFaceFW/fineweb-edu.
 - Bettina Messmer, Vinko Sabolčec, and Martin Jaggi. Enhancing multilingual llm pretraining with model-based data selection. *arXiv*, 2025. URL https://arxiv.org/abs/2502.10361.
 - Sabrina J Mielke, Zaid Alyafeai, Elizabeth Salesky, Colin Raffel, Manan Dey, Matthias Gallé, Arun Raja, Chenglei Si, Wilson Y Lee, Benoît Sagot, et al. Between words and characters: A brief history of open-vocabulary modeling and tokenization in nlp. *arXiv preprint arXiv:2112.10508*, 2021.
 - Guilherme Penedo, Hynek Kydlíček, Alessandro Cappelli, Thomas Wolf, and Mario Sasko. DataTrove: large scale data processing. URL https://github.com/huggingface/datatrove.

- Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for the finest text data at scale, 2024. URL https://arxiv.org/abs/2406.17557.
 - Guilherme Penedo, Hynek Kydlíček, Vinko Sabolčec, Bettina Messmer, Negar Foroutan, Amir Hossein Kargaran, Colin Raffel, Martin Jaggi, Leandro Von Werra, and Thomas Wolf. Fineweb2: One pipeline to scale them all–adapting pre-training data processing to every language. *arXiv* preprint *arXiv*:2506.20920, 2025.
 - Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
 - Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
 - Aquia Richburg, Ramy Eskander, Smaranda Muresan, and Marine Carpuat. An evaluation of subword segmentation strategies for neural machine translation of morphologically rich languages. In *Proceedings of the Fourth Widening Natural Language Processing Workshop*, pp. 151–155, 2020.
 - Angelika Romanou, Negar Foroutan, Anna Sotnikova, Sree Harsha Nelaturu, Shivalika Singh, Rishabh Maheshwary, Micol Altomare, Zeming Chen, Mohamed A. Haggag, Snegha A, Alfonso Amayuelas, Azril Hafizi Amirudin, Danylo Boiko, Michael Chang, Jenny Chim, Gal Cohen, Aditya Kumar Dalmia, Abraham Diress, Sharad Duwal, Daniil Dzenhaliou, and et al. INCLUDE: evaluating multilingual language understanding with regional knowledge. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025.*OpenReview.net, 2025. URL https://openreview.net/forum?id=k3gCieTXeY.
 - Phillip Rust, Jonas Pfeiffer, Ivan Vulić, Sebastian Ruder, and Iryna Gurevych. How good is your tokenizer? on the monolingual performance of multilingual language models. *arXiv preprint arXiv:2012.15613*, 2020.
 - Craig W Schmidt, Varshini Reddy, Haoran Zhang, Alec Alameddine, Omri Uzan, Yuval Pinter, and Chris Tanner. Tokenization is more than compression. *arXiv preprint arXiv:2402.18376*, 2024.
 - Craig W. Schmidt, Varshini Reddy, Chris Tanner, and Yuval Pinter. Boundless byte pair encoding: Breaking the pre-tokenization barrier, 2025. URL https://arxiv.org/abs/2504.00178.
 - Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In Katrin Erk and Noah A. Smith (eds.), *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1715–1725, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/P16-1162. URL https://aclanthology.org/P16-1162/.
 - Jean Seo, Jaeyoon Kim, SungJoo Byun, and Hyopil Shin. How does a language-specific tokenizer affect llms? *arXiv preprint arXiv:2502.12560*, 2025.
 - Chaofan Tao, Qian Liu, Longxu Dou, Niklas Muennighoff, Zhongwei Wan, Ping Luo, Min Lin, and Ngai Wong. Scaling laws with vocabulary: Larger models deserve larger vocabularies. *Advances in Neural Information Processing Systems*, 37:114147–114179, 2024.
 - Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
 - MosaicML NLP Team. Introducing mpt-7b: A new standard for open-source, ly usable llms, 2023. URL www.mosaicml.com/blog/mpt-7b. Accessed: 2023-03-28.
 - NLLB Team, Marta R Costa-Jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*, 2022.

- Mathurin Videau, Badr Youbi Idrissi, Daniel Haziza, Luca Wehrstedt, Jade Copet, Olivier Teytaud, and David Lopez-Paz. Meta Lingua: A minimal PyTorch LLM training library, 2024. URL https://github.com/facebookresearch/lingua.
- Boxin Wang, Chejian Xu, Shuohang Wang, Zhe Gan, Yu Cheng, Jianfeng Gao, Ahmed Hassan Awadallah, and Bo Li. Adversarial GLUE: A multi-task benchmark for robustness evaluation of language models. *CoRR*, abs/2111.02840, 2021a. URL https://arxiv.org/abs/2111.02840.
- D Wang, Y Li, J Jiang, Z Ding, G Jiang, J Liang, and D Yang. Tokenization matters! degrading large language models through challenging their tokenization (no. arxiv: 2405.17067). arxiv, 2024.
- Xiao Wang, Qin Liu, Tao Gui, Qi Zhang, Yicheng Zou, Xin Zhou, Jiacheng Ye, Yongxin Zhang, Rui Zheng, Zexiong Pang, Qinzhuo Wu, Zhengyan Li, Chong Zhang, Ruotian Ma, Zichu Fei, Ruijian Cai, Jun Zhao, Xingwu Hu, Zhiheng Yan, Yiding Tan, Yuan Hu, Qiyuan Bian, Zhihua Liu, Shan Qin, Bolin Zhu, Xiaoyu Xing, Jinlan Fu, Yue Zhang, Minlong Peng, Xiaoqing Zheng, Yaqian Zhou, Zhongyu Wei, Xipeng Qiu, and Xuanjing Huang. TextFlint: Unified multilingual robustness evaluation toolkit for natural language processing. In Heng Ji, Jong C. Park, and Rui Xia (eds.), Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: System Demonstrations, pp. 347–355, Online, August 2021b. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-demo.41. URL https://aclanthology.org/2021.acl-demo.41/.
- Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. *arXiv* preprint *arXiv*:2109.00859, 2021c.
- BigScience Workshop, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, et al. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. Google's neural machine translation system: Bridging the gap between human and machine translation. *CoRR*, abs/1609.08144, 2016. URL http://arxiv.org/abs/1609.08144.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. mt5: A massively multilingual pre-trained text-to-text transformer, 2021. URL https://arxiv.org/abs/2010.11934.
- Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. Byt5: Towards a token-free future with pre-trained byte-to-byte models. *Transactions of the Association for Computational Linguistics*, 10:291–306, 2022.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv* preprint *arXiv*:2505.09388, 2025.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In Anna Korhonen, David R. Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pp. 4791–4800. Association for Computational Linguistics, 2019. doi: 10.18653/V1/P19-1472. URL https://doi.org/10.18653/v1/p19-1472.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt

Table 2: Comprehensive Overview of Selected Tokenizers—Part A: Basic Properties

Tokenizer	Method	Vocab. Size	OOV Handling	Language(s)	Pretokenization
ByT5	Bytes	259 (XS)	Bytes	LA.	None (raw bytes)
TokenMonster	Custom	32,000 (S)	Ignores Unknowns	English-Only	None (boundaries are learned)
Phi-3	BPE	32,064 (S)	Byte-fallback	Multilingual	SentencePiece
GPT-2	BPE	50,257 (M)	Byte-fallback	English-Only	GPT-2
Comma	BPE	64,000 (M)	Byte-fallback	Multilingual	GPT-4
mBERT	WordPiece	110,000 (M)	[UNK]	Multilingual	BERT
Llama-3.2	BPE	128,256 (M)	Byte-fallback	Multilingual	GPT-4
Tekken	BPE	130,000 (M)	Byte-fallback	Multilingual	GPT-4*
Qwen-3	BPE	151, 646 (L)	Byte-fallback	Multilingual	GPT-4*
GPT-4o	BPE	200,000 (L)	Byte-fallback	Multilingual	GPT-4o
BLOOM	BPE	250,680 (L)	Byte-fallback	Multilingual	BLOOM
Aya	BPE	255,029 (L)	Byte-fallback	Multilingual	GPT-2
Gemma-2	Unigram	256, 128 (L)	Byte-fallback	Multilingual	SentencePiece
XGLM	Unigram	256,008 (L)	Byte-fallback	Multilingual	SentencePiece

¹ Vocabulary bucket is indicated in ().

Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. Opt: Open pre-trained transformer language models, 2022.

Table 3: Comprehensive Overview of Selected Tokenizers—Part B: Processing Details. See Appendix A for detailed explanations of tokenization processing terminologies and methodologies.

Tokenizer Name	Numbers	Contractions	Unicode Norm.	Whitespace	Zerowidth chars
ByT5	N/A	N/A	None	N/A	3 Bytes
TokenMonster	Learned	Learned	NFD	Learned	Token
Phi-3	Split	Learned	None	Manual	Token
GPT-2	Group	GPT-2	None	Individual	Token
Comma	Group by 3	GPT-4	None	Learned	Token
mBERT	Learned	Composed	None	Normalized	Normalized/Removed
Llama-3.2	Group by 3	GPT-4	None	Learned	Token
Tekken	Split	GPT-4*	None	Learned	Token
Qwen-3	Split	GPT-4	NFC	Learned	Token
GPT-40	Group by 3	Learned	None	Learned	Token
BLOOM	Learned	Learned	None	Learned	Token
Aya	Split	GPT-2	NFC	Learned	Token
Gemma-2	Split	Learned	None	Manual	Token
XGLM	Learned	Learned	NFKC	Normalized	Normalized/Removed

A TOKENIZER PROCESSING GLOSSARY

PRETOKENIZATION

BERT	Pre-tokenization splits are based on whitespace and punctuation.
GPT-2	Pre-tokenization splits are done on whitespace and transitions between letters, numbers, and punctuation.
GPT-4	GPT-4 pre-tokenization follows GPT-2's approach, but it also creates a new token after 3 contiguous digits. Note that Qwen 3 uses the same pretokenization as GPT-4, but does not split numbers into groups of three.
GPT-40	GPT-40 pre-tokenization follows that of GPT-4, but specific contractions—('s, 'd, 'm, 't, 'll, 've, 're)—are not split from the preceding word. Note that Tekken uses the same pre-tokenization methods as GPT-40, but without special case handling

of the specific english contractions.

