

Tokenization Falling Short: hTe Cusre of Tkoeniaztion

Anonymous ACL submission

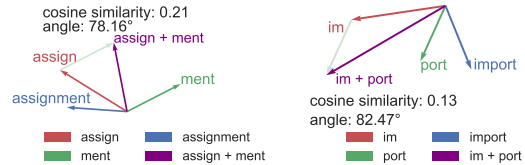
Abstract

Language models typically tokenize raw text into sequences of subword identifiers from a predefined vocabulary, a process inherently sensitive to typographical errors, length variations, and largely oblivious to the internal structure of tokens—issues we term *the curse of tokenization*. In this study, we delve into these drawbacks and demonstrate that large language models (LLMs) remain susceptible to these problems. This study systematically investigates these challenges and their impact on LLMs through three critical research questions: (1) complex problem solving, (2) token structure probing, and (3) resilience to typographical variation. Our findings reveal that scaling model parameters can mitigate the issue of tokenization; however, LLMs still suffer from biases induced by typos and other text format variations. Our experiments show that subword regularization such as BPE-dropout can mitigate this issue. We will release our code and data to facilitate further research.

1 Introduction

Tokenization is a fundamental step in the preprocessing pipeline of large language models (LLMs) (OpenAI, 2023; Anil et al., 2023; Touvron et al., 2023; Chai et al., 2023; Lozhkov et al., 2024), converting raw text into a sequence of subword units derived from a predefined vocabulary (Sennrich et al., 2016; Kudo and Richardson, 2018). This process, while effective in many scenarios, presents significant challenges that can hinder the performance and robustness of LLMs. These challenges include sensitivity to typographical errors (Cao et al., 2023), length variations (Aghajanyan et al., 2022), and a lack of awareness of the internal structure of tokens (Brown et al., 2020)—collectively termed *the curse of tokenization*.

Typographical errors, such as minor misspellings or misplaced characters, can drastically affect the tokenization process. Unlike humans, who can



(a) cosine (“assignment”, “assign” + “ment”). (b) cosine (“import”, “im” + “port”).

Figure 1: Compositional challenges in token embeddings. (a) “assignment” decomposed into “assign” and “ment” shows a cosine similarity of 0.21 and an angle of 78.16°. (b) “import” decomposed into “im” and “port” shows a cosine similarity of 0.13 and an angle of 82.47°. These results indicate that existing LLMs do not accurately capture surface form composition.

easily overlook these errors and understand the intended meaning, LLMs can misinterpret or fail to recognize these variations, leading to degraded performance (Cao et al., 2023). This typo-sensitivity reveals a crucial gap in current tokenization methods, which do not sufficiently mimic human reading capabilities.

Another critical issue is the length unawareness of current tokenization approaches. LLMs often struggle to accurately represent the organizational structure of text, being insensitive to the number of characters or words (Aghajanyan et al., 2022). This insensitivity affects their ability to understand and process text effectively, particularly in tasks requiring a nuanced understanding of text length and compositional structure.

Furthermore, existing tokenization methods are often blind to the internal structure of tokens. The decoupled embedding space and lookup table approach **fail to account for the hierarchical composition of language**, spanning characters, subwords, and words, as depicted in Figure 1. This lack of integration across different levels of token composition limits the model’s ability to fully grasp semantic relationships and differences.

Case insensitivity in certain languages adds another layer of complexity, where variations in cap-

069 italization can lead to different token representa-
070 tions, further complicating the model’s understand-
071 ing and processing of text.

072 To address these challenges, we conducted a
073 comprehensive study examining the limitations of
074 current tokenization methods and their impact on
075 LLM performance. Our study is guided by three
076 critical research questions (RQs):

077 **RQ1: Complex Problem Solving (§3).** As a pilot
078 experiment, we firstly investigate the performance
079 of LLMs on complex problems that are sensitive to
080 tokenization, involving anagram task and complex
081 mathematical language understanding.

082 **RQ2: Token Structure Probing (§4).** We study
083 the token structural tasks such as case manipulation,
084 length counting, and length-sensitive tasks to probe
085 the token structural understanding of LLMs.

086 **RQ3: Typographical Variation (§5).** We de-
087 signed a robust set of evaluation benchmarks
088 termed TKEval on top of various datasets such as
089 MMLU, TruthfulQA, GSM8K, and HumanEval,
090 covering diverse tasks and linguistic phenomena.
091 These benchmarks allow us to systematically test
092 and analyze the LLMs’ resilience to tokenization.

093 Our findings highlight that while scaling model
094 parameters can enhance the robustness of tokeniza-
095 tion, LLMs still suffer from biases introduced by
096 typographical errors and text format variations. We
097 demonstrate the persistent nature of these tokeniza-
098 tion challenges.

099 **Contribution 1)** We provide a comprehensive
100 analysis of the problem known as the *curse of tok-*
101 *enization*, detailing its impact on language model
102 performance and introducing systematic evaluation
103 benchmarks to assess these issues.

104 **2)** By evaluating various scales of LLMs, includ-
105 ing Llama3, Mistral, and GPT-4 families, across
106 thirteen distinct tasks, we demonstrate that even
107 state-of-the-art models struggle with handling ty-
108 pographical variations. Specifically, LLMs exhibit
109 greater sensitivity to character-level variations com-
110 pared to subword-level variations.

111 **3)** We demonstrate that regularized tokenization
112 approaches, such as BPE-dropout with moderate
113 drop rates, can enhance the model’s resilience to
114 the discussed issues.

115 **4)** To facilitate further research, we will release
116 our evaluation code and benchmarks, enabling the
117 research community to build upon our findings and
118 develop more robust models.

2 Related Work 119

2.1 Tokenization 120

Tokenization Approach Conventional language
121 models (Radford et al., 2018; Chai et al., 2020)
122 typically tokenize input text into a sequence of
123 tokens by splitting it into smaller subwords. Tradi-
124 tional tokenization approaches include frequency-
125 based methods such as Byte Pair Encoding (BPE;
126 Sennrich et al., 2016) and probability-based meth-
127 ods like WordPiece (Schuster and Nakajima, 2012).
128 BPE merges tokens based on bigram frequency, re-
129 lying on subword pair co-occurrence to greedily
130 merge neighboring pairs. In contrast, WordPiece
131 can be viewed as a language-modeling based BPE
132 variant. It select the unit pair that maximizes the
133 bigram likelihood of training data at utmost, rather
134 than choose the most frequent pair. 135

136 Unigram Language Model (Kudo, 2018) prunes
137 tokens based on unigram LM perplexity, treating
138 the segmentation process as a probabilistic mix-
139 ture of characters, subwords, and words, reduc-
140 ing subwords by evaluating likelihood reduction.
141 Additionally, some tokenization methods handle
142 text at the byte level (Xue et al., 2022) or char-
143 acter level (Sutskever et al., 2011; Clark et al.,
144 2022). Conventional LLMs often use byte-level
145 BPE (BBPE) for base vocabulary construction, rep-
146 resenting any text with a moderate vocabulary size
147 and avoiding the out-of-vocabulary (OOV) prob-
148 lem. For a detailed introduction to tokenization,
149 readers can refer to Chai (2021).

Tokenization-Free Approach Tokenization ap-
150 proaches often suffer from the *vocabulary bottle-*
151 *neck*, where there is a trade-off between vocabulary
152 size and diverse language coverage in multilingual
153 scenarios. To address this issue, Rust et al. (2023);
154 Chai et al. (2024) introduced a tokenization-free
155 approach that renders raw text as visual text images
156 for masked language modeling and autoregressive
157 pre-training. This method demonstrates robust mul-
158 tilingual generalization capabilities compared to
159 subword tokenization approaches. 160

2.2 Perturbation Probing 161

162 Several studies have investigated the behavior of
163 language models under input perturbations at vari-
164 ous levels, including character-level (Nishino et al.,
165 2019), subword-level (Abdou et al., 2022), and
166 word-level (Sinha et al., 2021) scrambling. Despite
167 these efforts, intrinsic evaluations of perturbing
168 LLM inputs remain under-explored.

Cao et al. (2023) proposed examining scrambled sentence recovery and scrambled QA with context corruptions. In contrast, our study conducts a comprehensive evaluation of both character- and subword-level perturbations, along with noise injection. We evaluate a wide range of LLMs across various tasks to provide a detailed comparison and inspire future research in tokenization and robust model performance.

