
HOW DOES A LANGUAGE-SPECIFIC TOKENIZER AFFECT LLMs?

Jean Seo, Jaeyoon Kim, SungJoo Byun, Hyopil Shin

Seoul National University

{seemdog, toscour345, byunsj, hpshin}@snu.ac.kr

ABSTRACT

The necessity of language-specific tokenizers intuitively appears crucial for effective natural language processing, yet empirical analyses on their significance and underlying reasons are lacking. This study explores how language-specific tokenizers influence the behavior of Large Language Models predominantly trained with English text data, through the case study of Korean. The research unfolds in two main stages: (1) the development of a Korean-specific extended tokenizer and (2) experiments to compare models with the basic tokenizer and the extended tokenizer through various Next Token Prediction tasks. Our in-depth analysis reveals that the extended tokenizer decreases confidence in incorrect predictions during generation and reduces cross-entropy in complex tasks, indicating a tendency to produce less nonsensical outputs. Consequently, the extended tokenizer provides stability during generation, potentially leading to higher performance in downstream tasks.

1 Introduction

The training data for Large Language Models (LLMs) today is heavily skewed towards English and code. For instance, Llama-2 [Touvron et al., 2023] from Meta-AI¹ is trained with 89.7% English, 8.38% code, while other languages make up the remaining 2%. Although these models are claimed to be multilingual, the strong emphasis on English affects the tokenizer’s ability to effectively process languages other than English, consequently impacting the model’s multilingual performance. As opposed to the tokenization of English text where text is segmented into appropriate units that capture linguistic information well, they often segment other languages—especially non-alphabetic ones—into bytes. Excessive byte-level segmentation can lead to issues such as limiting the maximum input and output length. Additionally, a byte token can be decoded into subwords of different languages depending on the neighboring tokens. This makes efficient representation embedding challenging, as a single entry may need to absorb too much information.

Moreover, the optimal tokenizing strategy for each language can vary depending on linguistic features, such as whether the language is agglutinative or inflective. This means the tokenizer of a certain LLM might not be fundamentally suitable for some languages. However, using a completely new tokenizer for a specific language is challenging because it would require retraining the language model from scratch. Therefore, many language-specific fine-tuned LLMs use tokenizer extension, incorporating a larger vocabulary size and additional merge rules for the target language. While extending the tokenizer also necessitates further training to align the embeddings, it is much more efficient than training the model from scratch with an entirely new vocabulary.

The necessity of language-specific tokenizers is supported by works such as Goldman et al. [2024], Alrefaie et al. [2024], which demonstrate that tokenization strategy and vocabulary size affect the performance of downstream tasks. However, deeper analysis is needed to understand exactly how tokenizers impact LLM performance. Such insights could explain why different tokenizers lead to varying model performance and increase the explainability of black-box LLMs. This, in turn, would provide evidence to inform decisions when building better models in the future.

In this research, we conduct an in-depth analysis of the effect language-specific tokenizers have on LLM performance. For an intrinsic analysis, we employ a Next Token Prediction (NTP) task with varying levels and develop corresponding metrics to evaluate the influence of different tokenizers. We experiment with Korean, a non-English, non-alphabetic language. The experimental results show that the extended tokenizer produces more sensible outputs more stably, with lower confidence levels in incorrect token generation and lower cross entropy in complex tasks. This underscores the

¹<https://www.meta.ai/>

importance of tokenizer extension, supporting claims from previous works that language specific tokenizers lead to elevated downstream task performance [Alrefaie et al., 2024, Ali et al., 2024].

The contributions of this study are:

- Intrinsically analyzing the effect of language-specific tokenizers on how LLMs work
- Proposing a novel framework for the intrinsic evaluation of tokenizers with Next Token Prediction task, using accuracy, confidence level, and cross entropy loss as metrics

2 Background and Related Work

Traditionally, Language Model (LM) tokenizers were primarily optimized for English, often resorting to byte-level segmentation for non-English languages. Segmenting text into byte-levels has a critical drawback: it greatly limits the maximum input and output length of the model. To address this issue, recent LLMs such as Llama-3 [AI@Meta, 2024] and Gemma [Team et al., 2024] adopt a different tokenization strategy, exponentially expanding the vocabulary size to 128,256 and 256,128 tokens, respectively. Therefore, a common approach has been developed to extend an LLM tokenizer to be specialized in one additional language besides English, expanding vocabulary and merge rules of only this target language. This approach is demonstrated by several Korean-extended models [L, 2024, Junbum Lee, Taekyoon Choi, 2024, L. Junbum, 2023, Kim et al., 2024b].

Prior research on tokenizers has predominantly focused on evaluating their effectiveness for specific languages [Rust et al., 2021, Toraman et al., 2023] or domains [Gopalakrishnan et al., 2023, Barrett and Weber, 2011]. Tokenization evaluation is essential for enhancing common knowledge, analyzing linguistic analysis, distinguishing processing types, and improving tokenizer design based on the presented dimensions [Habert et al., 1998]. Nonetheless, there is a clear lack of research directly examining their influence on LLMs. The few studies about evaluating the effect of tokenizers on LLMs include Goldman et al. [2024], Fujii et al. [2023], which are based on the downstream performance. Alrefaie et al. [2024] examined how different tokenization strategies and vocabulary sizes affect the performance of Arabic language models in downstream tasks. Ali et al. [2024] studied the influence of tokenizer choice on LLM downstream performance by training mono- and multilingual LLMs. However, it is uncertain whether downstream task performance can accurately measure the effectiveness of tokenizers. There remains a need for intrinsic and analytical studies to understand why performance differences occur. Such investigations are imperative to mitigate the opacity of black-box models and provide a rationale for decision-making towards the development of superior models.

Original Text

이번 방학 때 뭐해?

What is your plan for this vacation?

13 Characters

<s>

—

이
i

번
beon

—

방
bang

학
hag

—

때
ttae

—

뭐
mwo

해
hae

?

Base: 19 Tokens

<s>

—

이

E
B

2

8
8

—

방

학

—

E
B

9
5

8
C

—

E
B

A
D

9
0

해

?

Extended: 8 Tokens

<s>

—이번

—방

학

—때

—뭐

해

?

this

vacation

in

what

do

Figure 1: An example comparing how the Llama-2 base tokenizer and our extended tokenizer process the same sentence. The example sentence consists of 13 characters, including the special token <s> and whitespace. The base tokenizer segments 3 of these characters into bytes, resulting in a total of 19 tokens. In contrast, our extended tokenizer appropriately tokenizes the sentence into meaningful units of subwords, resulting in only 8 tokens.