² OOV = Out-of-vocabulary

³ LA. = Language-agnostic

and periods.

numbers and punctuation.

810

811

812

813

814

BLOOM

SentencePiece

815 NUMBERS PROCESSING 816 817 **Split** Numbers are deterministically broken down into individual digits which are each 818 treated as single tokens. 819 Group Numbers are deterministically split from adjoining text during pre-tokenization. 820 The learning algorithm then determines which numbers become single tokens and 821 which are further tokenized. 822 Group by 3 823 Similar to **Group**, but contiguous digits are split into groups of 3 during pretokenization. Again, the learning algorithm then determines which numbers are 824 single tokens. For example, "username12345" is pre-tokenized into "username", 825 "123", and "45", but "123" is not a token in V yielding a final token stream of 826 "username", "1", "23", "45". 827 828 Learned Numbers are not automatically segmented from surrounding text. Thus, the learn-829 ing algorithm determines token boundaries for letters and numbers jointly. This can result in tokens that include both characters and digits. 830 831 832 CONTRACTIONS PROCESSING 833 GPT-2 A selected number of English contractions ('s, 'd, 'm, 't, 'll, 've, 're) are manually 834 split into their own tokens. The learning algorithm then decides if they should be 835 their own token or if it should be broken down further. This makes it impossible 836 to have a token like "I'll". 837 GPT-4 Uses GPT-4's contraction processing method. The name set of contractions are 838 explicitly handled, but the regex is implemented differently. Note that Tekken 839 uses the GPT-4 regex without special casing english contractions; however, it still 840 results in splitting contractions from the base during pre-tokenization. 841 Learned Contractions are not manually split from the base word; the learning algorithm 842 decides if the contraction should be its own token or a composition. 843 844 Composed The pre-tokenization splits all contractions into multiple tokens (base, apostro-845 phe, and contraction, e.g., he'll \rightarrow "he", "", "ll"), which cannot be merged back 846 together in the learning algorithm. 847 848 UNICODE NORMALIZATION 849 850 None No Unicode normalization is applied; characters are processed exactly as they 851 appear in the input. Note that this can result in \mathcal{V} containing multiple tokens that 852 are visually the same, but differ in their underlying bytes, for example two "é" 853 tokens, but one is represented by a single code point while the other is represented as the composition of "e" and "'" 854 855 **NFD** Normalization Form Decomposed: Unicode characters are decomposed into their 856 constituent parts (base characters + combining marks separately). 857 **NFC** Normalization Form Composed: Unicode characters are composed into their 858 canonical combined form (base characters + combining marks merged when pos-859 sible). 860 861 **NFKC** Normalization Form Compatibility Composed: Similar to NFC but also applies compatibility mappings, converting visually similar characters to their canonical 862 equivalents before composition. Note that this can result in lossy detokenization as characters like "2" are mapped to "2".

Pre-tokenization splits are done based on whitespace and punctuation like commas

Pre-tokenization splits are done on whitespace, and at transitions between letters,

WHITESPACE TREATMENT

Normalized Whitespace like tabs, newlines, and contiguous spaces are normalized to a single

space. This results in lossy detokenization and often stops the downstream model

from understanding domains with meaningful whitespace such as code.

Learned Each piece of contiguous whitespace is segmented into a single token during pre-

tokenization, then the learning algorithm decides how to subdivide them into individual tokens. This results in whitespace being preserved and allows for lossless

detokenization.

Manual The handling of whitespaces during pre-tokenization matches **Learned**, but pre-

defined whitespace tokens of various sizes are used instead of learning them from the data. This results in whitespace being preserved and allows for lossless detok-

enization.

Individual Whitespace is preserved, but each individual whitespace character is represented

as its own token. This yields long token sequences for whitespace heavy inputs. This results in whitespace being preserved and allows for lossless detokenization.

ZERO-WIDTH CHARACTERS

3 Bytes Zero-width characters are maintained in their original 3-byte representation.

Token Zero-width characters are preserved and assigned as new tokens in the vocabulary.

Normalized/Removed Zero-width characters are either normalized to standard equivalents or completely removed.

B MODEL TRAINING

B.1 MODEL INITIALIZATION

We use the same initialization strategy as the Llama-1B configuration, however, we first create a shared initialization where the size of the embedding table—and the final output layer—is the size of the super vocabulary, $|E_{\rm sv}|=|\mathcal{SV}|$. Each model then uses the parameter values from this shared initialization for most layers. The embedding table for an individual model, E, is initialized by selecting the appropriate rows from the super vocabulary embedding table. Thus after initialization, $E(x)=E_{\rm sv}(sv(X))$. This results in a shared initialization for all models, including the initial embedding value for any shared tokens.

B.2 MODEL PERFORMANCE

We evaluate all models on standard English reasoning tasks (HellaSwag (Zellers et al., 2019), ARC Easy/Challenge (Clark et al., 2018), PIQA (Bisk et al., 2020)), multilingual natural language inference (XNLI (Conneau et al., 2018) in English, Turkish, and Chinese), reading comprehension (Belebele (Bandarkar et al., 2024) in English, Italian, Farsi, Turkish, and Chinese), and a multilingual reasoning benchmark (Include Base 44 (Romanou et al., 2025) in Chinese, Italian, and Turkish) in Fig. 1. Although models achieve sufficient performance on easier English reasoning tasks, their performance on the multilingual benchmarks hardly exceeds the random baseline. Note that models with larger vocabulary (Aya, XGLM, mBERT, Gemma-2, GPT-4o, and Llama-3.2) tend to perform better on the downstream tasks with TokenMonster and Tekken falling slightly behind.

C Intrinsic Tokenization Efficiency Metrics

Tokenizers exhibit varying degrees of compactness when segmenting text into tokens, resulting in notable disparities in model performance across languages and domains. To systematically evaluate these differences, we analyze several metrics across our selected pretrained tokenizers, focusing on our five languages.

We compute three primary intrinsic efficiency metrics using 10,000 parallel random samples from Flores200 (Team et al., 2022), split into "real" words via language-specific word-level tokenizers from the DataTrove library (Penedo et al.):

- Subword fertility (SF): is the mean number of tokens used to represent each "real" text word. This reflects how aggressively a tokenizer segments words. The theoretical minimum is 1, implying that the tokenizer's vocabulary encompasses every word in the reference text (Penedo et al., 2025).
- Parity: evaluates whether a tokenizer processes equivalent sentences fairly across languages. Achieved when the ratio of tokenized lengths $\frac{|T(s_A)|}{|T(s_B)|} \approx 1$ for parallel sentence sets s_A and s_B from languages A and B (Ali et al., 2024).
- Proportion of continued words (PCW): is the proportion of "real" text words that require two or more tokens for encoding. This metric indicates how frequently a tokenizer splits words. A score of 0 means no splitting occurs, while a score of 1 means every word is split (Rust et al., 2020).

The intrinsic metrics reflect a tokenizer's efficiency in processing a language and are critical factors in tokenizer selection, as they directly impact an LM's computational cost, context window utilization, and representation quality. Table 4 reveals substantial disparities in how our tokenizers handle our target languages. ByT5 and tokenizers with smaller vocabularies (TokenMonster, and Phi-3) exhibit significantly higher subword fertility and PCW scores, particularly for non-English languages—ByT5 requires 7.72 tokens per word in Farsi compared to 4.40 in English. Multilingual-specialized tokenizers (mBERT, XGLM) demonstrate superior language parity, with XGLM achieving near-optimal parity scores (1.18 average) and mBERT showing the lowest average subword fertility (1.54).

Notably, vocabulary size alone does not guarantee efficiency; Qwen-3 and Gemma-2, despite having large vocabularies (>150K), show comparable or worse performance than smaller vocabulary tokenizers like mBERT on certain metrics. We also observe higher fertility and PCW scores for morphologically rich languages (Turkish, Farsi) compared to English.

D TOKSUITE BENCHMARK DETAILS

D.1 QUESTION STYLE AND DIFFICULTY

The TokSuite benchmark comprises straightforward multiple-choice text completion questions. Below we present the canonical English questions that form our English subset, which are subsequently translated into Farsi (FA), Italian (IT), Turkish (TR), and Chinese (ZH). The fourteen models demonstrate strong performance on the canonical questions in English and Italian (Fig. 3), while the canonical accuracy on Farsi, Turkish, and Chinese is slightly behind. Higher PCW scores in these three languages (see Table 2) suggest that the models are likely to consume less information measured in raw bytes in these languages.