3 RQ1: Complex Problem Solving

Complex problem-solving tasks are critical benchmarks for evaluating the complex reasoning and comprehension capabilities of LLMs. We explore the LLM’s ability to perform intricate operations on tokenized inputs, as the tokenization process is fundamental to determine how the raw text is segmented and processed, directly impacting the model’s interpretation and prediction.

Anagram solving and mathematical language comprehension were selected to elucidate the relationship between tokenization quality and the model’s performance on complex problem-solving. Anagram tasks require models to decode and rearrange jumbled letters into coherent words, emphasizing the importance of precise token boundaries and recognition accuracy. On the other hand, mathematical language comprehension, particularly expressed with \LaTeX -formatted expressions, demands an exact interpretation of specialized symbols and structured notation, challenging the tokenization process’s robustness.

3.1 Anagram Task

Task Description and Settings The anagram task tests the model’s ability to unscramble a sequence of jumbled characters to form a valid word. This task evaluates the model’s handling of surface-form manipulations and its understanding of char-level compositions within a word. The complexity arises from the need to identify potential word candidates from mixed characters and reassemble them correctly. We present a task example in §A.1. Specifically, we include two tasks:

- **Cycled Letters in Word (CL; Srivastava et al., 2022)** – The model is given a word with its letters cycled, and is expected to generate the original word (e.g., “remo” → “more”).
- **Word Unscrambling (WU) (Srivastava et al., 2022)** – The model is given a randomly scrambled word, and must recover the original word (e.g.,

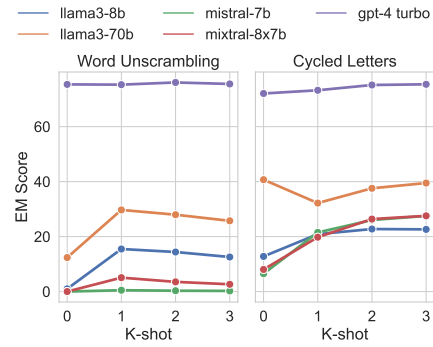


Figure 2: K -shot performance on Word Unscrambling (WU) and Cycled Letters (CL) tasks. The plots illustrate that increasing the number of demonstration examples (K -shot) does not consistently enhance performance. However, models with larger parameter sizes generally exhibit better performance across both tasks.

“nad” → “and”).

We employ exact match (EM) scores for evaluation. Unless otherwise specified, we use the inference-time temperature of 0 for all LLMs in following experiments, to assure the results reproducible.

Results and Analysis The experimental results reveal that larger models demonstrate better performance on the anagram task, yet they remain susceptible to tokenization errors. Specifically, models struggled with longer anagrams (see Figure 3) or those containing uncommon letter combinations, indicating that while scaling improves token recognition, inherent tokenization flaws persist.

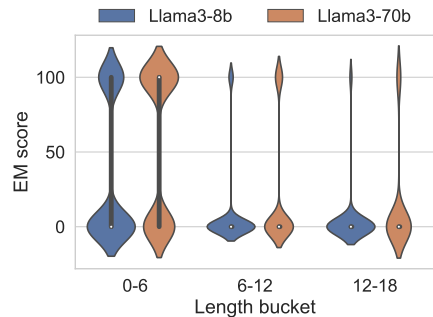


Figure 3: Violin plot illustrating the relationship between the length of scrambled words and the Exact Match (EM) score of Llama3-8B and Llama3-70B on the word unscrambling task under one-shot evaluation. The models tend to correctly reorder anagrams of shorter lengths, while struggling with longer words.

The violin plot in Figure 3 highlights the performance differences between Llama3-8B and Llama3-70B on the word unscrambling task across various word lengths. Notably, Llama3-70B consistently outperforms Llama3-8B, especially in the 0-6 and 6-12 character buckets, where it shows higher and denser EM scores. This trend indicates that as the model parameter size increases from 8B to 70B, the ability to accurately reorder scram-

bled words improves. However, both models struggle with longer words (12-18 characters), though Llama3-70B maintains a slight edge.

Our results indicate several key trends. Firstly, the performance improves significantly as model size scales from 8B to 70B parameters, as evidenced by the performance of Llama3-70B compared to Llama3-8B (AI@Meta, 2024). Secondly, we observe that while the dense Mistral-7B model performs poorly, the sparse Mixtral-8x7B model (an MoE sparse model) shows improved performance due to its parameter size scaling. Lastly, GPT-4 turbo, a much more powerful model, achieves state-of-the-art results, clearly outperforming all other models across all shot conditions. This sensitivity underscores the need for more robust tokenization that can handle typographical variations without degrading performance.

3.2 Mathematical Language (L^AT_EX) Comprehension

Task Description and Settings The mathematical language comprehension task evaluates the model’s ability to read and comprehend mathematics written in L^AT_EX, the typesetting language used by professional mathematicians. This task assesses the models’ capability to interpret complex mathematical expressions and accurately tokenize symbols and operators within structured L^AT_EX format.

We employ **Identify Math Theorems (IMT; Srivastava et al., 2022)** for evaluation, and use perplexity to measure the model’s confidence in different given choices. The dataset comprises 54 problems divided into nine sections, each representing a major area of mathematics research or advanced pedagogy. We present the input-output format in §A.1.

Setting	0-Shot	1-Shot	2-Shot	3-Shot
GPT-3 (6B) ^a	33.96	28.30	33.96	28.30
GPT-3 (200B) ^a	32.08	30.19	33.96	30.19
Llama2-7b	37.70	34.00	35.80	37.70
Llama3-8b	41.51	45.28	45.28	35.85
Llama3-70b	62.26	79.25	69.81	71.70
Mistral-7b	47.20	43.40	37.70	37.70
Mixtral-8x7b	49.10	56.60	64.20	62.30

Table 1: Few-shot results on **Identifying Math Theorems**, with exact match (EM) scores reported as percentages. ^a[Srivastava et al., 2022]

Results and Analysis Our evaluation results of various models are shown in Table 1. The results demonstrate that while larger models generally per-

form better on LaTeX-formatted mathematical content, the relationship between the number of in-context examples and model performance is not linear. The Llama3-70B model consistently outperformed other models, achieving a score of 62.26% in the zero-shot setting and improving to 79.25% with one-shot learning. However, additional in-context examples led to fluctuating performance, with scores of 69.81% in the two-shot setting and 71.70% in the three-shot setting. This indicates that increasing the number of in-context demonstrations does not consistently enhance performance and may lead to variability depending on the specific examples provided.

Other models exhibited similar trends. GPT-3-200B (Srivastava et al., 2022), despite its larger parameter count, did not show significant improvement over the smaller GPT-3-6B model, suggesting that simply increasing the model size does not guarantee better performance in L^AT_EX comprehension tasks. The comparison between dense and sparse models revealed that the dense Mistral-7B model performed poorly across all shot conditions, while the sparse Mixtral-8x7B model demonstrated better results, especially in few-shot scenarios.

4 RQ2: Token Structure Probe

Tokenization is a key preprocessing step in LLMs, yet it introduces several significant challenges, which we defined as the curse of tokenization. These challenges include length unawareness, case insensitivity, and a lack of awareness of the internal structure of tokens. Tokenization transforms text into sequences of token identifiers, often obscuring the surface form and internal structure of the original text. This conversion can lead to deficiencies in the model’s ability to understand and process textual data accurately.

The curse of tokenization manifests in several ways, which refers to the inherent challenges:

- A) **Length Unawareness:** Models struggle to recognize the organizational structure of text, such as the number of characters or words.
- B) **Case Insensitivity:** Variations in capitalization can lead to different token identifiers and representations, complicating the model’s processing of text.
- C) **Blindness to Internal Structure:** The decoupled embedding space and lookup table approach used in LLMs fail to preserve the hierarchical and relational structure within tokens,

obscuring the surface form and internal relationships between characters and subwords.

To address these challenges, we construct a set of probing tasks to evaluate the model’s understanding of token structure. These tasks are divided into intra-token (§4.1) and inter-token probes (§4.2).