3 Extended Tokenizer

Well-known open-source models, such as Mistral [Jiang et al., 2023] Llama-2 [Touvron et al., 2023], Yi [AI et al., 2024], MT5 [Xue et al., 2021], Giraffe [Niemeyer and Geiger, 2021], Solar [Kim et al., 2024a] use SentencePiece BPE. Meanwhile, models including Falcon [Almazrouei et al., 2023], Llama-3- [AI@Meta, 2024], GPT-Neox [Black et al.,

Tokenizer	Extension	% of Unknown Tokens	Average Number of Tokens
Mistral-7B-Instruct-v0.1	N	14.36	242.71
Mistral-ko-7B-v0.1	Y	1.0	128.56
Llama-2-7b	N	48.61	325.27
Llama-2-ko-7b	Y	0.98	128.92
Llama2-extended(Ours)	Y	0.96	125.91

Table 1: The percentage of unknown tokens and the average number of tokens per sentence in our Korean test sentences evaluated across various tokenizers. Unknown tokens refer to tokens that are segmented into bytes. We used different text data for tokenizer evaluation than the text data used for tokenizer extension. Extended tokenizers show a significantly lower percentage of unknown tokens, demonstrating less text segmented into bytes, and a reduced average number of tokens, indicating a longer maximum input and output length. Amongst the extended tokenizers, our Llama2-extended tokenizer shows the least percentage of unknown tokens and the shortest average number of tokens of sentences.

2022], Bloom [Workshop et al., 2023], Qwen [Bai et al., 2023] use ByteLevel BPE as their tokenization strategy. This shows that SentencePiece BPE and ByteLevel BPE are currently the most common tokenization methods for LLMs.

In this study, we employ TinyLlama [Zhang et al., 2024] to investigate the effect of language-specific tokenizers. TinyLlama, which has 1.1 billion parameters, is relatively small and trained with barely any amount of Korean text. It uses the Llama-2 tokenizer which is based on SentencePiece BPE, combining SentencePiece [Kudo and Richardson, 2018] with the BPE algorithm. The main advantage of SentencePiece BPE is that predicted tokens can be directly joined into sentences. Additionally, it offers the Byte Fallback option, which processes unknown input sequences into bytes based on the UTF-8 encoding of characters, enhancing its ability to handle unknown tokens.

Despite its ability to handle unseen tokens by splitting them into bytes, SentencePiece BPE can cause issues when processing languages underrepresented in the training data or vocabulary. This is because a byte can be decoded into various words or subwords depending on its context, making efficient learning of representations difficult. Byte-level segmentation also reduces the model’s maximum input and output sequence length.

To address these issues, we build a language-specific tokenizer through vocabulary and merge rule extension. For our case study on Korean, we train a SentencePiece BPE tokenizer on Korean texts and append the new Korean vocabulary and merge rules to those of TinyLlama’s tokenizer. Figure 1 shows that the extended tokenizer segments Korean sentences into bytes less frequently and into larger subwords with lexical and grammatical meaning. Furthermore, Table 1 demonstrates the superior performance of the extended tokenizer compared to other LLM tokenizers.

4 Next Token Prediction

In order to deeply analyze the effect of tokenizers on LLMs, we introduce a NTP(Next Token Prediction) task along with corresponding evaluation metrics. This approach transcends mere observation of downstream task performances, enabling a more comprehensive analysis of tokenizer efficacy, possibly explaining how LLMs with certain tokenizers yield better performance compared to others.

4.1 Task

Unlike Encoder-only models such as BERT [Devlin et al., 2019] which refer to both preceding and subsequent tokens to predict masked tokens, LLMs rely solely on previous tokens for causal language modeling. Therefore, for the NTP task, we meticulously curated test sentences featuring a target token, ensuring that only one answer is possible given the preceding tokens. Three human annotators manually went through all test sentences and sorted out the final ones with masked tokens that admit only one possible answer. Our Next Token Prediction task comprises six distinct tasks, categorized into two difficulty levels (easy and hard) and three target units (token, character, and word).

4.1.1 Difficulty

The difficulty level of the NTP task depends upon the amount of input text provided to the model. Each test sentence is divided into three segments: the sequence preceding the target, the target itself, and the sequence following the target. In the easy version, the full sentence and the sequence preceding the target are provided as input, whereas only the

sequence preceding the target is given in the hard version. Given that the model has already been exposed to the answer in the easy version, we presume predicting the target token to be considerably simpler than in the hard version. Further details regarding the two versions of the NTP task can be found in Figure 2.

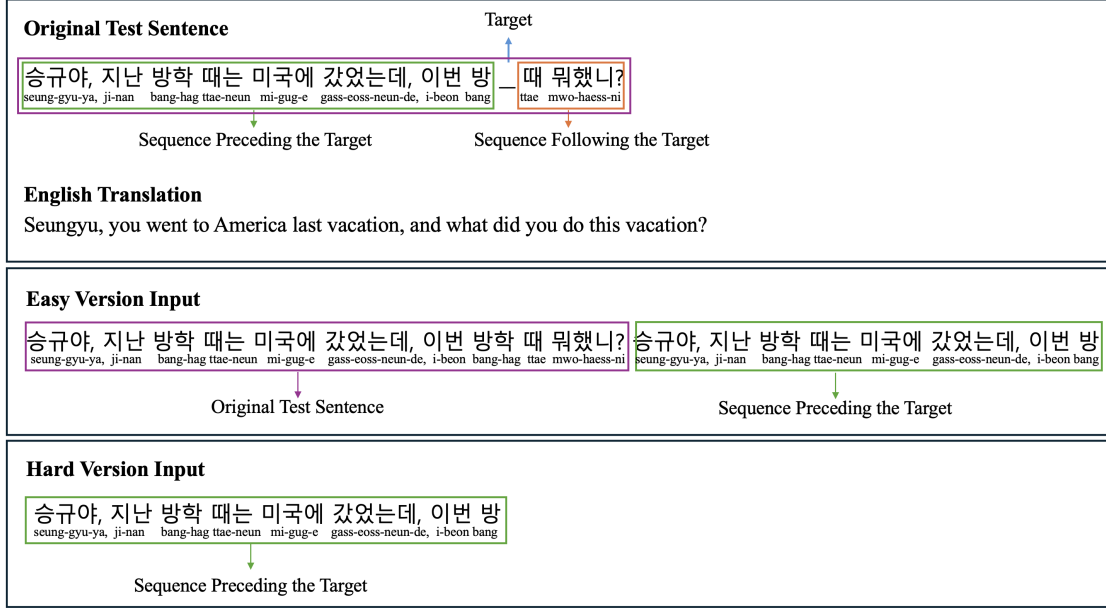


Figure 2: Example of the NTP task input in easy and hard version respectively. In the easy version, the answer is already provided to the model since the input includes the original test sentence. In contrast, in the hard version, only the sequence preceding the target is provided, making it more difficult for the model to predict the correct token.