- 300 Dr Smith is a doctor. Occupation of Dr Smith is: **doctor**, teacher, judge, lawyer
- 301 The color of the sky is: **blue**, red, green, yellow
- 302 The price of this house is 1,028,415 dollars. The cost of this house is: **1,028,415 dollars**, 1.028,415 dollars, 1,028,415 dollars
- 303 Today's date is 29/08/2025. Today is: **29/08/2025**, 19/08/2025, 26/08/2025, 29/09/2025
- 304 The number of continents on Earth is: 7, 5, 6, 8
- 305 The capital city of Iran is: **Tehran**, Mashhad, Baghdad, Isfahan
- 306 The number of days in a week is: 7, 5, 6, 8
- 307 The number of hours in a day is: **24**, 20, 25, 30
- 308 The number of legs a cow has is: 4, 8, 3, 5
 - 309 The number of minutes in 2 hours is: **120**, 100, 140, 90
 - 310 The number of months in a year is: **12**, 10, 11, 13

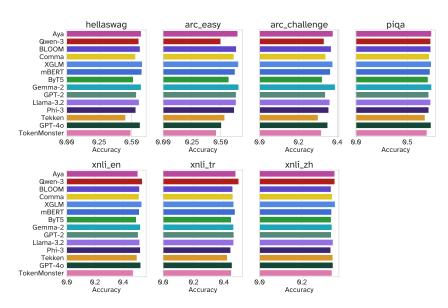


Figure 1: Model Performance on Multilingual Benchmarks

Table 4: Multilingual Tokenizers Comparison on Flores200 Using Intrinsic Tokenizer Efficiency Metrics. sf denotes subword fertility, pcw denotes proportion of continued words, and parity is measured against English parallel samples. Summary statistics report average values across all languages. Lower is better for all metrics. Bold font highlights the best performance in each row. Models are ordered from smallest to largest vocabulary size, left to right. Vocabulary size is categorized as XS, S, M, and L for <1K,1K-50K,50K-150K, and >150K tokens, respectively.

	1		4.											
		Potenty.	Se					^					,	٠.
	,	2	ò,	.0	a ²		- Nama	, [,] &	, ç	, Je		7		, ' z
m		68°	Tity	وعربه	Carara		- Salar	Legger S. Legger	Order.	ER LA		30	C SERVE	
Tokenizer Vocab, Size	XS	S	S	M	M	M M	∼ M	M	L	L	⋄ L	T.	L	¬⊽ L
vocab. Size	10			171	171	171	171		ட	ட		L		
English sf	4.40	1.75	1.24	1.30	1.44	1.15	1.26	1.35	1.28	1.24	1.31	1.19	1.14	1.23
English pcw	0.87	0.56	0.16	0.23	0.34	0.10	0.20	0.27	0.21	0.20	0.25	0.15	0.11	0.21
Chinese sf	5.00	4.92	3.44	3.54	2.45	1.68	1.49	1.64	1.21	1.44	1.16	1.23	1.28	2.19
Chinese pcw	0.98	0.97	0.97	0.82	0.58	0.55	0.35	0.41	0.16	0.32	0.13	0.18	0.21	0.87
Chinese parity	0.94	4.99	2.03	3.21	1.94	1.40	1.29	1.43	1.02	1.27	0.93	1.05	1.09	1.15
Turkish sf	6.49	4.31	3.20	3.20	3.29	1.99	2.38	2.44	2.58	2.33	2.71	2.17	2.23	1.69
Turkish pcw	0.49	0.80	0.76	0.76	0.78	0.52	0.72	0.73	0.74	0.71	0.72	0.68	0.69	0.52
Turkish parity	1.12	3.34	2.11	2.45	2.21	1.37	1.39	1.50	1.63	1.43	1.98	1.21	1.39	1.12
rannon panty	1,11	0.0.		2		1.07	1.07	1.00	1100	11.0	1.70	1.21	1.07	
Farsi sf	7.72	7.74	4.77	4.91	4.43	1.53	1.94	1.92	2.45	1.93	2.01	1.85	1.83	1.36
Farsi pcw	0.95	0.94	0.93	0.90	0.90	0.31	0.58	0.58	0.67	0.57	0.58	0.53	0.53	0.28
Farsi parity	1.72	9.45	4.08	5.35	4.31	1.38	1.52	1.47	2.63	1.55	1.80	1.48	1.45	1.21
Italian sf	4.78	2.50	1.64	1.99	2.05	1.34	1.81	1.77	1.83	1.71	1.75	1.61	1.54	1.36
Italian pcw	0.84	0.63	0.42	0.57	0.59	0.23	0.55	0.53	0.55	0.52	0.51	0.47	0.41	0.32
Italian parity	1.19	2.30	1.48	2.02	1.87	1.28	1.62	1.40	1.64	1.47	1.63	1.31	1.33	1.24
Avg sf	5.79	4.39	2.90	3.19	2.93	1.54	1.78	1.82	1.87	1.73	1.79	1.61	1.60	1.56
Avg si Avg pcw	0.90	0.78	0.62	0.66	0.64	0.34	0.48	0.50	0.47	0.46	0.44	0.40	0.39	0.46
Avg pew Avg parity	1.27	5.31	2.54	3.44	2.74	1.36	1.46	1.45	1.73	1.43	1.59	1.26	1.32	1.18
	1 1.27	0.01		2.11	, .	1.50	1.10	1.10	1.75	1.10	1.07	1.20		

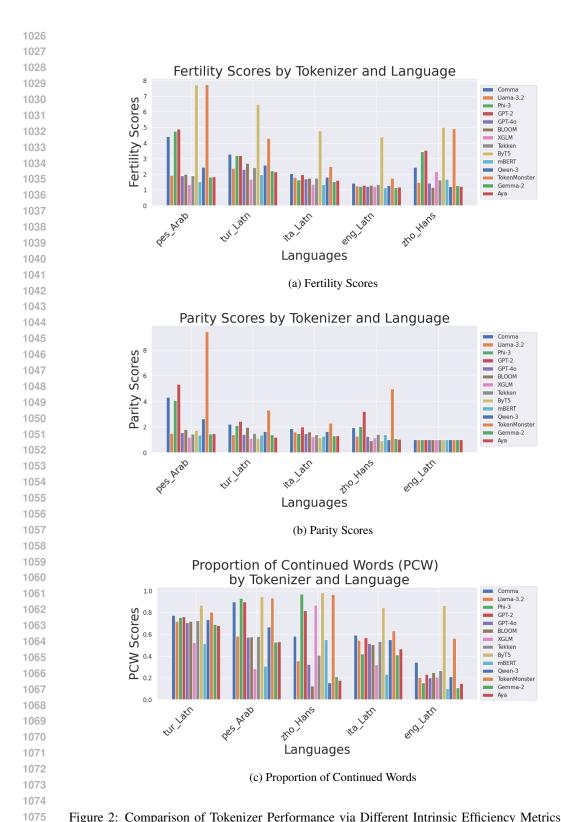


Figure 2: Comparison of Tokenizer Performance via Different Intrinsic Efficiency Metrics Using Flores200 dataset. (a) shows subword fertility, (b) shows parity against English, and (c) shows the proportion of continued words across languages and tokenizers. Lower values indicate better comparative performance.

1080 311 The number of seconds in a minute is: **60**, 50, 100, 30 1081 312 The number of sides a hexagon has is: 6, 5, 7, 8 1082 313 The number of sides a triangle has is: 3, 2, 4, 5 1084 314 In "I work at Apple", Apple is: company, person, city, fruit 315 In "I work at Google", Google is: company, person, city, fruit 1087 316 In "Microsoft released a new update", Microsoft is: company, person, place, date 1088 317 In "The cat sat on the mat", the subject is: the cat, sat, the mat, on 1089 322 The gas humans need to breathe to live is: **oxygen**, methane, helium, hydrogen 1090 1091 323 10% of 100 is: **10**, 5, 15, 20 324 25% of 80 is: **20**, 15, 25, 30 1093 326 Chad's capital is: **N'Djamena**, Moundou, Abéché, Ngama 1094 1095 327 The capital of France is: **Paris**, London, Berlin, Rome 328 The capital of Japan is: **Tokyo**, Kyoto, Osaka, Hiroshima 329 The capital of Turkey is: **Ankara**, İstanbul, İzmir, Bursa 1099 330 The chemical formula for water is: **H2O**, CO2, NaCl, O2 1100 331 The intent in "What time does the store close?" is: get information, make purchase, book 1101 appointment, file complaint 1102 332 The largest mammal in the world is: **blue whale**, dolphin, giraffe, bear 1103 1104 333 The unit of measurement for temperature in the International System is: **Kelvin**, Celsius, 1105 meter, Rankine 1106 334 The country whose space agency is NASA is: United States, Russia, China, Japan 1107 335 The language spoken in Brazil is: **Portuguese**, Spanish, French, Italian 1108 1109 336 The metal with chemical symbol 'Fe' is: **iron**, lead, zinc, gold 1110 337 The organ in the human body that pumps blood is: **heart**, liver, lungs, kidneys 1111 338 The planet closest to the Sun in our solar system is: Mercury, Venus, Mars, Earth 1112 1113 339 The largest planet in the Solar System is: Jupiter, Earth, Saturn, Mars 1114 340 The process that allows plants to produce their own food using sunlight is: **photosynthesis**, 1115 respiration, digestion, fermentation 1116 341 The author who wrote the play "Romeo and Juliet" is: William Shakespeare, Charles 1117 Dickens, Mark Twain, Jane Austen 1118 1119 342 What bees produce is: **honey**, milk, silk, wax 1120 343 What plants need from the air to make food is: **carbon dioxide**, nitrogen, hydrogen, helium 1121

D.2 BENCHMARK COMPOSITION

1122

112311241125

1126 1127

1128

1129

1130 1131

11321133

In Table 5, we list the composition of the categories and perturbations in TokSuite. The multilingual parallel dataset comprises 80% of the dataset, while the remaining part covers math, STEM, and general questions.