4.1 Intra-Token Probing

Task Description and Settings To measure the capability of LLMs, we devise intra-token probing tasks related to length, case, and counting problems. These tasks evaluate the model’s performance on the internal structure of tokens or word, specifically including four tasks:

- **Character Count (CC)** – The model is asked to count the number of occurrences of a specific character within the given word (*e.g.*, the character appears twice in the word “undertake” → the answer is: “e”).
- ***N*-th Character (NC)** – The model is expected to output the *n*-th character of the given word (*e.g.*, 4-th character of the word “dual” → “l”).
- ***N*-th Character Reverse (NCR)** – The model must identify the *n*-th character from the end of a word (*e.g.*, 2nd character from the end of the word “dual” → “a”).
- **Case Conversion (CCV)** – This task involves converting the characters within a word to different cases (uppercase, lowercase) or converting the word into title case.

For each task, we conduct many-shot evaluation (0-3 shot) and report the EM score to test the model’s ability to understand and manipulate the internal structure of tokens and words at a granular level, revealing the extent to which the tokenization process could preserve this information. Detailed test examples are provided in Appendix A.2.

Results and Analysis We evaluated CC, NC, NCR, and CCV tasks across different models and shot settings. The results are presented in Figure 4.

The CC task reveals that larger models exhibit superior performance, particularly GPT-4 turbo, which achieves near-perfect accuracy across all shot conditions. Smaller models, such as Llama3-8B, show significant improvement with few-shot learning, indicating that exposure to examples greatly enhances their performance. For example, Llama3-8B’s accuracy improves from 0% in the zero-shot setting to 81% in the three-shot setting. This demonstrates that increased model size and few-shot learning contribute positively to CC tasks.

The NC task underscores the difficulty models

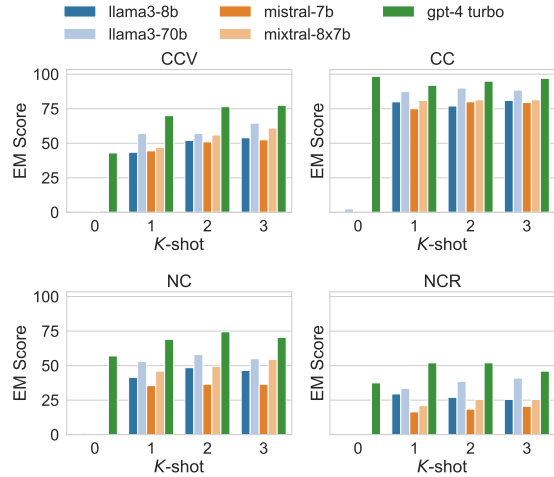


Figure 4: *K*-shot performance on intra-token probing tasks (CCV, CC, NC, NCR). The plots demonstrate that increasing the number of demonstration examples (*K*-shot) generally results in an improvement from zero-shot to one-shot, with performance stabilizing thereafter.

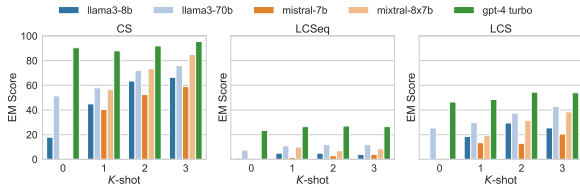
face in accurately identifying specific characters within words. GPT-4 turbo again leads in performance, while smaller models show substantial improvement with increased shots. Llama3-70B, for instance, improves from 1% in the zero-shot setting to 55% in the three-shot setting. This indicates that while larger models perform better, few-shot learning plays a crucial role in enhancing the model’s ability to identify specific characters.

Identifying characters from the end of the word, or reverse character identification, proves more challenging. GPT-4 turbo achieves the highest performance with a score of 52% in the one-shot setting, though overall accuracy is lower compared to other tasks. Smaller models like Llama3-8B show moderate improvements with additional shots, but their performance remains relatively low. This highlights the complexity of reverse character identification and the need for more advanced tokenization strategies to address this challenge.

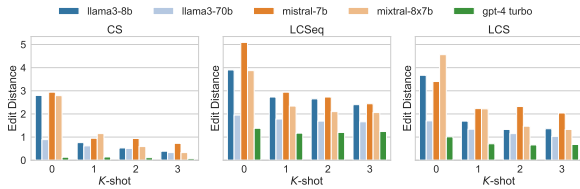
4.2 Inter-Token Probing

Task Description and Settings To evaluate the capabilities of LLMs in understanding and manipulating relationships between multiple tokens, we devised inter-token probing tasks. These tasks focus on identifying common patterns and sequences across tokens, and they assess the model’s ability to recognize and process such relationships. Specifically, we include three tasks:

- **Common Substrings (CS)** – The model identifies multiple common substrings between two given words.
- **Longest Common Substrings (LCS)** – The



(a) EM score performance on k -shot evaluation.



(b) Edit distance on k -shot evaluation.

Figure 5: K -shot performance on various inter-token probing tasks. For edit distances, lower is better.

model identifies the longest continuous common substring for two given words.

- **Longest Common Subsequences (LCSeq)** – The model identifies the longest common subsequence (not necessarily continuous) between two given words.

For each task, we conduct many-shot evaluations (0-3 shot) and report the EM and edit distance (ED) score to test the model’s ability to understand and manipulate relationships between tokens at a higher level. For CS tasks, the model’s response is considered correct if it generates one of the multiple possible common substrings. Detailed test examples are provided in Appendix A.3.

Results and Analysis The results are presented in Figure 5. The results for CS tasks indicate that larger models, such as Llama3-70B and GPT-4 Turbo, perform significantly better than smaller models. GPT-4 Turbo achieves the highest accuracy across all shot settings, demonstrating the model’s robustness in identifying continuous substrings. Notably, Llama3-70B also shows strong performance, particularly in the three-shot setting. Sparse models like Mixtral-8x7B exhibit notable improvements compared to dense models, highlighting the effectiveness of sparse architectures in handling complex token relationships.

For LCS tasks, GPT-4 Turbo leads in performance, achieving high accuracy across all shot settings. Llama3-70B and Mixtral-8x7B show considerable improvements with increased shots, indicating that exposure to more examples helps models better identify multiple common substrings. Dense models like Mistral-7B lag behind, reinforcing the advantage of sparse architectures in such tasks.

The LCSeq task reveals that even the best-performing models face challenges with non-continuous patterns. While Llama3-70B and GPT-4 Turbo demonstrate superior performance, there is a significant drop in accuracy compared to CS tasks. Few-shot learning significantly enhances the performance of smaller models, such as Llama3-8B, which improves from 1% in zero-shot to 4% in three-shot settings. This underscores the importance of few-shot examples in aiding models to recognize and process non-continuous patterns.

5 RQ3: Typographical Variation

To evaluate the robustness of LLMs to typographical variations, we constructed tasks that introduce character-level and token-level typographical errors into the input text. These tasks are designed to test the models’ ability to maintain semantic understanding despite the presence of such errors. The datasets include MMLU (Hendrycks et al., 2021), TruthfulQA (Lin et al., 2022), GSM8K (Cobbe et al., 2021), and HumanEval (Chen et al., 2021), ensuring diverse coverage.

Task Description and Settings The primary goal of these tasks is to assess the LLMs’ resilience to typographical errors at both the character and token levels, examining whether these models can maintain semantic understanding when faced with such perturbations. For **character-level typographical variation**, we employed n -gram shuffling within word boundaries (with n set to 2, 3, 5) with a 50% probability, and n -gram noise, which involve adding, deleting, and replacing characters, spaces, and punctuation marks to simulate spelling noise. This corruption occurs with a 30% probability, including insertion, deletion, or substitution operations. **Token-level typographical variation** was introduced by shuffling tokens within n -grams of sizes 2, 3, and 5, with a 50% probability of permutation or typo generation, similar to the character-level method. We include four tasks:

- **Character-Level Permutation:** Shuffling characters within word boundaries using n -grams of sizes 2, 3, and 5, with a 50% probability.
- **Character-Level Noise:** Adding, deleting, replacing random characters to simulate spelling noise, each with a 10% probability.
- **Token-Level Permutation:** Randomly reordering tokens using n -grams of sizes 2, 3, 5, with a 50% probability.
- **Token-Level Noise:** Adding, deleting, replacing