4.1.2 Target Unit

In our comparison between two distinct tokenizers, the base tokenizer and the extended tokenizer, we categorize the target unit into three types: (1) token level, where both tokenizers tokenize the target into a single token; (2) character level, where the extended tokenizer tokenizes the target into a single token while the base tokenizer splits it into multiple tokens; and (3) word level, where both tokenizers tokenize the target into several tokens. Figure 3 provides a detailed illustration of these three variations in the target units for the NTP task.

4.2 Evaluation Metric

To assess the intrinsic impact of tokenizers in LLMs through the NTP task, we employ three metrics: accuracy, confidence level, and cross-entropy loss.

4.2.1 Accuracy

We assess whether the model correctly predicts the target in a binary manner. When the target consists of a single token, determining accuracy is straightforward, as it can be classified as either correct or incorrect with ease. However, if the target comprises multiple tokens, we consider it correct only if the model predicts all tokens accurately. For instance, if the target is tokenized into three tokens and the model correctly predicts two tokens but incorrectly predicts one, the model’s prediction is deemed incorrect. Evaluating the model’s predictions binaryly, we compute the average accuracy across all test sentences.

4.2.2 Confidence Level

Building upon prior research that incorporates the notion of confidence level in deep learning, such as Inoue [2019], we examine the model’s confidence level during the generation of each target token. This approach allows for a more intrinsic evaluation and analysis of the effect of tokenizers on LLMs. Similar to Schuster et al. [2022], Niehues and Pham [2019], where confidence level is used specifically in Language Modeling, we consider the probability calculated using the softmax layer as the confidence level at each time step. Formally, given an input token sequence $\hat{X} = x_1, \dots, x_{i-1}$,

Token Level <div style="text-align: right;">Target ↑</div> 지수가 편지를 봉투__ 넣었다. ji-su- ga pyeon-ji-leul bong-tu-__ neoh-eoss-da. Jisoo put the letter into the envelope.	Target: 예 e <i>예: into</i> Base tokenizer: 예 Extended tokenizer: 예
Character Level <div style="text-align: right;">Target ↑</div> 아침에 씻었는데 또 샤워__해? a-chim-e ssiss-eoss-neun-de tto sya-__-hae? You already showered in the morning, are you going to shower again?	Target: 워 워 <i>샤워: shower</i> Base tokenizer: <0xEC>, <0x9B>, <0x8C> Extended tokenizer: 워
Word Level <div style="text-align: right;">Target ↑</div> 올해 말에 지구의 종__ 온대. ol-hae. mal-e. ji-gu-ui jong-__ on-dae. The end of the world is coming at the end of this year.	Target: 말이 mal-i <i>종말이: end</i> Base tokenizer: <0xEB>, <0xA7>, <0x90>, 이 Extended tokenizer: 말, 이

Figure 3: Example of the NTP task input in the three different target units. At the token level, the target is tokenized into a single token by both tokenizers. At the character level, the target is segmented into 3 tokens by the base tokenizer but is treated as a single token by the extended tokenizer. At the word level, both tokenizers split the target into multiple tokens.

the model predicts the next token x_i by computing the probabilities of all possible tokens in the vocabulary V . These probabilities are obtained using the softmax function applied to the logits (the raw scores) produced by the model. Let z_j denote the logit corresponding to the j -th token in the vocabulary V , where $j \in \{1, \dots, |V|\}$. The softmax probability for the j -th token is given by:

$$P(x_i = j \mid X) = \frac{\exp(z_j)}{\sum_{k=1}^{|V|} \exp(z_k)}$$

The confidence level at each time step, which is the maximum softmax probability, is then calculated as:

$$\text{Confidence}(X) = \max_{j \in \{1, \dots, |V|\}} P(x_i = j \mid X) = \max_{j \in \{1, \dots, |V|\}} \frac{\exp(z_j)}{\sum_{k=1}^{|V|} \exp(z_k)} \quad (1)$$

To ensure more fair and accurate comparison between tokenizers with varying vocabulary sizes, we conduct normalization and additionally compare the normalized confidence levels. This necessity stems from the inherent characteristics of the softmax layer, wherein larger embedding sizes (or vocabulary sizes in this context) often lead to decreased probabilities assigned to each predicted token. To normalize the confidence level at time step t , we first calculate the average confidence level of the previous $t - 1$ time steps:

$$\overline{\text{Confidence}}(t - 1) = \frac{1}{t - 1} \sum_{k=1}^{t-1} \text{Confidence}(k)$$

Then, we divide the confidence level at the current time step by this average:

$$\text{NormalizedConfidence}(t) = \frac{\text{Confidence}(t)}{\overline{\text{Confidence}}(t - 1)}$$

By using the normalized confidence level, we can compare the effects of the base and extended tokenizers while minimizing the impact of differing vocabulary sizes.

4.2.3 Cross Entropy Loss

Another metric we employ to compare the impact of different tokenizers on LLMs is Cross Entropy Loss. Cross Entropy Loss can be intuitively interpreted as a measure of how uncertain the model is when predicting the next token.

5 Experiments

5.1 Models and Data

To evaluate the isolated impact of the extended tokenizer, we train two models with all settings identical except for the tokenizer. Due to cost issues, we utilize a relatively small LLM, Tinyllama which has 1.1 billion trainable parameters as the base model. Our primary objective is to assess the influence of a tokenizer featuring a language-specific extended vocabulary—in this case, Korean—thus we compile a bilingual Korean-English pretraining dataset. Additionally, as our focus is on observing intrinsic effects rather than downstream task performance, we opt for a pretrained model instead of instruction-tuned model. The dataset comprises Korean news data from the Modu Corpus², English data from Wikipedia³, and Alpaca dataset Taori et al. [2023]. Given that the dataset used to train Tinyllama [Soboleva et al., 2023, Li et al., 2023] barely includes any Korean and our aim is to scrutinize the intrinsic impact of tokenizers on the LLM’s ability to generate Korean, we maintain a dataset ratio of Korean to English of 9:1, enabling the model to become more adept with Korean. Furthermore, to investigate whether the quantity of training data correlates with the intrinsic ability of Korean token generation, we save intermediate checkpoints and evaluate them accordingly.