344 In "Can you please book a flight to Paris?", the person wants to: make a booking, go

E DETAILED BENCHMARK RESULTS

In this section, we provide case studies for each category in Section 4.1.

shopping, file a complaint, cancel reservation

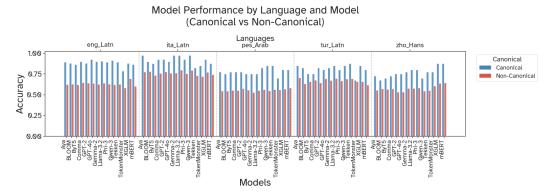


Figure 3: Accuracies of models on canonical versus perturbed questions across the English (eng_Latn), Italian (ita_Latn), Farsi (pes_Arab), Turkish (tur_Latn), and Chinese (zho_Hans) TokSuite subsets.

Table 5: Benchmark statistics by language and domain

Language/Domain	Total Examples	Perturbations
English	1,180	42 types
Chinese	485	18 types
Turkish	638	21 types
Italian	1,088	19 types
Farsi	747	15 types
Math	189	5 types
STEM	614	25 types
General	89	4 types

Table 6: Tokenization robustness under different input mediums, granular version of Input in Table 1. Values represent relative performance drop (canonical-perturbed)/canonical; lower values indicate greater robustness. Traditional: Traditional Chinese characters instead of simplified.

Model	Roman	ization	Number Romaniza- tion	English	Keyboard	Arabic Keyboard	Traditional	Avg
	FA	ZH	FA	TR	IT	FA	ZH	
TokenMonster	0.46	0.58	-0.10	-0.04	0.21	0.25	0.02	0.20
Comma	0.42	0.59	0.21	0.03	0.24	0.42	0.04	0.28
GPT-40	0.57	0.67	-0.03	0.22	0.09	0.43	0.03	0.28
Llama-3.2	0.60	0.66	-0.23	0.24	0.14	0.53	0.09	0.29
BLOOM	0.63	0.48	0.08	0.21	0.15	0.40	0.10	0.29
Aya	0.55	0.62	0.01	0.06	0.16	0.55	0.12	0.29
ByT5	0.61	0.46	0.21	0.13	0.15	0.39	0.18	0.30
Tekken	0.59	0.61	0.00	0.17	0.20	0.44	0.18	0.31
Gemma-2	0.40	0.52	0.28	0.24	0.19	0.47	0.18	0.32
Phi-3	0.58	0.66	0.25	0.06	0.24	0.39	0.09	0.33
XGLM	0.59	0.63	0.13	0.29	0.19	0.41	0.10	0.34
mBERT	0.44	0.60	0.42	0.22	0.18	0.50	0.10	0.35
GPT-2	0.61	0.67	0.31	0.30	0.16	0.32	0.11	0.35
Qwen-3	0.68	0.64	0.19	0.15	0.19	0.47	0.18	0.36
Avg	0.55	0.60	0.12	0.16	0.18	0.43	0.11	0.31

E.1 ORTHOGRAPHIC & SCRIPT CHALLENGES

Table 6 examines tokenization robustness under orthographic and script challenges, focusing on input medium variations where users employ non-native keyboards or writing systems. For Chinese romanization, we write the full question and choices in Pinyin without tone markers—as if the user only has access to an English keyboard—with spaces between individual groups that constitute a character for easy segmentation. However, this segmentation aid does not improve tokenization robustness, as models still exhibit substantial performance degradation (0.60 average drop) when processing romanized Chinese text compared to native scripts. In Table 7, the errors due to input systems (like homoglyphs and zero-width characters) are presented.

Diacritics perturbations systematically vary the presence and accuracy of accent marks and diacritical symbols. We test how tokenizers handle optional diacritics, where text remains valid with or without marks (e.g., marks placed above or below letters to clarify pronunciation and phonetic details such as short vowels (fata for /a/, kasra for /e/, amma for /o/), or sukūn for the absence of vowels in Farsi), wrong accents such as using é instead of è in Italian. Table 8 expands on diacritics perturbations, examining how tokenizers handle optional Farsi diacritics that are used to clarify pronunciation and phonetic details, Chinese tonal variations in the Pinyin format, and incorrect accent placement in Italian text. We see that Tokenmonster, which decomposes these markers shows strong performance with Farsi optional diacritics.

In Table 14, we also present Pinyin input without these optional whitespaces and observe that the space-removed version causes less performance degradation. In Table 8, we further demonstrate that adding tone markers significantly improves performance. While we expected adding spaces to help models identify corresponding Chinese characters better, the opposite occurred—this likely reflects how native speakers typically write Pinyin without spaces in practice, making the spaced version appear more artificial to models trained on naturalistic text data.

Orthographic and Grammatical Errors Table 9 reveals that orthographic and grammatical errors create varying challenges depending on the morphological complexity of the language. Token-Monster demonstrates the strongest, while character-level approaches like ByT5 show competitive performance across multiple categories.

Orthographic Errors Imagine perturbing the word "week" to "weak" in the question, "The number of days in a week is". This change breaks 6/14 models despite both words existing as distinct tokens with separate embeddings. This suggests that tokenization robustness depends not merely on vocabulary coverage but on the semantic stability of token representations.

Grammatical Errors Consider the Turkish locative suffix variants "saatteki" for the root saat (in the hour) versus the incorrect "saatdeki" as part of the canonical question "The number of minutes in 2 hours is" (TR: 2 saatteki dakika sayısı).

This example demonstrates how agglutinative languages amplify tokenization brittleness: a single phoneme change (/t/ to /d/) can completely restructure token boundaries. This reflects the curse of multilinguality, where tokenizers trained predominantly on English struggle with morphologically complex languages, sometimes producing cleaner segmentation—with meaningful morphemes—for incorrect forms than correct ones (as Gemma-2 and BLOOM below). English grammatical errors on the other hand—with wrong prepositions, subject-verb agreement, etc—tend to change token boundaries less and we observe a less striking performance degradation in Table 9.

Assimilation error ("saatteki" vs. "saatdeki"):

- BLOOM, Gemma-2: sa, atte, ki vs. saat, de, ki (meaningful morphemes after error)
- XGLM: saat, teki vs. saat, deki (clean morpheme separation)
- Llama-3.2: sa, atte, ki vs. sa, at, deki (inconsistent segmentation)
- mBERT: saat, ##tek, ##i vs. saat, ##deki (subword fragmentation changes)
- Qwen-3: sa, atte, ki vs. sa, at, de, ki (boundary reorganization)
- TokenMonster: sa, at, tek, ivs. sa, a, td, ek, i(severe fragmentation)

Table 7: Tokenization robustness under errors from input mediums. Values represent relative performance drop (canonical-perturbed)/canonical; lower values indicate greater robustness.

Model	Homoglyphs EN	Zero-wid	th chars.	Avg
	1			
mBERT	0.08		0.00	0.06
Phi-3	0.03	0.21	-0.06	0.06
TokenMonster	0.09	0.18	-0.06	0.07
BLOOM	0.12	0.17	-0.07	0.07
XGLM	0.03	0.19	0.03	0.08
ByT5	0.06	0.32	-0.11	0.09
Comma	0.05	0.32	-0.07	0.10
GPT-40	0.14	0.23	-0.03	0.11
Aya	0.28	0.23	-0.14	0.12
Gemma-2	0.15	0.27	0.03	0.15
Llama-3.2	0.12	0.30	0.03	0.15
GPT-2	0.13	0.23	0.13	0.16
Tekken	0.13	0.29	0.10	0.17
Qwen-3	0.11	0.38	0.11	0.20
Avg	0.11	0.24	-0.01	0.11

Table 8: Tokenization robustness to diacritics, granular version of Diacr in Table 1 and wrong accents in Italian. Values represent relative performance drop (canonical-perturbed)/canonical; lower values indicate greater robustness.

Model	Dia	critics	Wrong	Avg
	FA	ZH	accents IT	
BLOOM	0.33	0.37	0.08	0.26
TokenMonster	0.21	0.45	0.17	0.28
GPT-2	0.42	0.50	-0.02	0.30
Qwen-3	0.41	0.43	0.10	0.31
ByT5	0.42	0.46	0.06	0.31
mBERT	0.31	0.57	0.06	0.31
Gemma-2	0.43	0.42	0.10	0.32
Phi-3	0.39	0.53	0.05	0.32
Tekken	0.47	0.48	0.07	0.34
Aya	0.45	0.48	0.10	0.34
XGLM	0.44	0.54	0.11	0.36
GPT-4o	0.47	0.57	0.08	0.37
Comma	0.39	0.48	0.30	0.39
Llama-3.2	0.60	0.50	0.16	0.42
Avg	0.41	0.49	0.10	0.33

Table 9: Tokenization robustness under orthographic and grammatical errors. Values represent relative performance drop (canonical-perturbed)/canonical; lower values indicate greater robustness.