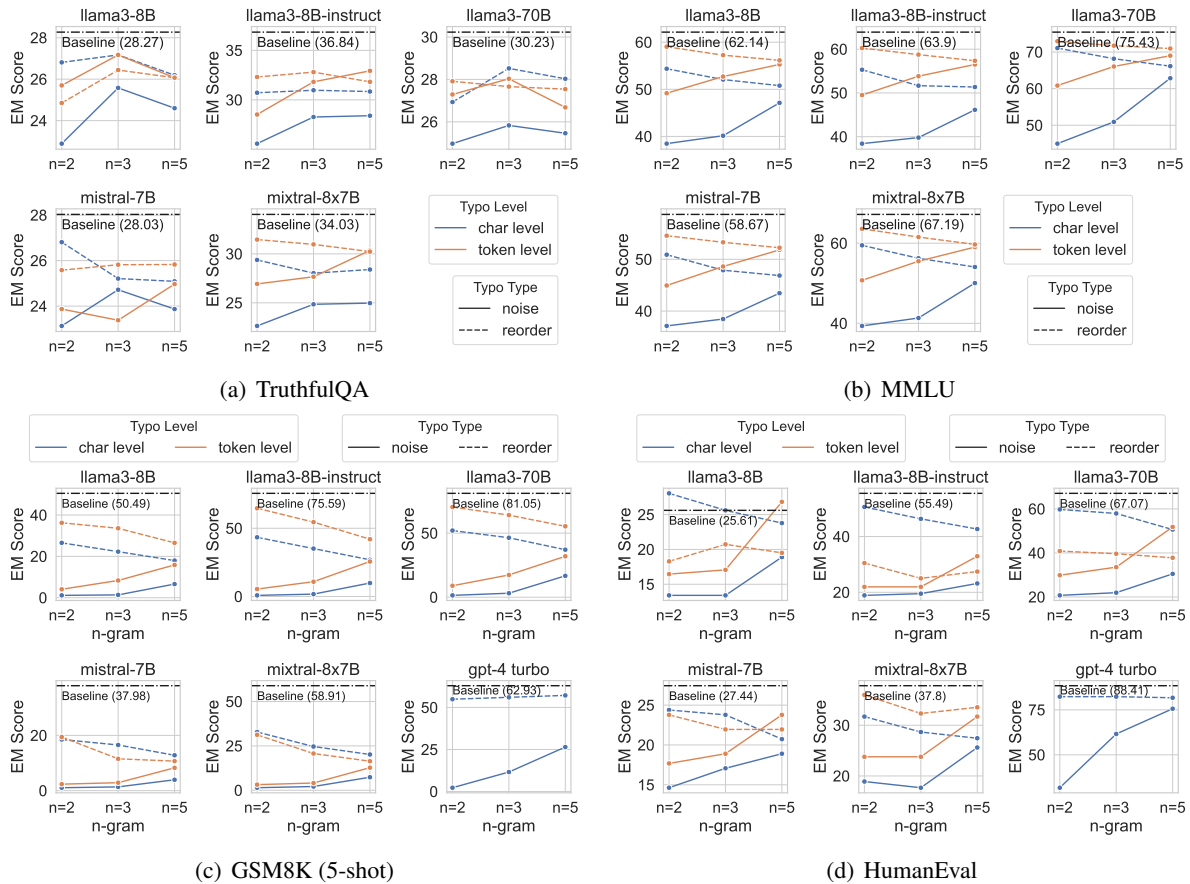


Figure 6: Performance comparison for various models across different n -gram sizes ($n=2,3,5$) and typographical variations on (1) **TruthfulQA**, (b) **MMLU**, (c) **GSM8K**, and (d) **HumanEval**. The typographical variations include character-level (blue) and token-level (orange) perturbations, with noise (solid line) and reorder (dashed line) types. Baseline performance is indicated with a dotted line at the top of each plot.

tokens to inject noise, with a 30% probability.

For evaluation, we report pass@1 for HumanEval using a temperature of 0.2 and a top- p of 0.95. For others, we used a temperature of 0. GSM8K was evaluated using a 5-shot setting, while MMLU, TruthfulQA, and HumanEval were assessed in a zero-shot setting. We measured performance on MMLU and TruthfulQA using perplexity for multiple-choice selection¹.

Results and Analysis Figures 6 presents a comprehensive analysis of the impact of typographical variations on various LLMs, specifically focusing on character-level and token-level perturbations across different n -gram sizes ($n=2, 3, 5$). The evaluation covers a range of datasets, including TruthfulQA, MMLU, GSM8K (5-shot), and HumanEval.

Across all datasets and models, there is a consistent trend showing that LLMs are much more sensitive to noise (solid lines) than to reordering (dashed lines). Noise injection, which involves adding, deleting, or replacing characters or tokens,

¹Since GPT-4 Turbo does not support perplexity computation, it was excluded from the evaluation for these two tasks.

leads to more pronounced variations and generally degrades overall performance. This is evident from the lower EM scores for noise perturbations compared to reorder perturbations.

Despite the challenges posed by n -gram reordering within word boundaries, GPT-4 Turbo maintained high accuracy across all n -gram sizes. Character-level n -gram noise injection, simulating realistic spelling noise, further tested the models' robustness. The results indicate that all models experienced evident performance degradation, regardless of the parameter sizes, highlighting their sensitivity to typographical noise.

For most models and datasets, as the n -gram size of noise injection increases from $n=2$ to 5, the performance tends to stabilize or improve. This trend suggests that models can better handle larger n -gram noise injection, likely because the context within larger n -grams provides more semantic coherence compared to smaller n -grams.

At the token level, models were subjected to n -gram permutations similar to those applied at the character level. The results indicated that models

generally performed better with token-level permutations than with character-level shuffles and noise injection. This suggests that token-level errors may be less disruptive to the overall semantic structure of the input text.

6 Does BPE-dropout Matter?

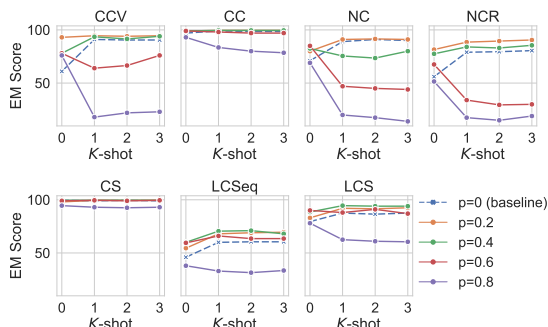


Figure 7: K -shot performance on various tasks (CCV, CC, NC, NCR, CS, LCSeq, and LCS) using the Mistral-7B model fine-tuned with a BPE-dropout tokenizer at different dropout rates p , ranging from 0 to 0.8. The baseline without BPE-dropout (*i.e.*, $p = 0$) is depicted with a dashed line. It demonstrates that introducing a moderate amount of variability during tokenization improves the model’s generalization capabilities, mitigating *the curse of tokenization* issues.

To further enhance the robustness of our models, we explore regularized tokenization approach, BPE-dropout (Provilkov et al., 2020), which randomly drops BPE merges during tokenization. This technique allows text sequences to be tokenized in more diverse ways, promoting robustness to various token combinations and increasing the likelihood of encountering smaller tokens. Intuitively, this diversity benefits the model’s understanding of internal token structures.

Training Setup For the training data, we synthesized a dataset consisting of 111k examples specifically designed for RQ2. More detailed training settings are provided in Appendix B.3.

Results and Analysis Figure 7 shows the impact of varying BPE-dropout rates on the Mistral-7B model’s performance across multiple K -shot settings and tasks. The baseline performance, with a dropout rate of $p = 0$, shows robust results across several tasks, particularly in CS and CC. These tasks are relatively straightforward, involving simple token manipulations that do not significantly challenge the model’s capacity to generalize from zero-shot to few-shot scenarios. The high performance in these tasks suggests that the Mistral-7B model, even without BPE-dropout, is adept at handling simpler token relationships. However, it is

important to note that the baseline performance does not uniformly extend to more complex tasks.

In contrast, tasks such as LCSeq reveal relatively low performance across all models, irrespective of the BPE-dropout rate. This suggests inherent difficulties in these tasks that stem from the requirement to identify non-continuous and intricate token patterns. The consistent under-performance indicates that LCSeq tasks pose a significant challenge to the model’s ability to generalize, likely due to the increased complexity in recognizing and processing longer and fragmented sequences.

Interestingly, the introduction of a moderate BPE-dropout rate ($p = 0.2$) frequently surpasses the baseline, highlighting the benefits of inducing variability during tokenization. This moderate dropout rate enhances the model’s generalization capabilities by preventing overfitting and promoting a more robust learning process. Notably, in tasks such as CCV, NC, and LCS, the $p = 0.2$ model consistently achieves higher EM scores, underscoring the benefits of incorporating tokenization regularization.