5.2 Experimental Settings

Each model was trained using identical hyperparameters: batch size = 4, linear learning rate scheduler, and weight decay = 0.01. Training was conducted using 3 A100 GPUs, with each checkpoint taking approximately 27 hours to be trained.

6 Result and Analysis

As mentioned in Section 5.1, we saved a total of four checkpoints at different training steps. We primarily used the final checkpoint, which is the most trained, to compare the three metrics: accuracy, normalized confidence level, and cross-entropy loss. The results for accuracy, normalized confidence level, and cross-entropy loss are presented in Figures 4, 5, and 6.

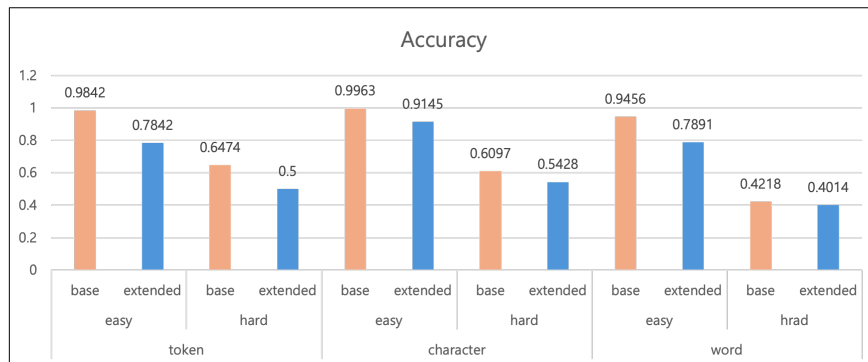


Figure 4: Accuracy

Accuracy As opposed to our expectation, the model with the base tokenizer shows higher accuracy than the one with the extended tokenizer, though this difference decreases at larger unit levels: token > character > word, with the hard, word level showing almost the same accuracy for both models. This suggests that if the NTP task targets a larger number of tokens, the accuracy ranking might reverse. The accuracy on the easy level is higher than on the hard level across all unit levels for both the models as predicted. The expected accuracy rankings were token, character, word. However, only the base tokenizer model on the hard level followed this ranking. On the easy level, the base tokenizer

²<https://kli.korean.go.kr/corpus/main/requestMain.do?lang=en>

³<https://www.wikipedia.org/>

model ranked as character, token, word, while the extended tokenizer model ranked as character, word, token on the easy level and character, token, word on the hard level. Notably, the accuracy ranking of the token level for the extended tokenizer model tend to be lower than for the base tokenizer model. This discrepancy likely arises because the extended tokenizer includes a larger vocabulary, where tokens can overlap (e.g., ‘my’ and ‘mystery’). This overlap can lower accuracy when there are various larger tokens that encompass the target token.

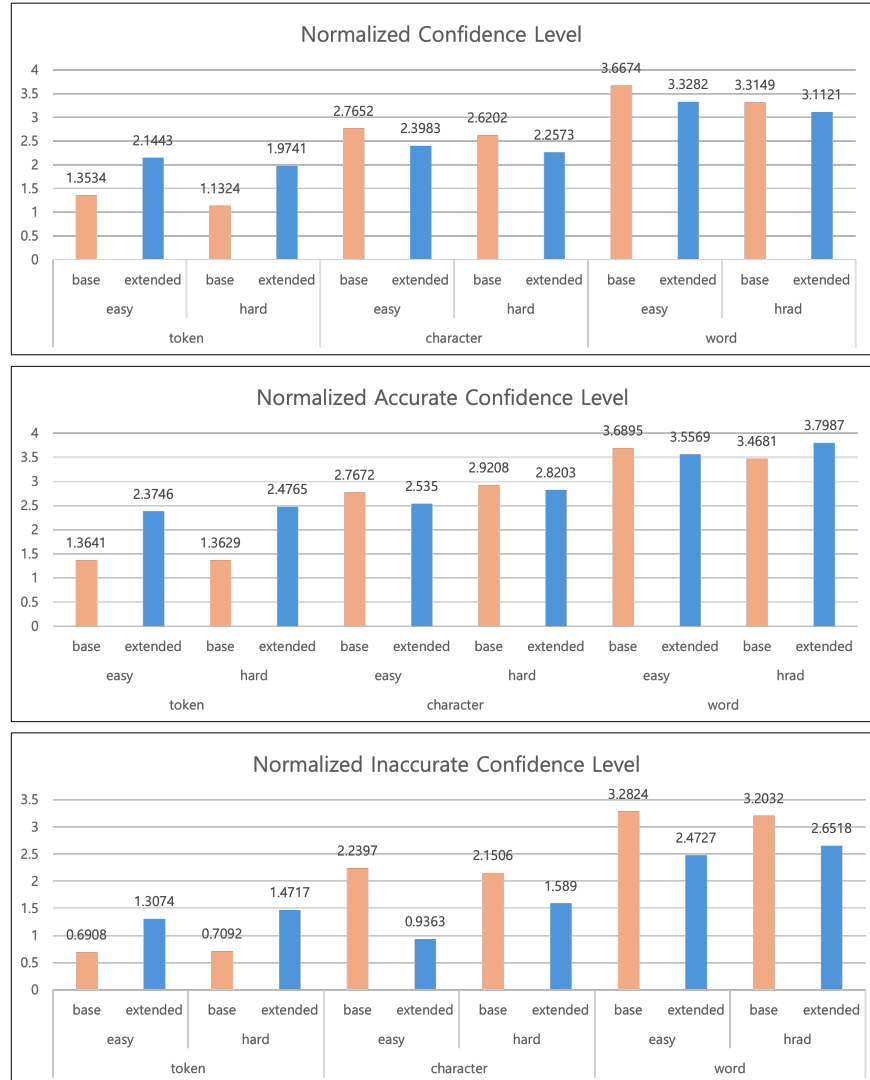


Figure 5: Normalized Confidence Level

Normalized Confidence Level As shown in Fig 5, the confidence level is consistently lower for the hard level tasks compared to the easy level tasks across all unit levels and both models. This indicates that the models are more certain about their predictions in easier tasks, which aligns with our expectation. The model with the extended tokenizer exhibited higher confidence in the token-level task than the base tokenizer model. However, the opposite trend was observed in the character-level and word-level tasks.