Model	Or	thographic E	Errors	G	rammatical l	Errors	Phonetic	Avg
	EN	TR	IT	EN	TR	IT	IT	_
TokenMonster	0.10	0.04	0.04	0.06	0.03	-0.03	0.04	0.04
ByT5	0.06	0.10	0.08	0.00	-0.01	0.04	0.02	0.04
GPT-40	0.12	0.13	0.08	0.00	0.05	-0.01	0.02	0.06
Comma	0.09	0.20	0.06	-0.03	0.13	0.01	0.04	0.07
Llama-3.2	0.14	0.18	0.13	0.05	0.07	0.03	0.02	0.09
Tekken	0.24	0.23	-0.01	0.08	0.21	-0.07	-0.01	0.09
GPT-2	0.08	0.30	0.10	0.05	0.12	0.01	0.09	0.11
BLOOM	0.18	0.24	0.05	0.03	0.21	-0.01	0.07	0.11
Qwen-3	0.17	0.18	0.12	0.08	0.15	0.05	0.02	0.11
Phi-3	0.18	0.22	0.13	0.11	0.09	-0.02	0.07	0.11
Aya	0.21	0.21	0.13	0.03	0.07	0.02	0.14	0.11
mBERT	0.15	0.41	0.08	0.03	0.22	-0.02	0.04	0.13
XGLM	0.13	0.32	0.12	0.03	0.23	-0.02	0.15	0.14
Gemma-2	0.18	0.30	0.12	0.05	0.29	0.07	0.09	0.16
Avg	0.14	0.22	0.09	0.04	0.13	0.00	0.06	0.10

```
• GPT-40: s, aat, te, ki vs. s, aat, de, ki (character-level consistency)
• Tekken: sa, atte, ki vs. sa, at, deki (partial boundary preservation)
```

• GPT-2: sa, at, te, kivs. sa, at, d, eki (fine-grained segmentation)

Turkish final-obstruent devoicing error ("ineğin" → "inekin") in the word cow's (possesive)

```
• BLOOM: ine, \zeta, \S, in vs. in, ekin
```

• XGLM: in, e, ğ, in vs. in, ekin

• Llama-3.2: ine, ζ , \S , in vs. ine, kin

• mBERT: [UNK] vs. in, ##ekin (unknown token fallback)

• Qwen-3: ine, Ç\$, in vs. ine, kin

• TokenMonster: ine, g, i, n vs. ine, kin (diacritic decomposition)

• Gemma-2: ine, ğ, in vs. ine, kin

• GPT-40: ine, ğ, in vs. ine, kin

• Tekken: ine, ğ, in vs. ine, kin

• GPT-2: ine, ğ, in vs. ine, kin

Register and style variations compound tokenization challenges. Consider using emoji substitution in "The capital of Japan is" by replacing "Japan" with the Japanese flag.

Table 10: Tokenization robustness under different register and style variations. Values represent relative performance drop (canonical-perturbed)/canonical; lower values indicate greater robustness. Abb.: abbreviations, Word Ord.: word reordering, emoji: emoji substituion, char. subs.: character substitution, repet.: letter repetition for emphasis

Model	W	eb Se	arch		Abb.	Wo	rd Ord.	Phone	tic	Co	lloqui	ial	Emoji	Char. Subs.	Repet.	Avg
	EN	TR	IT	EN	IT	EN	TR	IT	EN	FA	TR	ZH	EN	EN	EN	
TokenMonster	0.26	0.07	0.38	0.32	0.04	0.06	-	0.04	0.11	0.00		0.04	0.25	-	0.22	0.11
mBERT	0.33	0.25	0.23	0.27	0.07	0.08	0.01 0.18	0.04	0.15	0.09	0.00 0.12	0.18	0.29	0.07	0.18	0.16
GPT-40	0.36	0.34	0.53	0.18	0.09	0.05	0.03	0.02	0.20	0.10	0.12	0.15	0.16	0.08	0.21	0.17
ByT5	0.40	0.30	0.29	0.28	0.11	0.06	0.12	0.02	0.15	0.19	0.14	0.16	0.32	0.01	0.11	0.17
Comma	0.43	0.33	0.43	0.32	0.08	_	0.03	0.04	0.12	0.13	0.14	0.19	0.23	0.04 0.01	0.13	0.17
						0.03										
BLOOM GPT-2	0.41 0.29			0.24 0.20	0.09 0.16	0.12 0.13	0.20 0.15	0.07 0.09		0.20			0.20 0.26	0.00	0.17 0.28	0.18 0.19
XGLM	0.29	0.32	0.30	0.29	0.16	0.03	0.17	0.15	0.20	0.22	0.17	0.15	0.33	0.05 0.01	0.08	0.19
Llama-3.2 Tekken	0.38 0.49			0.30 0.29	0.13	0.10 0.05	0.14 0.19	0.02		0.17			0.25 0.26	0.06 0.01	0.27 0.20	0.20
TERREII								0.01						0.01		
Aya	0.42	0.38	0.33	0.28	0.24	0.08	0.20	0.14	0.17	0.13	0.11	0.15	0.11	0.03	0.32	0.20
Qwen-3	0.32	0.41	0.49	0.26	0.03	0.08	0.17	0.02	0.14	0.32	0.17	0.16	0.14	0.08	0.36	0.21
Gemma-2	0.50			0.25	0.28	0.08	0.15	0.09		0.07			0.18	0.04	0.20	0.22
Phi-3	0.43	0.31	0.02	0.20	0.04	0.11	0.15	0.07	0.24	0.21	0.19	0.23	0.33	0.05	0.28	0.22
Avg	0.38	0.32	0.40	0.26	0.11	0.07	0.13	0.06	0.16	0.15	0.13	0.16	0.24	0.01	0.21	0.19

Emoji handling reveals differences: Most modern tokenizers like Gemma-2, GPT-40, Tekken, GPT-2, and Qwen-3 have emojis in their vocabulary, correctly parse the Japanese flag emoji into two tokens as the corresponding regional indicators ([J] and [P]). Aya on the other hand has a standalone token for the flag emoji. BLOOM, Llama-3.2, and TokenMonster use byte-fallback, XGLM and mBERT resort to unknown tokens. The coverage of emojis translate into good performance in the Emoji substitution perturbations (see Table 10).

Linguistic Variety Table 11 examines how tokenizers handle linguistic diversity including historical spellings, code-switching, dialects, and colloquial expressions. TokenMonster demonstrates remarkable consistency across varied linguistic phenomena (0.08 average drop), while most models struggle significantly with certain types of variation. In Table 12, we group the models based on their vocabulary size (see Table 2) to investigate potential correlations with vocabulary size, as larger vocabularies theoretically provide more comprehensive dictionaries.

Counterintuitively, vocabulary size shows little to no correlation with linguistic robustness—byte-level model (ByT5) demonstrates superior consistency despite operating without traditional vocabulary constraints, while some large-vocabulary tokenizers exhibit significant brittleness. We observe that larger vocabulary size doesn't always produce a lexically-rich vocabulary. Modern tokenizers may actually compound the problem by learning multiple variants of common words (Gemma-2 has distinct tokens for "hello", "hello", "Hello", and "Hello"), reducing the effective vocabulary. While this multiplicity has efficiency gains it could make models sensitive to stylistic variations that should be semantically equivalent.

Historical spelling variants ("capitall"³, "Japane") demonstrate systematic fragmentation patterns where tokenizers often segment archaic or non-standard spellings along morphological boundaries:

- Most tokenizers: capit, all and Jap, ane (consistent morpheme-like splitting)
- mBERT: capital, ##l and Japan, ##e (subword suffix handling)
- XGLM: capital, land Japan, e (clean separation)

Colloquial expressions reveal deeper challenges in world knowledge representation. The question "Turkey's capital turns out to be" with the correct answer "Ankara" illustrates how informal phrasing can disrupt factual recall: as it breaks 3 models. This suggests that tokenizers' handling of casual discourse markers and words ("turns out to be") may interfere with models' access to factual knowledge. The pattern indicates that linguistic variety challenges extend beyond mere tokenization to fundamental issues of how models integrate linguistic style with semantic content.

E.2 MORPHOLOGICAL CHALLENGES

Table 13 examines how tokenizers handle morphological variations including derivations, inflections, and contractions across English, Turkish, and Italian. Morphological perturbations reveal fundamental inconsistencies in how tokenizers segment related word forms—contractions like "Google's" versus decomposed forms, or Italian elision patterns where "dell'Italia" and "d'Italia" receive dramatically different tokenization despite identical meaning. These inconsistencies suggest that current tokenization approaches lack coherent strategies for handling morphologically related forms, potentially leading models to develop disparate semantic representations for linguistically equivalent expressions. For example while BLOOM learns contractions, GPT-2 and GPT-40 use a regex-based search.