Our analysis reveals that higher dropout rates ($p = 0.6$ and $p = 0.8$) exhibit relatively lower performance across most tasks. This decline can be attributed to insufficient training, as the dataset was trained for only five epochs. The higher dropout rates introduce greater tokenization variation, which necessitates additional training compute to achieve convergence. The lack of adequate training epochs likely hindered these models from fully leveraging the potential benefits of higher BPE-dropout rates. We provided detailed training analysis in Appendix §D.

7 Conclusion

In this study, we investigated the challenge of *the curse of tokenization*, comprehensively evaluating mainstream LLMs across thirteen tasks sensitive to conventional subword tokenization. Our findings reveal that while larger models and increased shot counts can partially mitigate these issues, LLMs still struggle with understanding internal structures and token compositions. Moderate BPE-dropout can alleviate some of these challenges, whereas larger drop rates lead to performance degradation. We encourage the research community to develop more flexible approaches to further address these limitations and enhance model robustness.

624	Limitations		
625	While our study provides valuable insights into		
626	the robustness and performance of large language		
627	models (LLMs) under various tokenization and ty-		
628	pographical variation scenarios, several limitations		
629	should be acknowledged:		
630	Data Diversity and Size The training data syn-		
631	thesized for RQ2 (token structure probe) consists		
632	of approximate 30k examples. While this dataset		
633	size is substantial, it may not fully capture the di-		
634	versity and complexity of real-world text. Future		
635	work could benefit from expanding the dataset size		
636	and incorporating a wider range of linguistic phe-		
637	nomena.		
638	Evaluation Metrics Our evaluation primarily re-		
639	lies on metrics such as pass@1 for HumanEval and		
640	accuracy for MMLU, TruthfulQA, and GSM8K.		
641	While these metrics provide valuable insights, they		
642	may not fully capture the nuanced performance		
643	of LLMs in real-world applications. Incorporat-		
644	ing additional metrics that assess other aspects of		
645	model performance, such as robustness to out-of-		
646	distribution data and interpretability, could provide		
647	a more comprehensive evaluation.		
648	Subword Regularization Although our use of		
649	BPE-dropout shows promising improvements in		
650	model robustness and accuracy, the approach in-		
651	troduces randomness into the tokenization process.		
652	This randomness can lead to variability in model		
653	performance, making it challenging to ensure con-		
654	sistent improvements across different datasets and		
655	tasks. Further research is needed to optimize the		
656	BPE-dropout technique and evaluate its long-term		
657	impact on model performance.		
658	Typographical Variation Our study focuses on		
659	character-level and token-level typographical varia-		
660	tions, but it does not address other common types		
661	of text perturbations, such as grammatical errors,		
662	semantic variations, or contextual inconsistencies.		
663	Exploring the effects of these additional types of		
664	variations could provide a more holistic understand-		
665	ing of LLM robustness.		
666	Generalizability The findings from our evalu-		
667	ation on specific datasets (MMLU, TruthfulQA,		
668	GSM8K, and HumanEval) may not generalize to		
669	all types of text and tasks. Further studies are		
670	needed to assess the generalizability of our findings		
671	across a broader range of datasets and real-world		
	scenarios, such as extending to multilingual evalu-	672	
	ation (Peng et al., 2024).	673	
	Ethical Consideration	674	
	Bias and Fairness Tokenization strategies can	675	
	introduce or exacerbate biases present in the train-	676	
	ing data. Our study, which involves diverse to-	677	
	kenization techniques like BPE-Dropout, should	678	
	include thorough bias assessments to ensure that	679	
	these methods do not perpetuate unfair or discrim-	680	
	inatory outcomes. Mitigating bias is crucial for	681	
	creating fair and equitable AI systems.	682	
	Transparency and Interpretability Our re-	683	
	search involves complex tokenization processes	684	
	that can obscure the decision-making of LLMs.	685	
	Enhancing the transparency of these models by	686	
	providing clear explanations of how tokenization	687	
	impacts model behavior is essential. This trans-	688	
	parency helps build trust and allows users to un-	689	
	derstand and identify potential issues in language	690	
	model predictions.	691	
	Privacy and Data Security The datasets used	692	
	for training and evaluating tokenization methods	693	
	often contain sensitive information. Ensuring data	694	
	anonymization and compliance with data protec-	695	
	tion regulations is critical to protecting user privacy.	696	
	Our study adheres to strict data security protocols	697	
	to prevent any misuse of sensitive information.	698	
	References	699	
	Mostafa Abdou, Vinit Ravishankar, Artur Kulmizev, and	700	
	Anders Søgaard. 2022. Word order does matter and	701	
	shuffled language models know it . In <i>Proceedings</i>	702	
	<i>of the 60th Annual Meeting of the Association for</i>	703	
	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	704	
	pages 6907–6919, Dublin, Ireland. Association for	705	
	Computational Linguistics.	706	
	Armen Aghajanyan, Dmytro Okhonko, Mike Lewis,	707	
	Mandar Joshi, Hu Xu, Gargi Ghosh, and Luke Zettle-	708	
	moyer. 2022. HTLM: hyper-text pre-training and	709	
	prompting of language models . In <i>The Tenth Inter-</i>	710	
	<i>national Conference on Learning Representations,</i>	711	
	<i>ICLR 2022, Virtual Event, April 25-29, 2022</i> . Open-	712	
	Review.net.	713	
	AI@Meta. 2024. Llama 3 model card .	714	
	Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-	715	
	Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan	716	
	Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Mil-	717	
	lican, David Silver, Slav Petrov, Melvin Johnson,	718	
	Ioannis Antonoglou, Julian Schrittwieser, Amelia	719	
	Glaese, Jilin Chen, Emily Pitler, Timothy P. Lilli-	720	
	crap, Angeliki Lazaridou, Orhan Firat, James Molloy,	721	

722	Michael Isard, Paul Ronald Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs, Anaïs White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakub Sygnowski, and et al. 2023. Gemini: A family of highly capable multimodal models . <i>CoRR</i> , abs/2312.11805.	779
723		780
724		781
725		782
726		783
727		784
728		785
729		786
730		787
731		788
732		789
733		790
734	Loubna Ben Allal, Niklas Muennighoff, Logesh Kumar Umapathi, Ben Lipkin, and Leandro von Werra. 2022. A framework for the evaluation of code generation models. https://github.com/bigcode-project/bigcode-evaluation-harness .	791
735		792
736		793
737		794
738		795
739		796
740	Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. <i>Advances in neural information processing systems</i> , 33:1877–1901.	797
741		798
742		799
743		800
744		801
745		802
746	Qi Cao, Takeshi Kojima, Yutaka Matsuo, and Yusuke Iwasawa. 2023. Unnatural error correction: GPT-4 can almost perfectly handle unnatural scrambled text . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 8898–8913. Association for Computational Linguistics.	803
747		804
748		805
749		806
750		807
751		808
752		809
753	Yekun Chai. 2021. Word Tokenization for Pre-trained Models. https://cyk1337.github.io/notes/2021/11/29/Subword-Tokenization-in-NLP/ .	810
754		811
755		812
756	Yekun Chai, Shuo Jin, and Xinwen Hou. 2020. Highway transformer: Self-gating enhanced self-attentive networks . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 6887–6900, Online. Association for Computational Linguistics.	813
757		814
758		815
759		816
760		817
761		818
762	Yekun Chai, Qingyi Liu, Jingwu Xiao, Shuohuan Wang, Yu Sun, and Hua Wu. 2024. Dual modalities of text: Visual and textual generative pre-training .	819
763		820
764		821
765	Yekun Chai, Shuohuan Wang, Chao Pang, Yu Sun, Hao Tian, and Hua Wu. 2023. ERNIE-code: Beyond English-centric cross-lingual pretraining for programming languages . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 10628–10650, Toronto, Canada. Association for Computational Linguistics.	822
766		823
767		824
768		825
769		826
770		827
771		828
772	Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared Kaplan, Harrison Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz	829
773		830
774		831
775		832
776		833
777		834
778		835
		836
		837
	Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating large language models trained on code . <i>CoRR</i> , abs/2107.03374.	799
		800
		801
		802
		803
	Jonathan H. Clark, Dan Garrette, Iulia Turc, and John Wieting. 2022. Canine: Pre-training an efficient tokenization-free encoder for language representation . <i>Transactions of the Association for Computational Linguistics</i> , 10:73–91.	804
		805
		806
		807
		808
		809
		810
		811
		812
	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems . <i>CoRR</i> , abs/2110.14168.	813
		814
		815
		816
		817
		818
	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. 2023. A framework for few-shot language model evaluation .	819
		820
		821
		822
		823
		824
		825
		826
	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.	827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b . <i>CoRR</i> , abs/2310.06825.	838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900

952 et al. 2022. [Beyond the imitation game: Quantifying](#)
953 [and extrapolating the capabilities of language models.](#)
954 *CoRR*, abs/2206.04615.