Normalized Accurate Confidence Level For a more detailed analysis, we further compare the confidence levels in cases where the model predicted correctly and incorrectly, referring to these as the normalized accurate confidence level and normalized inaccurate confidence level, respectively. For the normalized accurate confidence level, the differences between the base tokenizer model and the extended tokenizer model were generally smaller than the overall confidence level difference. In the word-level hard task, ranking even reversed. Additionally, for the character-level task, the base tokenizer model showed higher normalized accurate confidence for the hard task than the easy task. Notably, the

extended tokenizer model had a higher normalized accurate confidence level in the hard tasks than in the easy tasks across all three tasks, indicating increased certainty in its predictions for the harder versions.

Normalized Inaccurate Confidence Level The base tokenizer model showed less normalized inaccurate confidence than the extended tokenizer model in the token-level task. Conversely, for the character-level and word-level tasks, the base tokenizer model exhibited significantly higher normalized inaccurate confidence compared to the extended tokenizer model. This suggests that the base tokenizer model is more likely to produce nonsensical outputs with high certainty, whereas the extended tokenizer model is more cautious in its predictions.

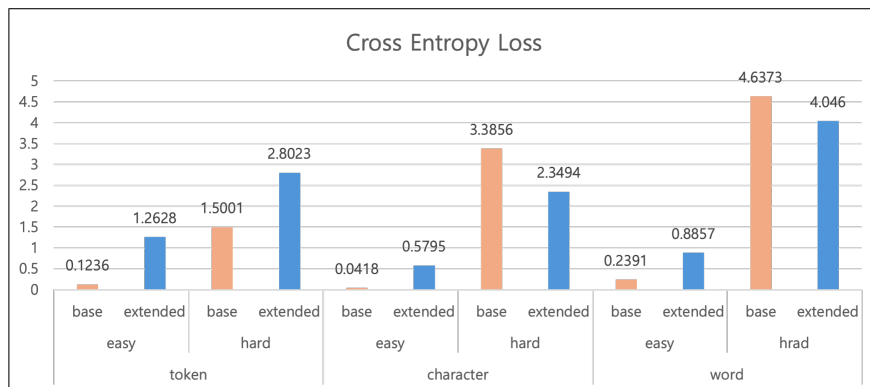


Figure 6: Cross Entropy Loss

Cross Entropy Loss Cross entropy loss measures the confusion of the model when generating an output. As expected, both models exhibit higher cross entropy loss in the hard tasks compared to the easy tasks across all unit levels. In the easy tasks, the cross entropy loss of the extended tokenizer model is greater than that of the base tokenizer model. However, in the hard versions of the character-level and word-level tasks, the extended tokenizer model shows lower cross entropy loss. This suggests that the extended tokenizer model performs better in more complex tasks, as character-level and word-level tasks are considered more challenging than token-level tasks.

Summing up the results, while the extended tokenizer model may not necessarily be more accurate than the base tokenizer model in granular tasks, such as predicting a small number of next tokens, it appears to exhibit greater stability. The extended tokenizer model shows less confidence when generating inappropriate tokens and handles difficult tasks with less perplexity compared to the base tokenizer model. This increased stability in generating sensible outputs is likely to positively influence the downstream task performance of language models.

To examine if there are differences in results based on the training steps, we analyzed the changes in accuracy, confidence level, and cross entropy loss at each checkpoint for each task. The results in Appendix A indicate that after a certain number of training steps, the results tend to converge, suggesting that the results can be generalized at least to TinyLlama based models.

7 Conclusion

In this work, we investigate the influence of language-specific extended tokenizers on Large Language Models, focusing on a Next Token Prediction Task with two difficulty levels and three unit levels. Using accuracy, confidence level, and cross entropy loss as evaluation metrics, we observe that the extended tokenizer provides more stable generation, exhibiting less confidence in inaccurate token generation and lower cross entropy loss on complex tasks. Our study is the first to explore the internal impact of tokenizers and suggest the possible underlying mechanism of how language-specific tokenizers could enhance downstream task performance.

Limitation and Future Work

Since this study experiments with a relatively small LLM, containing only 1.1 billion parameters, the results cannot be directly generalized to larger models. Additionally, as this is a case study on Korean, it is uncertain whether the same results would apply to other languages. We believe that similar intrinsic research on the effect of different tokenizers on LLMs should be further investigated, varying both the sizes of the models and the languages studied.

References

01. AI, :, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. Yi: Open foundation models by 01.ai, 2024.
- AI@Meta. Llama 3 model card, 2024. URL https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md.
- Mehdi Ali, Michael Fromm, Klaudia Thellmann, Richard Rutmann, Max Lübbering, Johannes Leveling, Katrin Klug, Jan Ebert, Niclas Doll, Jasper Schulze Buschhoff, Charvi Jain, Alexander Arno Weber, Lena Jurkschat, Hammam Abdelwahab, Chelsea John, Pedro Ortiz Suarez, Malte Ostendorff, Samuel Weinbach, Rafet Sifa, Stefan Kesselheim, and Nicolas Flores-Herr. Tokenizer choice for llm training: Negligible or crucial?, 2024.
- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, M  rouane Debbah,   tienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, Daniele Mazzotta, Badreddine Noun  , Baptiste Pannier, and Guilherme Penedo. The falcon series of open language models, 2023.
- Mohamed Taher Alrefaie, Nour Eldin Morsy, and Nada Samir. Exploring tokenization strategies and vocabulary sizes for enhanced arabic language models, 2024.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report, 2023.
- Neil Barrett and Jens Weber. Building a biomedical tokenizer using the token lattice design pattern and the adapted viterbi algorithm. *BMC bioinformatics*, 12 Suppl 3:S1, 06 2011. doi: 10.1186/1471-2105-12-S3-S1.
- Sidney Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, Usvsn Sai Prashanth, Shivanshu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. GPT-NeoX-20B: An open-source autoregressive language model. In Angela Fan, Suzana Ilic, Thomas Wolf, and Matthias Gall  , editors, *Proceedings of BigScience Episode #5 – Workshop on Challenges & Perspectives in Creating Large Language Models*, pages 95–136, virtual+Dublin, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.bigscience-1.9. URL <https://aclanthology.org/2022.bigscience-1.9>.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
- Takuro Fujii, Koki Shibata, Atsuki Yamaguchi, Terufumi Morishita, and Yasuhiro Sogawa. How do different tokenizers perform on downstream tasks in scriptio continua languages?: A case study in japanese, 2023.
- Omer Goldman, Avi Caciularu, Matan Eyal, Kris Cao, Idan Szpektor, and Reut Tsarfaty. Unpacking tokenization: Evaluating text compression and its correlation with model performance, 2024.
- Seethalakshmi Gopalakrishnan, Victor Zitian Chen, Wenwen Dou, and Wlodek Zadrozny. Reducing tokenizer’s tokens per word ratio in financial domain with t-mufin bert tokenizer. In *FINNLP*, 2023. URL <https://api.semanticscholar.org/CorpusID:263610065>.
- Beno  t Habert, Gilles Adda, Maria Decker, Philippe Boula De Mareuil, Silvana Ferrari, Olivier Ferret, Gabriel Illouz, and P. Paraubeck. Towards tokenization evaluation. In *International Conference on Language Resources and Evaluation*, 1998. URL <https://api.semanticscholar.org/CorpusID:63999585>.
- Hiroshi Inoue. Adaptive ensemble prediction for deep neural networks based on confidence level, 2019.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. Mistral 7b, 2023.
- Junbum Lee, Taekyoon Choi. gemma-ko-7b, 2024. URL <https://huggingface.co/beomi/gemma-ko-7b>.
- Dahyun Kim, Chanjun Park, Sanghoon Kim, Wonsung Lee, Wonho Song, Yunsu Kim, Hyeonwoo Kim, Yungi Kim, Hyeonju Lee, Jihoo Kim, Changbae Ahn, Seonghoon Yang, Sukyung Lee, Hyunbyung Park, Gyoungjin Gim, Mikyoung Cha, Hwalsuk Lee, and Sunghun Kim. Solar 10.7b: Scaling large language models with simple yet effective depth up-scaling, 2024a.