English Contractions: "Google is"→ "Google's"

- BLOOM, Llama-3.2, Qwen-3, Gemma-2, GPT-2, GPT-40, Tekken,: Google, 's (separate marker)
- **XGLM**, **mBERT**: Google, ', s (fragmentation)
- **TokenMonster:** google, 's (lowercase normalization)

Italian Ellisions The Italian contraction "L'intento" (the intent) demonstrates varying approaches to handling elided articles:

```
• BLOOM: L', int, ento
• XGLM: L, ', inten, to
• Llama-3.2: L, 'int, ento
• mBERT: L, ', intento
```

³https://www.oed.com/search/dictionary/?scope=Entries&q=capitall

Table 11: Tokenization robustness under linguistic variety. Values represent relative performance drop (canonical-perturbed)/canonical; lower values indicate greater robustness. Hist.: historical spelling, equiv. exp.: equivalent expressions, sim. words: similar words

Model	Hist.		Cod	le swite	ch		Diale	ects		Equ	iiv. ex	р.	Si	m. wo	rds	Avg
	EN	FA	TR	IT	ZH	FA	TR	IT	EN	FA	TR	ZH	EN	TR	IT	
TokenMonster	0.09	0.07	0.00	0.00	0.03	0.22	0.09	0.17	0.14	0.07	0.04	0.03	0.03	-0.06	0.22	0.08
ByT5	0.06	0.03	0.04	0.06	-0.04	0.29	0.15	0.15	0.02	0.13	0.06	0.04	0.08	-0.08	0.24	0.08
Comma	0.21	0.10	0.13	0.06	0.03	0.30	0.04	0.06	-0.05	0.10	0.06	0.03	0.08	-0.02	0.28	0.09
BLOOM	0.25	-0.07	0.16	-0.03	-0.04	0.31	0.19	0.14	0.05	0.07	0.14	-0.07	0.09	0.13	0.26	0.11
mBERT	0.11	0.09	0.16	0.03	0.09	0.30	0.31	0.12	-0.05	0.06	0.04	0.06	0.02	0.23	0.05	0.11
Tekken	0.21	0.12	0.16	-0.03	0.03	0.37	0.14	-0.02	0.17	0.15	0.06	0.03	0.05	0.18	-0.01	0.11
GPT-40	0.08	-0.03	0.10	-0.08	0.07	0.29	0.10	0.14	0.14	-0.03	-0.03	0.13	0.05	0.29	0.44	0.11
XGLM	0.18	0.09	0.21	0.06	-0.03	0.30	0.15	0.02	0.17	0.03	0.10	0.09	0.08	0.16	0.10	0.11
Gemma-2	0.31	0.17	0.05	0.05	0.10	0.33	0.23	0.07	0.17	0.00	0.07	-0.10	0.04	0.08	0.40	0.13
Aya	0.21	0.03	0.13	0.08	0.03	0.30	0.18	0.14	0.27	0.16	0.10	0.00	0.07	0.10	0.23	0.14
GPT-2	0.18	0.10	0.18	0.06	0.20	0.28	0.23	0.23	0.07	0.10	0.14	0.03	0.09	0.08	0.10	0.14
Llama-3.2	0.25	0.03	0.13	0.03	0.09	0.24	0.05	0.17	0.10	0.03	0.17	0.19	0.09	0.16	0.40	0.14
Qwen-3	0.32	0.21	0.18	0.05	0.04	0.34	0.18	0.11	0.02	0.24	0.17	-0.07	0.09	0.22	0.15	0.15
Phi-3	0.32	0.12	0.16	0.09	0.13	0.35	0.10	0.23	-0.05	0.15	0.34	0.09	0.09	0.29	0.19	0.17
Avg	0.20	0.08	0.13	0.03	0.05	0.30	0.15	0.12	0.08	0.09	0.11	0.03	0.07	0.13	0.22	0.12

Table 12: Tokenization robustness under linguistic variety. Same as Table 11 but grouped under vocabulary size. Values represent relative performance drop (canonical-perturbed)/canonical; lower values indicate greater robustness. Hist.: historical spelling, equiv. exp.: equivalent expressions, sim. words: similar words

Vocab Size	Hist.	Code switch			Dialects			Equiv. exp.			Si	Sim. words				
	EN	FA	TR	IT	ZH	FA	TR	IT	EN	FA	TR	ZH	EN	TR	IT	
X-Small	0.06	0.03	0.04	0.06	-0.04	0.29	0.15	0.15	0.02	0.13	0.06	0.04	0.08	-0.08	0.24	0.08
Medium	0.19	0.09	0.15	0.03	0.09	0.30	0.15	0.12	0.05	0.09	0.10	0.07	0.07	0.13	0.17	0.12
Large	0.23	0.07	0.14	0.02	0.03	0.31	0.17	0.10	0.14	0.08	0.09	0.00	0.07	0.16	0.26	0.13
Small	0.21	0.10	0.09	0.05	0.08	0.29	0.10	0.20	0.04	0.11	0.20	0.06	0.06	0.13	0.20	0.13
Avg	0.17	0.07	0.11	0.04	0.04	0.30	0.14	0.14	0.06	0.10	0.11	0.04	0.07	0.08	0.22	0.11

```
1458
            • Qwen-3: L, 'int, ento
1459
            • TokenMonster: 1', intent, o
1460
            • Gemma-2: L, ', int, ento
1461
             • GPT-40: L, 'int, ento
1462
1463
             • Tekken: L, 'int, ento
1464
            • GPT-2: L, ', intent, o
1465
       "dell'Italia" vs. "d'Italia":
1466
1467
            • BLOOM: d, ell, ', Italia vs. d', Italia
1468
            • XGLM: dell, ', Italia vs. d, ', Italia
1469
            • Llama-3.2, Qwen-3: d, ell, 'It, alia vs. d, 'It, alia (fragments "Italia")
1470
             • mBERT: dell, ', Italia vs. d, ', Italia (length-dependent)
1472
            • TokenMonster: dell, ', ita, lia vs. d, ', ita, lia (lowercase + frag-
1473
1474
             • Gemma-2: dell, ', Italia vs. d, ', Italia (clean separation)
            • GPT-40: d, ell, ', Italia vs. d, ', Italia (inconsistent decomposition)
1476
             • Tekken: d, ell, 'Italia vs. d, 'Italia (treats apostrophe differently)
1477
            \bullet GPT\text{-}2\text{:} d, ell, ', It, alia vs. d, ', It, alia (fragments country name)
1478
```

Table 13: Tokenization robustness under morphological challenges, granular version of Morphological in Table 1. Values represent relative performance drop (canonical-perturbed)/canonical; lower values indicate greater robustness.

Model	Contractions EN IT	Compounds EN	Derivations TR	Inflections EN TR	Avg
Comma	0.23 0.18	0.09	-0.11	0.02 0.02	0.07
TokenMonster	0.30 0.16	0.17	-0.12	0.02 -0.09	0.07
GPT-2	0.33 -0.08	0.09	0.05	0.02 0.13	0.09
Aya	0.27 -0.03	0.19	0.02	0.05 0.06	0.10
Gemma-2	0.27 -0.03	0.14	0.02	0.12 0.06	0.10
mBERT	0.26 -0.14	0.09	0.18	0.15 0.06	0.10
Qwen-3	0.31 0.12	0.09	0.02	0.10 0.06	0.12
GPT-4o	0.26 0.26	0.12	-0.04	0.07 0.06	0.12
ByT5	0.30 -0.03	0.15	0.09	0.21 0.05	0.13
BLOOM	0.20 -0.01	0.16	0.11	0.14 0.16	0.13
XGLM	0.26 0.02	0.07	0.11	0.25 0.06	0.13
Llama-3.2	0.29 0.12	0.16	0.02	0.14 0.11	0.14
Tekken	0.36 -0.04	0.14	0.08	0.17 0.18	0.15
Phi-3	0.28 0.07	0.14	0.09	0.25 0.08	0.15
Avg	0.28 0.04	0.13	0.04	0.12 0.07	0.11

E.3 Noise

We observe that tokenizers that segment text into complete word tokens tend to exhibit greater vulnerability to noise errors, as single character perturbations can cause familiar words to fragment into unfamiliar subword combinations, whereas tokenizers using smaller subword units maintain more consistent segmentation patterns.

Noise in Chinese subset For keyboard proximity errors in Chinese characters are replaced with phonetically or positionally similar alternatives on the keyboard layout. For space removal, we use the Pinyin input without any spaces.

Typos Typographical errors demonstrate how different tokenization approaches handle character-level perturbations. For example, the word "doctor" with a typo becomes "doctro":

```
mBERT: doctor, doc, ##tro
Comma AI: do, ctor, Ldoc, tro
Llama-3: doctor, Ldo, ct, ro
Tekken: doctor, doct, ro
Aya Expanse: doctor, Ldoct, ro
GPT-40: doctor, doct, ro
GPT-2: doctor, doct, ro
ByT5: d, o, c, t, o, r, L, d, o, c, t, r, o
```

Similarly, for Turkish text "gün sayısı" (day count) with spacing errors becoming "güns ayısı":

```
mBERT: gün, sayısı, gün, ##s, ay, ##ısı
Comma AI: g, ün, Lsay, ı, s, ı, Lg, ü, ns, Lay, ı, s, ı
Tekken: g, ün, say, ısı, gün, s, ay, ısı
GPT-40: g, ün, say, ısı, gün, s, ay, ısı
Llama-3.2: gün, Lsayısı, Lgü, ns, Lay, ısı
GPT-2: g, ü, n, say, ı, s, ı, g, ü, ns, ay, ı, s, ı
Aya Expanse: gün, Lsayısı, Lgün, s, Lay, ısı
ByT5: Character-level segmentation (individual Unicode characters)
```

E.4 MATHEMATICAL & SCIENTIFIC EXPRESSIONS

Table 15 demonstrates that models generally struggle with the formatting and structural challenges inherent in scientific domains. When numerical values are replaced with their spelled-out equivalents (15 \rightarrow fifteen), we observe a consistent performance degradation even in English. The parallel multilingual basic arithmetic questions reveal that certain tokenizers may exhibit inductive biases favoring specific languages. For instance, Gemma-2's performance on Italian questions matches that of the canonical English questions, whereas it shows a 53% performance degradation in Farsi. Llama-3.2 demonstrates similar behavior with Turkish, while the Aya tokenizer, developed as part of a multilingual language model, exhibits the greatest robustness across languages. It should be noted, however, that this represents one of the few instances in our study where Aya tokenizer demonstrates clear multilingual advantages.