955 Ilya Sutskever, James Martens, and Geoffrey Hinton.
956 2011. Generating text with recurrent neural networks.
957 In *Proceedings of the 28th International Conference*
958 *on International Conference on Machine Learning*,
959 ICML'11, page 1017–1024, Madison, WI, USA. Om-
960 nipress.

961 Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-
962 bert, Amjad Almahairi, Yasmine Babaei, Nikolay
963 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti
964 Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-
965 Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,
966 Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,
967 Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-
968 thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan
969 Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,
970 Isabel Kloumann, Artem Korenev, Punit Singh Koura,
971 Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-
972 ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-
973 tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-
974 bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-
975 stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,
976 Ruan Silva, Eric Michael Smith, Ranjan Subrama-
977 nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-
978 lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,
979 Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,
980 Melanie Kambadur, Sharan Narang, Aurélien Ro-
981 driguez, Robert Stojnic, Sergey Edunov, and Thomas
982 Scialom. 2023. [Llama 2: Open foundation and fine-](#)
983 [tuned chat models.](#) *CoRR*, abs/2307.09288.

984 Linting Xue, Aditya Barua, Noah Constant, Rami Al-
985 Rfou, Sharan Narang, Mihir Kale, Adam Roberts,
986 and Colin Raffel. 2022. [Byt5: Towards a token-free](#)
987 [future with pre-trained byte-to-byte models.](#) *Trans.*
988 *Assoc. Comput. Linguistics*, 10:291–306.

Contents		989
1 Introduction	1	990
2 Related Work	2	991
2.1 Tokenization	2	992
2.2 Perturbation Probing	2	993
3 RQ1: Complex Problem Solving	3	994
3.1 Anagram Task	3	995
3.2 Mathematical Language (L ^A T _E X) Comprehension	4	996
4 RQ2: Token Structure Probe	4	997
4.1 Intra-Token Probing	5	998
4.2 Inter-Token Probing	5	999
5 RQ3: Typographical Variation	6	1000
6 Does BPE-dropout Matter?	8	1001
7 Conclusion	8	1002
A Task Examples	14	1003
A.1 Complex Problem Solving Examples	14	1004
A.2 Intra-Token Probing Examples	14	1005
A.3 Inter-Token Probing Examples	14	1006
B Experimental Settings	15	1007
B.1 Baselines	15	1008
B.2 Evaluation Settings	15	1009
B.3 Post-Training Details	15	1010
C Details of Probing Task Construction	15	1011
C.1 Token Structure Probing (RQ2)	15	1012
C.2 Typographical Variation Task (RQ3)	16	1013
D Detailed Results of BPE-dropout Post-Training	16	1014
E Additional Analysis	18	1015
E.1 Impact of Typographical Variations on Sequence Length	18	1016
E.2 Compositional Challenges in Token Embeddings	19	1017

1018

A Task Examples

1019

A.1 Complex Problem Solving Examples

1020

We present detailed examples of complex problem-solving tasks including word anagram and identifying math theorems as follows:

1021

1022

Anagram Task Format

Input: A string of jumbled characters (e.g., “moeh” for “home”).

Output: The correct unscrambled word (e.g., “home”).

1023

Identifying Math Theorems Task Format

Input: A \LaTeX -formatted mathematical theorem. E.g., “Let $f \in L^1(\mathbb{R})$ be an integrable function. The span of $\{f_a(x) = f(x+a) : a \in \mathbb{R}\}$ is dense in $L^1(\mathbb{R})$ if and only if \widehat{f} has no real roots..

A) Let $f \in L^1(\mathbb{R})$ be an integrable function. The span of $\{f_a(x) = f(x+a) : a \in \mathbb{R}\}$ is dense in $L^1(\mathbb{R})$ if and only if \widehat{f} has no real roots.

B) Let $f \in L^1(\mathbb{R})$ be an integrable function. The span of $\{f_a(x) = f(x+a) : a \in \mathbb{R}\}$ is dense in $L^1(\mathbb{R})$ if and only if \widehat{f} has no real roots. .

C) Let $f \in L^1(\mathbb{R})$ be an integrable function. The span of $\{f_a(x) = f(x+a) : a \in \mathbb{R}\}$ is dense in $L^1(\mathbb{R})$ if and only if \widehat{f} is irreducible over \mathbb{Q} .

D) Let $f \in L^1(\mathbb{R})$ be an integrable function. The span of $\{f_a(x) = f(x+a) : a \in \mathbb{R}\}$ is dense in $L^1(\mathbb{R})$ if and only if \widehat{f} has no repeated roots.

Output: The model must determine whether the theorem is true. If it is false, the model should provide the correct version; i.e., select the option “A”.

1024

1025

A.2 Intra-Token Probing Examples

1026

1027

1028

1029

For intra-token probing tasks, we provide Character Count (CC), N -th Character (NC), N -th Character Reverse (NCR), and Case Conversion (CCV) for illustration.

Character Count (CC)

Input: Which character appears 3 times in the word ‘messrs’?

Output: ‘s’.

1030

N -th Character (NC)

Input: What is the 4th character of the word ‘myron’?

Output: ‘o’.

1031

N -th Character Reverse (NCR)

Input: What is the 2nd character from the end of the word ‘pensioner’?

Output: ‘e’.

1032

Case Conversion (CCV)

Input: Which character appears 3 times in the word ‘messrs’?

Output: ‘s’.

1033

A.3 Inter-Token Probing Examples

We present examples of inter-token probing tasks, which involve identifying Common Substrings (CS), Longest Common Subsequences (LCSeq), and Longest Common Substrings (LCS). These tasks evaluate the model’s ability to analyze and compare internal structure across different inputs.

1034

1035

1036

1037

1038

1039

1040

Common Substrings (CS)

Input: What are the common substrings of ‘critical’ and ‘conscious’?

Output: ‘i’, ‘c’.

1041

Longest Common Subsequences (LCSeq)

Input: What are the longest common subsequences of ‘illustrate’ and ‘critical’?

Output: ‘ita’.

1042

Longest Common Substrings (LCS)

Input: What are the longest common substrings of ‘cow’ and ‘condition’?

Output: ‘co’.

1043

B Experimental Settings

B.1 Baselines

We include Llama3-8B, Llama3-8B-Instruct, Llama3-70B, Mistral-7B and Mixtral-8x7B for LLM evaluation.

Llama3 Llama3 (AI@Meta, 2024) series are one of the most powerful open-sourced models recently. Llama3-8B is a dense pretrained model with a vocab size of 128256, which needs few-shot examples to better follow instructions. Llama3-8B-Instruct is also evolved for diverse model types. Llama3-8B-Instrcut is a instruction-fine-tuned version of Llama3-8B, showing much improvement over Llama3-8B on benchmarks like HumanEval and TruthfulQA.

Mistral & Mixtral Mistral-7B (Jiang et al., 2023) is a dense model with a vocab size of 32000 released last year. Mixtral-8x7B (Jiang et al., 2024) is a sparse mixture-of-expert(MoE) model with 13B active parameters, whose performance greatly surpasses Mistral-7B and matches the performance of Llama2-70B.

GPT-4 Turbo The model version we evaluate is "gpt-4-1106-preview"². Compared with GPT-4, "gpt-4-1106-preview" yields stronger performance on following instructions, structured output and other abilities.

B.2 Evaluation Settings

We utilize lm-evaluation-harness (Gao et al., 2023) for the evaluation of GSM8K, MMLU, and TruthfulQA. For HumanEval, we adopt bigcode-evaluation-harness (Ben Allal et al., 2022). All models are tested under bfloat16 precision for higher efficiency.

To eliminate the impact of the Chain-of-Thought (CoT) prompt, GSM8K is evaluated using a 5-shot setting without CoT, and we report the "exact match" as the final metric. For HumanEval, we report pass@1 using a temperature of 0.2 and a top- p of 0.95. The maximum total length of the prompt and model output is set to 512 for Llama3 and 1024 for Mistral-7B and Mixtral-8x7B. We only apply corruption to the annotation to make sure that entry point can be found after corruption. For TruthfulQA, model performances are measured within the "MC1 (single-true)" setting.