- Seungduk Kim, Seungtaek Choi, and Myeongho Jeong. Efficient and effective vocabulary expansion towards multilingual large language models, 2024b.
- Taku Kudo and John Richardson. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing, 2018.
- Junbum L. Llama-3-open-ko. 2024. URL <https://huggingface.co/beomi/Llama-3-Open-Ko-8B>.
- L. Junbum. llama-2-ko-7b (revision 4a9993e), 2023. URL <https://huggingface.co/beomi/llama-2-ko-7b>.
- Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João Monteiro, Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee, Logesh Kumar Umapathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra Murthy, Jason Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang, Nour Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh, Sasha Luccioni, Paulo Villegas, Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee, Nadav Timor, Jennifer Ding, Claire Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex Gu, Jennifer Robinson, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder: may the source be with you!, 2023.
- Jan Niehues and Ngoc-Quan Pham. Modeling confidence in sequence-to-sequence models, 2019.
- Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields, 2021.
- Phillip Rust, Jonas Pfeiffer, Ivan Vulić, Sebastian Ruder, and Iryna Gurevych. How good is your tokenizer? on the monolingual performance of multilingual language models. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3118–3135, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.243. URL <https://aclanthology.org/2021.acl-long.243>.
- Tal Schuster, Adam Fisch, Jai Gupta, Mostafa Dehghani, Dara Bahri, Vinh Q. Tran, Yi Tay, and Donald Metzler. Confident adaptive language modeling, 2022.
- Daria Soboleva, Faisal Al-Khateeb, Robert Myers, Jacob R Steeves, Joel Hestness, and Nolan Dey. SlimPajama: A 627B token cleaned and deduplicated version of RedPajama. <https://www.cerebras.net/blog/slimpajama-a-627b-token-cleaned-and-deduplicated-version-of-redpajama>, 2023. URL <https://huggingface.co/datasets/cerebras/SlimPajama-627B>.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikula, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. Gemma: Open models based on gemini research and technology, 2024.
- Cagri Toraman, Eyup Halit Yilmaz, Furkan Şahi-nuç, and Oguzhan Ozelik. Impact of tokenization on language models: An analysis for turkish. *ACM Transactions on Asian and Low-Resource Language Information Processing*, 22(4): 1–21, March 2023. ISSN 2375-4702. doi: 10.1145/3578707. URL <http://dx.doi.org/10.1145/3578707>.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023.