Tokenization of scientific text: Consider the unit "cubic meters" expressed as m^3 , m^3 , m^3 , m^3 , and m^6 3 \\$. Despite semantic equivalence, tokenization patterns reveal increasing fragmentation:

• BLOOM:

```
Plain: m, ^3
LaTeX: $m, ^3, $
Braced: $m, ^{3, }$
Spaced: $m, ^{, 3, }$
XGLM:
Plain: m, ^, 3
LaTeX: $, m, ^, 3, $
Braced: $, m, ^, {, 3, }, $
Spaced: $, m, ^, {, 3, }, $
```

• Llama-3.2:

Table 14: Tokenization robustness under multi-lingual noise. Values represent relative performance drop (canonical-perturbed)/canonical; lower values indicate greater robustness.

Model	Keyboard Errors		(OCR	Char. Del.	Space	Removal		Typos		Avg			
	EN	FA	TR	IT	ZH	EN	ZH	EN	EN	ZH	EN	TR	IT	
Comma	0.05	0.29	0.15	0.18	0.17	0.12	0.10	0.10	0.14	0.55	0.20	0.04	0.25	0.18
ByT5	0.19	0.26	0.13	0.22	0.11	0.18	0.11	0.09	0.18	0.43	0.17	0.11	0.18	0.18
TokenMonster	0.22	0.18	0.15	0.13	0.16	0.10	0.26	0.04	0.13	0.58	0.08	0.09	0.25	0.18
GPT-2	0.20	0.16	0.29	0.16	0.27	0.15	0.23	0.09	0.16	0.50	0.18	0.22	0.20	0.22
Qwen-3	0.20	0.32	0.25	0.19	0.11	0.15	0.25	0.12	0.17	0.43	0.23	0.16	0.26	0.22
GPT-40	0.13	0.20	0.13	0.13	0.23	0.15	0.40	0.18	0.16	0.53	0.24	0.13	0.22	0.22
BLOOM	0.22	0.23	0.34	0.16	0.11	0.19	0.11	0.16	0.21	0.56	0.16	0.25	0.17	0.22
Gemma-2	0.17	0.23	0.21	0.22	0.19	0.17	0.29	0.16	0.15	0.52	0.14	0.13	0.30	0.22
Llama-3.2	0.12	0.30	0.26	0.21	0.19	0.10	0.28	0.17	0.20	0.56	0.08	0.22	0.24	0.22
XGLM	0.18	0.25	0.29	0.19	0.23	0.15	0.29	0.13	0.13	0.60	0.11	0.22	0.21	0.23
Tekken	0.23	0.29	0.33	0.12	0.26	0.20	0.29	0.11	0.12	0.52	0.11	0.21	0.20	0.23
Phi-3	0.15	0.27	0.22	0.20	0.22	0.20	0.22	0.21	0.18	0.53	0.20	0.20	0.21	0.23
mBERT	0.24	0.25	0.32	0.16	0.14	0.20	0.20	0.14	0.24	0.60	0.11	0.23	0.26	0.24
Aya	0.15	0.42	0.25	0.26	0.24	0.17	0.28	0.19	0.21	0.52	0.10	0.19	0.27	0.25
Avg	0.18	0.26	0.24	0.18	0.19	0.16	0.24	0.13	0.17	0.53	0.15	0.17	0.23	0.22

Table 15: Tokenization robustness under math and STEM related challenges. Values represent relative performance drop (canonical-perturbed)/canonical; lower values indicate greater robustness. LaTeX: LaTeX-style math formatting; Diag. scientific diagrams and notations; Unic.: Unicode formatted ASCII characters. NEN=non-English.

Model	LaTeX			Spelled	Out		Diag.		Mul	tilingua	ı	Unicode	Avg
	EN	EN	FA	TR	IT	ZH	EN	FA	TR	ΙΤ	ZH	EN	Ü
TokenMonster	0.23	0.28	0.49	0.07	0.33	0.31	0.11	0.29	0.00	0.14	0.00	0.08	0.19
Phi-3	0.25	0.34	0.39	0.14	0.47	0.23	0.22	0.29	0.00	0.00	0.24	0.11	0.22
Aya	0.23	0.32	0.35	0.41	0.47	0.26	0.38	0.07	0.00	0.00	0.00	0.21	0.23
mBERT	0.15	0.35	0.55	0.45	0.35	0.38	0.22	0.14	0.07	0.14	0.07	0.23	0.26
Llama-3.2	0.18	0.33	0.43	0.34	0.45	0.23	0.29	0.18	0.47	0.00	0.18	0.07	0.26
GPT-2	0.25	0.38	0.35	0.32	0.44	0.08	0.35	0.18	0.35	0.24	0.24	0.17	0.28
Tekken	0.27	0.37	0.33	0.36	0.38	0.31	0.44	0.18	0.24	0.12	0.24	0.15	0.28
BLOOM	0.25	0.29	0.24	0.47	0.40	0.20	0.11	0.41	0.35	0.24	0.29	0.19	0.29
Comma	0.23	0.36	0.54	0.17	0.47	0.26	0.29	0.39	0.28	0.17	0.22	0.19	0.30
ByT5	0.18	0.37	0.54	0.42	0.54	0.23	0.29	0.07	0.20	0.27	0.27	0.23	0.30
GPT-40	0.25	0.38	0.33	0.45	0.52	0.28	0.33	0.37	0.32	0.05	0.16	0.20	0.30
Gemma-2	0.22	0.35	0.33	0.32	0.53	0.40	0.37	0.53	0.35	0.00	0.18	0.23	0.32
Qwen-3	0.26	0.41	0.50	0.41	0.47	0.23	0.29	0.25	0.35	0.20	0.30	0.23	0.33
XGLM	0.30	0.35	0.46	0.41	0.53	0.30	0.29	0.27	0.33	0.20	0.20	0.27	0.33
Avg	0.23	0.35	0.42	0.34	0.45	0.26	0.29	0.26	0.24	0.13	0.18	0.19	0.28

```
1620
                - Plain: m, ^, 3
1621
                - LaTeX: $m, ^, 3, $
1622
                - Braced: $m, ^{, 3, }$
1623
                - Spaced: $m, ^{, , 3, }$
1624
             • mBERT:
1625
1626
                – Plain: m, ^, 3
1627
                - LaTeX: $, m, ^, 3, $
1628
                - Braced: $, m, ^, {, 3, }, $
1629
                - Spaced: $, m, ^, {, 3, }, $ (identical tokenization)
1630
             • Qwen-3:
                – Plain: m, ^, 3
1632
                - LaTeX: $m, ^, 3, $
1633
                - Braced: $m, ^{, 3, }$
1634
1635
                - Spaced: $m, ^{, , 3, }$
1636
             • TokenMonster:
1637
                - Plain: m, ^, 3
                - LaTeX: $, m^, 3$
1639
                - Braced: $, m^, {3}$
1640
                - Spaced: $, m^, {, 3, }$
1641
1642
```

Performance drops precipitously with formatting complexity: while all models correctly identified "volume" for plain text, only 8/14 succeeded with basic LaTeX formatting, 2/14 with braces, and just 2/14 with spaced braces. TokenMonster and Qwen-3 showed the highest robustness, maintaining correct answers through the spaced version.

This shows that even trivial whitespace differences in technical notation can cause catastrophic performance degradation, highlighting a critical vulnerability for applications that require strong mathematical reasoning.

Structural ASCII Art and Chemical Notation These examples demonstrate how tokenizers handle structured chemical representations, from simple formulas to ASCII molecular diagrams and systematic nomenclature. The input contains CH4, an ASCII diagram of methane, H2SO4, and the systematic name "Dihydrogen sulfur tetraoxide":

• BLOOM:

1643

1644

1645

1646

1647

1648

1649 1650

1651

1652

1653

1654 1655

1656

1657

1658

1659

1660 1661

1663

1664

1665

1668

1669

1670

1671

1672

1673

- Simple formulas: CH, 4 and H2, SO4
- ASCII structure: H, , H-C-H, , H (preserves structural elements)
- Systematic name: D, ih, yd, rogen, sulfur, tet, ra, oxide

• XGLM:

- Simple formulas: CH, 4 and H, 2, SO, 4
- ASCII structure: H, , H-, C, -, H, , H (fragments bonds)
- Systematic name: Di, hydro, gen, su, lfur, te, tra, oxide

• mBERT:

- Simple formulas: CH, ##4 and H, ##2, ##S, ##O, ##4
- ASCII structure: H, , H, -, C, -, H, , H (aggressive fragmentation)

• Gemma-2:

- Simple formulas: CH, 4 and H, 2, SO, 4
- ASCII structure: Uses special spacing tokens () for whitespace
- Systematic name: Di, hydrogen, sulfur, tetra, oxide

• **GPT-40**:

- Simple formulas: CH, 4 and H, 2, SO, 4
- ASCII structure: H, , H-C-H, , H (clean structural preservation)
- Systematic name: D, ih, yd, rogen, sulfur, tetra, oxide

• **GPT-2**:

- Simple formulas: CH, 4 and H, 2, SO, 4
- ASCII structure: H, , H-, C, -, H, , H
- Systematic name: D, ih, yd, rogen, sulfur, tet, ra, oxide

• Tekken:

- Simple formulas: CH, 4 and H, 2, SO, 4
- ASCII structure: H, , H-C-H, , H (preserves structure well)
- Systematic name: D, ihydro, gen, sulfur, tetra, oxide

TokenMonster:

- Simple formulas: ch, 4 and h2, so, 4 (lowercase normalization)
- ASCII structure: Complex Unicode handling with encoding artifacts
- Systematic name: di, hydrogen, sul, fur, tet, ra, ox, ide

While all models correctly identified CH4 as methane, only Llama and GPT-2 models correctly interpreted the ASCII molecular diagram. For H2SO4, all models succeeded, while spelled-out systematic nomenclature achieved 65% accuracy. The ASCII diagram failure is particularly revealing—the structured representation that humans easily recognize as methane becomes nearly incomprehensible to models when tokenized, despite containing identical chemical information. XGLM and mBERT normalize the whitespaces in the diagram, however they still fail to identify the molecule, maybe due to—characters. Gemma-2's special whitespace handling (___) and GPT-4o's clean structural preservation suggest different approaches to spatial formatting, yet neither prevented the semantic confusion in the ASCII representation.