²<https://platform.openai.com/docs/models/gpt-4-turbo-and-gpt-4>

For models after instruction-tuning, we do not apply chat templates except for TruthfulQA, as model outputs tend to be difficult to parse in chat mode.

B.3 Post-Training Details

Dataset We synthesized the dataset for RQ2 with a template-based method. Table 2 describes statistics of dataset splits for each task.

Task	Train	Test
Intra-Token Probing	111,070	800
Character Count	20,775	200
N-th Character	31,241	200
N-th Character Reverse	31,316	200
Case Conversion	27,738	200
Inter-Token Probing	14,400	600
Common Substrings	4,800	200
Longest Common Substrings	4,800	200
Longest Common Subsequences	4,800	200

Table 2: The dataset size for training and testing.

Training Details We employ the AdamW optimizer (Loshchilov and Hutter, 2017) with the hyperparameters of $\beta_1 = 0.9$, $\beta_2 = 0.95$. The peak learning rate is set to $5e-5$, and the minimum learning rate is set to $1e-6$. The learning rate warms up during the first 10% of training steps and then decays with a cosine scheduler.

Given the difference in data distribution resulting from BPE-dropout, we post-train Mistral-7B (Jiang et al., 2023) on the training split for 5 epochs with a global batch size of 16. Following the pre-training recipe, we concatenate all sequences and then chunk them into fixed context lengths of 4096 for our autoregressive post-training. We conduct data shuffling within the same epochs.

C Details of Probing Task Construction

C.1 Token Structure Probing (RQ2)

To create a comprehensive test set for evaluating the tokenization capabilities of LLMs, we followed a systematic data synthesis process. Initially, we manually collected a set of around 300 words from the web, ensuring a diverse representation of word structures. This collection included words with common suffixes, prefixes, and varying lengths to cover a broad range of token structures.

Next, we defined a set of rules to create tasks for both intra-token and inter-token evaluations.

1125 These rules were designed to test different aspects
 1126 of tokenization, such as character counting, character
 1127 identification, case conversion, and identifying
 1128 common substrings and subsequences.

1129 Then we generated the probing tasks. For intra-
 1130 token evaluations, we created tasks like Character
 1131 Count (CC), N -th Character (NC), N -th Character
 1132 Reverse (NCR), and Case Conversion (CCV). For
 1133 inter-token evaluations, we developed tasks such as
 1134 Common Substrings (CS), Longest Common Sub-
 1135 strings (LCS), and Longest Common Subsequences
 1136 (LCSeq). Each task was carefully crafted to test
 1137 the model’s ability to understand and manipulate
 1138 token structures at various levels.

1139 C.2 Typographical Variation Task (RQ3)

Task	Test
GSM8k	1319
MMLU	14,042
TruthfulQA	817
HumanEval	164

Table 3: Dataset statistics of typographical variation tasks (RQ3).

1140 To thoroughly evaluate the typographical varia-
 1141 tion (RQ3), we conduct corruption to the questions
 1142 of several benchmark datasets’ test split and keep
 1143 answers intact. Details of evaluation dataset are
 1144 summarized at Table 3.

1145 **Character-Level** We performed character-level
 1146 corruption within word boundaries, ensuring that
 1147 the positions of punctuation marks or spaces re-
 1148 mained unchanged after corruption. Even when
 1149 there were not enough characters to form a com-
 1150 plete n -gram, corruption was still applied.

1151 **Token-Level** We converted the input prompts
 1152 into tokens using tokenizers from different model
 1153 families. Corruption was then applied within an
 1154 n -gram range, without regard to word boundaries.
 1155 This approach simulates token-level typographical
 1156 variations, challenging the models to handle disrup-
 1157 tions at the token level.

1158 **Permutation** For permutation task, we randomly
 1159 shuffle tokens or characters within an n -gram range
 1160 with a 50% probability, as illustrated in Figure 8.
 1161 We evaluate by using n -gram range from 2 to 5,
 1162 in which various n -gram levels could assess the

1163 model’s performance under varying degrees of cor-
 1164 ruption.

1165 **Noise Injection** We consider three kinds of noise
 1166 that commonly encountered by humans: insertion
 1167 (10%), deletion (10%), and replacement (10%). For
 1168 insertion, we choose a token or character from the
 1169 current n -gram under 50% circumstances. Other-
 1170 wise, we randomly select a token or character from
 1171 the whole vocabulary (token or character) and add
 1172 it to a random position.

1173 When performing replacing, we replace a token
 1174 or character in the current n -gram with a token or
 1175 character randomly selected from the whole vocabu-
 1176 lary. Figure 8 illustrates our approach of injecting
 1177 noise to tokens using a tri-gram granularity.

D Detailed Results of BPE-dropout Post-Training

1178 The results, as shown in Figure 10, illustrate the
 1179 effect of BPE-dropout on EM scores across seven
 1180 tasks (CS, LCSeq, CCV, CC, NC, LCS, NCR) un-
 1181 der various post-training conditions (0-shot, 1-shot,
 1182 2-shot, and 3-shot) with dropout rates ranging from
 1183 0.0 to 0.8.

1184 We observe that moderate dropout rates (0.2 to
 1185 0.4) appear to improve the convergence of EM
 1186 scores across all tasks, particularly in the zero-shot
 1187 setting. This indicates that a certain level of vari-
 1188 ability introduced by dropout can help the model
 1189 generalize better when no additional examples are
 1190 provided. In the 1-shot, 2-shot, and 3-shot settings,
 1191 moderate dropout rates contribute to stabilizing per-
 1192 formance, suggesting that this level of dropout in-
 1193 troduces useful regularization without significantly
 1194 compromising token integrity.

1195 It is evident that higher dropout rates (0.6 and
 1196 0.8) lead to a noticeable decline in EM scores
 1197 across all tasks and post-training conditions. This
 1198 indicates that excessive dropout disrupts the tok-
 1199 enization process, resulting in subwords with fewer
 1200 merges that conflict with the original pre-training
 1201 of the models. Tasks such as NC, LCS, and NCR
 1202 show more significant drops in EM scores with
 1203 higher dropout rates, reflecting their complexity
 1204 and the challenge of maintaining token integrity
 1205 under substantial dropout.

1206 Task-specific performance varies under differ-
 1207 ent dropout conditions. CS and CC tasks exhibit
 1208 high robustness to dropout, maintaining relatively
 1209 stable EM scores even at moderate dropout rates.
 1210 This suggests that the nature of these tasks makes
 1211
 1212

Origin	<pre>def largest_divisor(n: int) -> int: """ For a given number n, find the largest number that divides n evenly, smaller than n >>> largest_divisor(15) 5 """</pre>
Char-Level Permutation (5-gram)	<pre>def largest_divisor(n: int) -> int: """ For a given number n, fdin teh raglets number htat dividse n leenvy, asllmer than n >>> largest_divisor(15) 5 """</pre>

Figure 8: Example of a char-level permutation(5-gram) prompt from HumanEval.

Origin	<pre>[8100, 50777, 2380, 369, 17954, 1475, 6693, 323, 293, 2094, 55404, 1354, 369, 1077, 4885, 1475, 1938, 449, 3116, 13, 3005, 31878, 279, 27410, 520, 279, 20957, 6, 3157, 7446, 369, 400, 17, 824, 7878, 37085, 19151, 13]</pre>
Token-Level Noise (3-gram)	<pre>[8100, 50777, 2380, 412, 17954, 1475, 6693, 323, 293, 2094, 4124, 55404, 1354, 369, 1077, 4885, 1475, 1938, 449, 75316, 13, 3005, 31878, 279, 27410, 520, 279, 20957, 6, 3157, 7446, 369, 400, 6134, 17, 824, 7878, 37085, 19151, 13]</pre>

Figure 9: Example of token-level noise (tri-gram) injected prompt from GSM8K. Red, Green, and Gray denote replacement, insertion, and deletion respectively.

1213 them less sensitive to the variability introduced by
1214 dropout. LCSeq and CCV tasks show moderate
1215 sensitivity to dropout, with a noticeable decline in
1216 performance at higher dropout rates. NC, LCS, and
1217 NCR tasks are more adversely affected by higher
1218 dropout rates, indicating their reliance on stable
1219 token sequences.