BigScience Workshop, :, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurençon, Yacine Jernite, Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue, Christopher Klam, Colin Leong, Daniel van Strien, David Ifeoluwa Adelani, Dragomir Radev, Eduardo González Ponferrada, Efrat Levkovizh, Ethan Kim, Eyal Bar Natan, Francesco De Toni, Gérard Dupont, Germán Kruszewski, Giada Pistilli, Hady Elsahar, Hamza Benyamina, Hieu Tran, Ian Yu, Idris Abdulmumin, Isaac Johnson, Itziar Gonzalez-Dios, Javier de la Rosa, Jenny Chim, Jesse Dodge, Jian Zhu, Jonathan Chang, Jörg Froberg, Joseph Tobing, Joydeep Bhattacharjee, Khalid Almubarak, Kimbo Chen, Kyle Lo, Leandro Von Werra, Leon Weber, Long Phan, Loubna Ben allal, Ludovic Tanguy, Manan Dey, Manuel Romero Muñoz, Maraim Masoud, María Grandury, Mario Šaško, Max Huang, Maximin Coavoux, Mayank Singh, Mike Tian-Jian Jiang, Minh Chien Vu, Mohammad A. Jauhar, Mustafa Ghaleb, Nishant Subramani, Nora Kassner, Nurulaqilla Khamis, Olivier Nguyen, Omar Espejel, Ona de Gibert, Paulo Villegas, Peter Henderson, Pierre Colombo, Priscilla Amuok, Quentin Lhoest, Rhea Harliman, Rishi Bommasani, Roberto Luis López, Rui Ribeiro, Salomey Osei, Sampo Pyysalo, Sebastian Nagel, Shamik Bose, Shamsuddeen Hassan Muhammad, Shanya Sharma, Shayne Longpre, Somaieh Nikpoor, Stanislav Silberberg, Suhas Pai, Sydney Zink, Tiago Timponi Torrent, Timo Schick, Tristan Thrush, Valentin Danchev, Vassilina Nikoulina, Veronika Laippala, Violette Lepercq, Vrinda Prabhu, Zaid Alyafeai, Zeerak Talat, Arun Raja, Benjamin Heinzerling, Chenglei Si, Davut Emre Taşar, Elizabeth Salesky, Sabrina J. Mielke, Wilson Y. Lee, Abheesht Sharma, Andrea Santilli, Antoine Chaffin, Arnaud Stiegler, Debajyoti Datta, Eliza Szczechla, Gunjan Chhablani, Han Wang, Harshit Pandey, Hendrik Strobelt, Jason Alan Fries, Jos Rozen, Leo Gao, Lintang Sutawika, M Saiful Bari, Maged S. Al-shaibani, Matteo Manica, Nihal Nayak, Ryan Teehan, Samuel Albanie, Sheng Shen, Srulik Ben-David, Stephen H. Bach, Taewoon Kim, Tali Bers, Thibault Fevry, Trishala Neeraj, Urmish Thakker, Vikas Raunak, Xiangru Tang, Zheng-Xin Yong, Zhiqing Sun, Shaked Brody, Yallow Uri, Hadar Tojarieh, Adam Roberts, Hyung Won Chung, Jaesung Tae, Jason Phang, Ofir Press, Conglong Li, Deepak Narayanan, Hatim Bourfoune, Jared Casper, Jeff Rasley, Max Ryabinin, Mayank Mishra, Minjia Zhang, Mohammad Shoeybi, Myriam Peyrounette, Nicolas Patry, Nouamane Tazi, Omar Sanseviero, Patrick von Platen, Pierre Cornette, Pierre François Lavallée, Rémi Lacroix, Samyam Rajbhandari, Sanchit Gandhi, Shaden Smith, Stéphane Requena, Suraj Patil, Tim Dettmers, Ahmed Baruwa, Amanpreet Singh, Anastasia Cheveleva, Anne-Laure Ligozat, Arjun Subramonian, Aurélie Névél, Charles Lovering, Dan Garrette, Deepak Tunuguntla, Ehud Reiter, Ekaterina Taktasheva, Ekaterina Voloshina, Eli Bogdanov, Genta Indra Winata, Hailey Schoelkopf, Jan-Christoph Kalo, Jekaterina Novikova, Jessica Zosa Forde, Jordan Clive, Jungo Kasai, Ken Kawamura, Liam Hazan, Marine Carpuat, Miruna Clinciu, Najoung Kim, Newton Cheng, Oleg Serikov, Omer Antverg, Oskar van der Wal, Rui Zhang, Ruochen Zhang, Sebastian Gehrmann, Shachar Mirkin, Shani Pais, Tatiana Shavrina, Thomas Scialom, Tian Yun, Tomasz Limisiewicz, Verena Rieser, Vitaly Protasov, Vladislav Mikhailov, Yada Pruksachatkun, Yonatan Belinkov, Zachary Bamberger, Zdeněk Kasner, Alice Rueda, Amanda Pestana, Amir Feizpour, Ammar Khan, Amy Faranak, Ana Santos, Anthony Hevia, Antigona Unldreaj, Arash Aghagol, Arezoo Abdollahi, Aycha Tammour, Azadeh HajiHosseini, Bahareh Behroozi, Benjamin Ajibade, Bharat Saxena, Carlos Muñoz Ferrandis, Daniel McDuff, Danish Contractor, David Lansky, Davis David, Douwe Kiela, Duong A. Nguyen, Edward Tan, Emi Baylor, Ezinwanne Ozoani, Fatima Mirza, Frankline Ononiwu, Habib Rezanejad, Hessie Jones, Indrani Bhattacharya, Irene Solaiman, Irina Sedenko, Isar Nejadgholi, Jesse Passmore, Josh Seltzer, Julio Bonis Sanz, Livia Dutra, Mairon Samagaio, Maraim Elbadri, Margot Mieskes, Marissa Gerchick, Martha Akinlolu, Michael McKenna, Mike Qiu, Muhammed Ghauri, Mykola Burynok, Nafis Abrar, Nazneen Rajani, Nour Elkott, Nour Fahmy, Olanrewaju Samuel, Ran An, Rasmus Kromann, Ryan Hao, Samira Alizadeh, Sarmad Shubber, Silas Wang, Sourav Roy, Sylvain Viguier, Thanh Le, Tobi Oyeade, Trieu Le, Yoyo Yang, Zach Nguyen, Abhinav Ramesh Kashyap, Alfredo Palasciano, Alison Callahan, Anima Shukla, Antonio Miranda-Escalada, Ayush Singh, Benjamin Beilharz, Bo Wang, Caio Brito, Chenxi Zhou, Chirag Jain, Chuxin Xu, Clémentine Fourier,

Daniel León Perifán, Daniel Molano, Dian Yu, Enrique Manjavacas, Fabio Barth, Florian Fuhrmann, Gabriel Altay, Giyaseddin Bayrak, Gully Burns, Helena U. Vrabec, Imane Bello, Ishani Dash, Jihyun Kang, John Giorgi, Jonas Golde, Jose David Posada, Karthik Rangasai Sivaraman, Lokesh Bulchandani, Lu Liu, Luisa Shinzato, Madeleine Hahn de Bykhovetz, Maiko Takeuchi, Marc Pàmies, Maria A Castillo, Marianna Nezhurina, Mario Sängler, Matthias Samwald, Michael Cullan, Michael Weinberg, Michiel De Wolf, Mina Mihaljcic, Minna Liu, Moritz Freidank, Myungsun Kang, Natasha Seelam, Nathan Dahlberg, Nicholas Michio Broad, Nikolaus Muellner, Pascale Fung, Patrick Haller, Ramya Chandrasekhar, Renata Eisenberg, Robert Martin, Rodrigo Canalli, Rosaline Su, Ruisi Su, Samuel Cahyawijaya, Samuele Garda, Shlok S Deshmukh, Shubhanshu Mishra, Sid Kiblawi, Simon Ott, Sinee Sang-aaroonsiri, Srishti Kumar, Stefan Schweter, Sushil Bharati, Tanmay Laud, Théo Gigant, Tomoya Kainuma, Wojciech Kusa, Yanis Labrak, Yash Shailesh Bajaj, Yash Venkatraman, Yifan Xu, Yingxin Xu, Yu Xu, Zhe Tan, Zhongli Xie, Zifan Ye, Mathilde Bras, Younes Belkada, and Thomas Wolf. Bloom: A 176b-parameter open-access multilingual language model, 2023.

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.

Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tynyllama: An open-source small language model, 2024.