E.5 STYLING & UNICODE CHALLENGES

Table 16: Tokenization robustness under Unicode formatting, NFKC normalization used by XGLM strips away all normalizations below. Values represent relative performance drop (canonical-perturbed)/canonical; lower values indicate greater robustness.

Model	Decorative Unicode	Fullwidth Characters	Scripted Text	Double Struck	Enclosed Characters	(Sup/sub) script	Avg
	EN	EN	EN	EN	EN	EN	
XGLM	0.07	0.07	0.02	0.12	0.19	0.08	0.09
ByT5	0.40	0.54	0.58	0.56	0.73	0.66	0.58
GPT-2	0.47	0.59	0.59	0.68	0.61	0.65	0.60
TokenMonster	0.36	0.62	0.57	0.64	0.72	0.70	0.60
Tekken	0.41	0.73	0.57	0.62	0.73	0.62	0.62
Gemma-2	0.53	0.54	0.67	0.62	0.68	0.66	0.62
GPT-4o	0.47	0.62	0.61	0.70	0.67	0.67	0.62
Phi-3	0.47	0.54	0.59	0.75	0.73	0.67	0.62
Aya	0.36	0.68	0.71	0.63	0.69	0.69	0.63
BLOOM	0.59	0.51	0.62	0.67	0.72	0.65	0.63
Qwen-3	0.60	0.67	0.69	0.62	0.57	0.64	0.63
mBERT	0.36	0.73	0.70	0.69	0.81	0.71	0.67
Llama-3.2	0.59	0.60	0.70	0.69	0.76	0.68	0.67
Comma	0.67	0.60	0.67	0.81	0.70	0.58	0.67
Avg	0.45	0.57	0.59	0.63	0.67	0.62	0.59

Using Unicode characters and applying styling to the questions (or all choices) causes performance degradation across all models (see Tables 16 and 17). Although some tokenizers maintain distinct tokens for certain styled characters, they nevertheless exhibit significant failure rates. These styling variations could potentially be mitigated through normalization techniques, such as the NFKC normalization employed by XGLM. However, this is not always desirable as these transformations are irreversible. We include the sample transformations in Fig. 4.

Table 17: Tokenization robustness under different styling formats. Values represent relative performance drop (canonical-perturbed)/canonical; lower values indicate greater robustness.

Model	Diacriticized EN	Lowercase EN	Capitalized EN	Upside Down EN	Spaced EN	Hyphenated EN	Avg
TokenMonster	0.60	0.01	-0.03	0.47	0.66	0.69	0.40
Aya	0.66	0.08	0.15	0.42	0.54	0.67	0.42
GPT-2	0.52	0.06	0.21	0.52	0.63	0.63	0.43
Tekken	0.57	0.03	0.16	0.60	0.63	0.61	0.43
Gemma-2	0.69	0.06	0.15	0.47	0.64	0.67	0.45
GPT-40	0.57	0.00	0.16	0.62	0.62	0.70	0.45
Phi-3	0.58	0.11	0.18	0.47	0.68	0.66	0.45
Comma	0.58	0.06	0.11	0.60	0.68	0.68	0.45
Llama-3.2	0.60	0.11	0.05	0.54	0.68	0.75	0.45
Qwen-3	0.58	0.09	0.11	0.67	0.53	0.76	0.46
ByT5	0.61	0.06	0.06	0.73	0.69	0.67	0.47
BLOOM	0.61	0.08	0.12	0.65	0.72	0.65	0.47
mBERT	0.64	0.09	0.16	0.80	0.59	0.65	0.49
XGLM	0.63	0.11	0.32	0.87	0.61	0.63	0.53
Avg	0.60	0.07	0.14	0.60	0.64	0.67	0.45

		Style	Text
Style	Text	Underline	Python
Full width	Python	Macron	P̄ȳt̄hōñ
Script	Python	Overline	Python
Enclosed/ Circled		Upside down	РАІчои
Enclosed/ Circled		Ring above	P ython
Enclosed/ Parenthesized	(P)(Y)(t)(h)(0)(n)	Diacritics	Pÿthøñ
Superscript	Python	Strikethrough	Python-
Subscript	pγthon	Strikethrough/ Forward slash	Ply/th/o/n/
Double struck	Python	Strikethrough/ Backward slash	Psython:

Figure 4: **Left:** Styling challenges that are normalized by NFKC, **Right:** Styling challenges that NFKC cannot

F EVALUATING INDUSTRY-LEVEL MODELS ON TOKSUITE BENCHMARK

Table 18: Tokenization robustness of original (industry) pre-trained models under multilingual text perturbations. Values represent relative performance drop (canonical-perturbed)/canonical; lower values indicate greater robustness. NEN=non-English.

	Input	Diacr.	Orth.	Errors	Morp	hological	N	oise	LaTeX	STEM	Unicode	Avg
Model	NEN	ZH,FA	EN	IT,TR	EN	TR	EN	NEN	EN	EN	EN	
bert-base-multilingual-cased	0.02	-0.18	0.04	-0.11	0.10	-0.04	-0.15	0.03	0.05	-0.51	-0.12	-0.08
xglm-564M	-0.26	-0.30	0.20	0.05	0.14	0.09	0.13	0.06	0.24	0.16	0.11	0.06
Phi-3-mini-4k-instruct	-0.14	0.13	0.07	-0.25	0.24	-0.26	0.08	-0.02	0.04	0.20	0.59	0.06
GPT-2	-0.30	0.00	0.14	0.00	0.13	0.11	0.18	-0.01	0.23	0.11	0.49	0.10
phi-1_5	-0.13	0.13	0.10	0.00	0.29	-0.17	0.18	-0.04	0.11	0.20	0.62	0.12
Llama-3.2-1B-Instruct	0.14	-0.25	0.04	0.05	0.27	0.13	0.10	0.16	0.04	0.31	0.62	0.15
gemma-2-2b-it	0.21	0.07	0.03	0.24	0.22	0.10	0.04	0.21	0.00	0.20	0.41	0.16
aya-expanse-8b	0.18	0.36	0.03	0.10	0.16	0.07	0.03	0.09	0.11	0.17	0.49	0.16
Qwen3-1.7B-Base	0.25	0.39	0.03	0.12	0.25	0.06	0.06	0.19	-0.02	0.23	0.52	0.19
Llama-3.2-1B	0.13	0.11	0.05	0.26	0.24	0.11	0.08	0.15	0.14	0.33	0.59	0.20
gemma-2-2b	0.30	0.30	-0.01	0.39	0.23	0.13	0.02	0.25	0.16	0.22	0.37	0.21
Avg	0.04	0.07	0.07	0.08	0.21	0.03	0.07	0.10	0.10	0.15	0.43	0.12

While direct comparisons between our models and their original pre-trained counterparts must be interpreted with caution due to fundamental differences in training data, model architectures, and coverage, several noteworthy patterns emerge (see Tables 18 and 1). It should be noted that these models are trained significantly longer than our controlled experiments—for example, Gemma-2-2B (Team et al., 2024) is trained on 2 trillion tokens.

Notably, model size does not appear to be the determining factor, as evidenced by Aya-Expanse-8B (Dang et al., 2024a) performing comparably to smaller models. Instruction-tuned models show marginally better robustness compared to their base counterparts, though the improvement is modest.

Industry models exhibit better overall robustness, with mBERT demonstrating negative degradation values, indicating improved performance on perturbed inputs. This performance gain could stem from training data or training procedure. However, they still struggle significantly with Unicode styling (0.43 average degradation), suggesting that even extensive real-world training data may not adequately cover such specialized character variations. Conversely, our controlled study isolates the effect of tokenization differences by maintaining identical initialization and training data across models, revealing that tokenization choices alone can account for substantial performance variations and more data doesn't always translate into robustness under input variations. The consistent patterns observed across both settings suggest that these robustness challenges are fundamental rather than artifacts of specific training regimes.

LARGE LANGUAGE MODEL USAGE

We used Claude throughout the research process for dataset design brainstorming, generating perturbation ideas, rephrasing sentences, summarizing related work, and assisting with literature review.