1220 The positive effects of moderate BPE-dropout
1221 include improved generalization and regularization.
1222 Moderate dropout rates (0.2 to 0.4) introduce ben-
1223 efiticial variability, enhancing the model’s ability
1224 to generalize, particularly in zero-shot scenarios
1225 where the model must rely solely on its pre-trained
1226 knowledge without additional examples. BPE-
1227 dropout acts as a regularizer, preventing the model
1228 from overfitting to specific token sequences seen
1229 during pre-training. This is especially useful in low-
1230 resource settings (e.g., 0-shot and 1-shot) where

overfitting can be a significant concern.

1231
1232 However, challenges arise with high dropout
1233 rates. Higher dropout rates lead to subwords with
1234 fewer merges, deviating from the token sequences
1235 the model encountered during pre-training. This
1236 disruption results in poorer performance, as the
1237 model struggles to reconcile the altered tokeniza-
1238 tion with its pre-trained representations. More com-
1239 plex tasks (NC, LCS, NCR) suffer more from high
1240 dropout rates, highlighting the need for careful tun-
1241 ing of dropout rates based on task complexity and
1242 tokenization stability requirements.

1243 The findings suggest that there is an optimal
1244 range for BPE-dropout rates that balances the ben-
1245 efits of regularization and improved generalization
1246 with the need to maintain token integrity. Prac-
1247 titioners should consider moderate dropout rates
1248 to leverage the positive effects while avoiding the

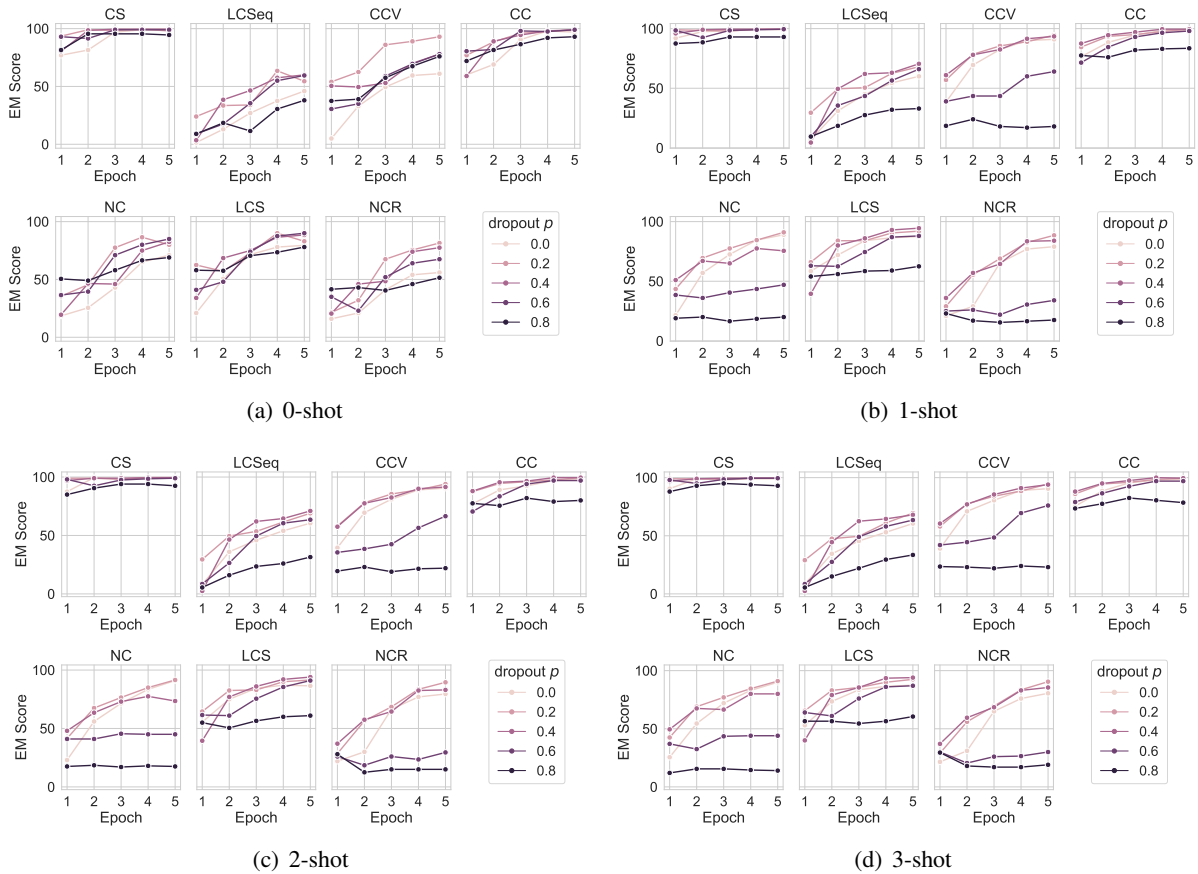


Figure 10: The impact of BPE-dropout on EM scores across seven tasks (CS, LCSeq, CCV, CC, NC, LCS, NCR) under different post-training conditions: (a) 0-shot, (b) 1-shot, (c) 2-shot, and (d) 3-shot. The dropout rates range from 0.0 to 0.8. The plots show that moderate dropout rates generally lead to improvements. Tasks such as CS and CC are more robust to dropout, maintaining higher scores even at moderate dropout rates, while tasks like NC, LCS, and NCR show significant performance drops with increasing dropout.

1249 pitfalls of excessive dropout. Different tasks exhibit
 1250 varying sensitivities to dropout, underscoring
 1251 the importance of task-specific dropout tuning to
 1252 achieve the best performance outcomes.

1253 E Additional Analysis

1254 E.1 Impact of Typographical Variations on 1255 Sequence Length

1256 Figure 11 shows a strong positive correlation between
 1257 token lengths before and after introducing
 1258 typographical errors across different tasks (GSM8K
 1259 and MMLU) and n -gram settings (2, 3, 5). We observe
 1260 that token lengths after introducing errors are
 1261 proportional to their original lengths. Besides, both
 1262 GSM8k and MMLU tasks exhibit similar patterns,
 1263 and the n -gram settings (2, 3, 5) do not significantly
 1264 alter this relationship. Most interestingly, the
 1265 slope for reorder errors is relatively larger than for
 1266 noise errors, indicating that **reorder errors tend
 1267 to result in a slightly greater increase in token**

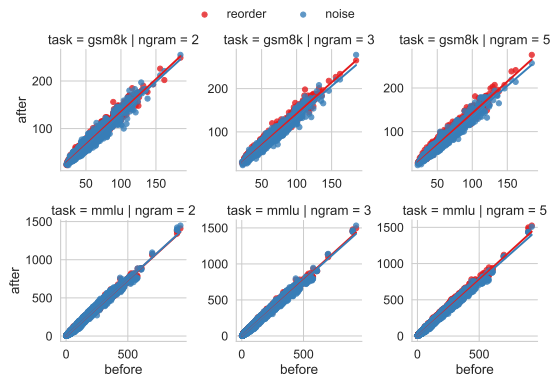


Figure 11: Scatter plots showing the positive correlation between token lengths before and after introducing typographical errors on MMLU and GSM8K, varying n -gram settings (2, 3, or 5). The x-axis represents the original token length, and the y-axis represents token length after adding errors.

length compared to noise errors.

E.2 Compositional Challenges in Token Embeddings

Figure 1 illustrates the compositional challenges faced by existing LLMs when handling subword units. Specifically, it presents cosine similarities and angular differences between embeddings of original words and their subword components.

In Figure 1(a), we observe the word "assignment" and its subword components "assign" and "ment." The cosine similarity between the composite embedding "assign + ment" and the original word "assignment" is relatively low at 0.51, with a significant angular difference of 59.24 degrees. This substantial disparity indicates that the learned token embeddings fail to capture the surface form composition accurately. The model does not recognize that "assign" combined with "ment" should semantically align closely with "assignment."

Similarly, Figure 1(b) shows the word "import" and its subword components "im" and "port". Here, although the cosine similarity between "im + port" and "import" is higher at 0.93, the angular difference is still notable at 21.17 degrees. This suggests that while the model captures some compositional aspects, it still struggles to fully integrate subword information to match the original word's embedding perfectly.

These observations highlight a critical limitation of existing LLMs: their learned token embeddings do not adequately capture the surface form composition. The inability to effectively combine subword units to represent the full word's meaning undermines the model's overall understanding and processing capabilities.