A Appendix

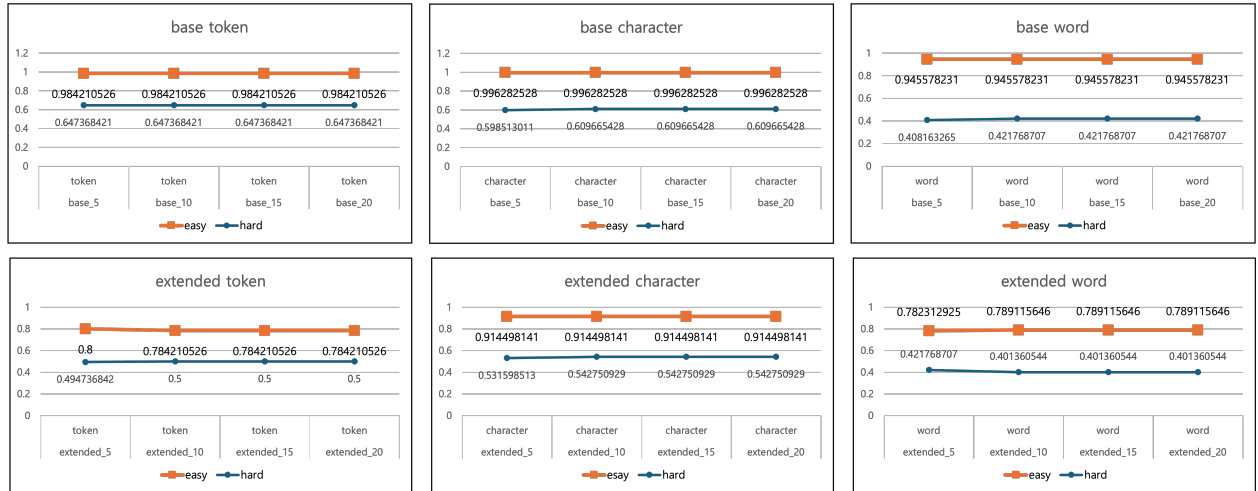


Figure 7: Accuracy of each checkpoint. The top three rows represent the base tokenizer model, while the bottom three rows represent the extended tokenizer model. The figures on the left show token-level results, the middle figures show character-level results, and the figures on the right show word-level results. Within each figure, the orange points indicate the easy level, and the blue points indicate the hard level. Additionally, within each figure, the leftmost points correspond to the least trained model, and the rightmost points correspond to the most trained model.

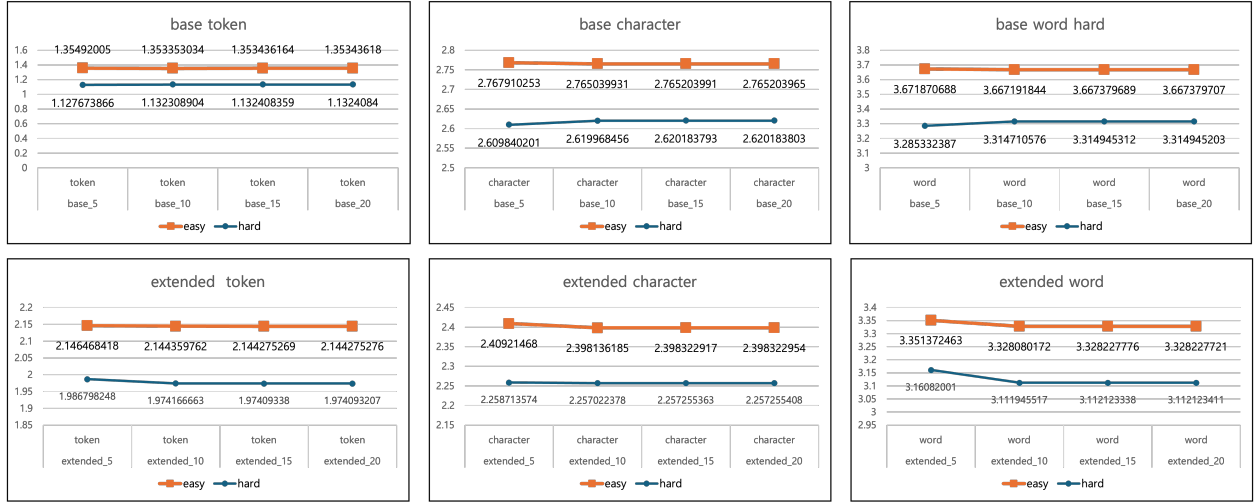


Figure 8: Normalized confidence level of each checkpoint.

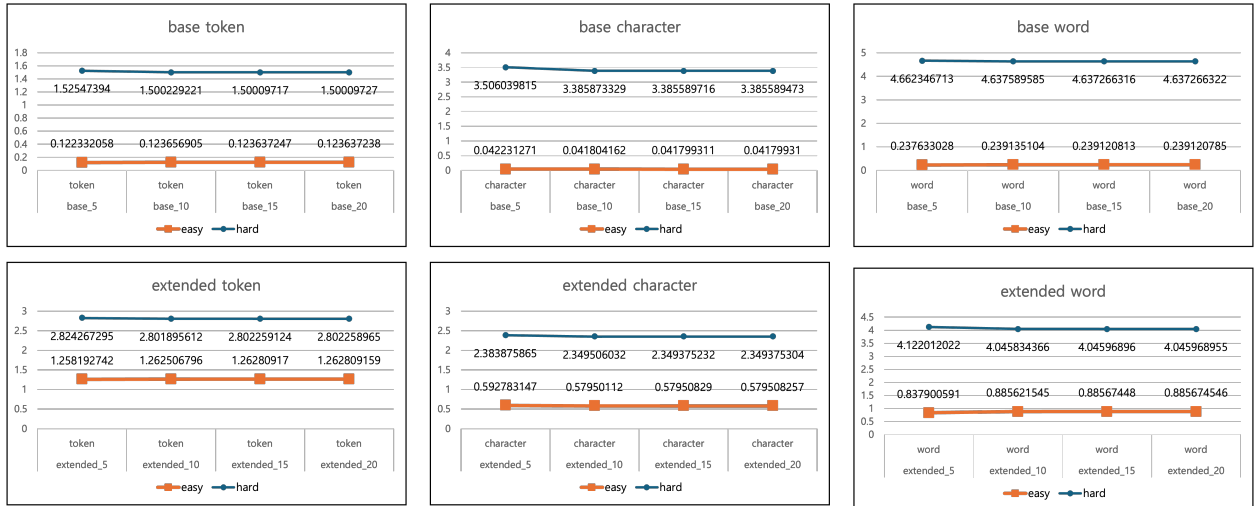


Figure 9: Cross entropy loss of of each checkpoint.