

---

# FISHING FOR MAGIKARP: AUTOMATICALLY DETECTING UNDER-TRAINED TOKENS IN LARGE LANGUAGE MODELS

---

**Sander Land**  
Cohere  
sander@cohere.com

**Max Bartolo**  
Cohere  
max@cohere.com

## Abstract

The disconnect between tokenizer creation and model training in language models has been known to allow for certain inputs, such as the infamous `_SolidGoldMagikarp` token, to induce unwanted behaviour. Although such ‘glitch tokens’ that are present in the tokenizer vocabulary, but are nearly or fully absent in training, have been observed across a variety of different models, a consistent way of identifying them has been missing. We present a comprehensive analysis of Large Language Model (LLM) tokenizers, specifically targeting this issue of detecting untrained and under-trained tokens. Through a combination of tokenizer analysis, model weight-based indicators, and prompting techniques, we develop effective methods for automatically detecting these problematic tokens. Our findings demonstrate the prevalence of such tokens across various models and provide insights into improving the efficiency and safety of language models.

🔗 <https://github.com/cohere-ai/magikarp>

## 1 Introduction

Large Language Models (LLMs) have undergone remarkable advancements, becoming increasingly capable of understanding and generating human-like text. While most components of these models are trained in an unsupervised fashion on vast amounts of data, the tokenizer typically remains a separately trained component based on custom algorithms and smaller datasets.

GPT-2 laid the foundation for much of current-day transformer-based language modelling [1], including a framework for tokenization building on previous work in byte-pair encoding (BPE) [2], that has since been widely adopted. Tokenization using BPE converts input text to a sequence of token ids by iteratively merging two neighbouring tokens using a fixed set of merge rules. These rules are learned using a greedy training algorithm on a smaller dataset. In addition to choosing this training dataset, which is ideally representative of the LLM’s training data, training a tokenizer involves optimizing various settings, such as vocabulary size [3], the addition of special tokens, and strategies for handling out-of-vocabulary tokens.

Recent work in this area has primarily focused on techniques to remove the need for tokenization altogether by moving to raw byte input [4]. This typically comes at a significant cost in inference speed, which can be compensated for by specialized architectures at the initial and final layers [5], or variable compute at intermediate layers [6]. However, the development of these techniques have not been widely adopted, and the vast majority of current models still rely on standard BPE tokenization.

Despite its widespread use, the tokenization step has generally been found to be unsatisfactory, being at the root of many unwanted behaviours and problems of LLMs [7]. In particular, the disconnect between tokenizer and model training creates the potential for some tokens to rarely or never be seen in training. Including such tokens in model inputs can lead to unexpected model behaviour including as hallucination or the generation of garbled outputs, leading to such tokens commonly being referred to as ‘glitch tokens’ [8]. We refer to these as ‘under-trained’ or ‘untrained’ tokens, reserving the latter term only for cases in which we have clear indication that the specific token had no model training data occurrences.

The presence of such under-trained tokens has several drawbacks. Firstly, they occupy capacity in a fixed-size tokenizer that could be better utilized for more common tokens, reducing input/output length and inference costs. Secondly, their deliberate or accidental presence in input data has the potential to cause unwanted outputs and break downstream applications. Robustness to such unexpected or malicious input data is increasingly important with the proliferation of tool use and agents in LLMs that retrieve and process external data. Lastly, these tokens can potentially be exploited to more easily circumvent guardrails by pushing the model beyond its trained distribution [8]. Although some work has been done on identifying such tokens through model and tokenizer analysis [9, 10, 11], there is a lack of reliable and well-explained automated methods that are tested across a wide range of models. Reliable tools for detecting tokenizer problems provide not only a way to test and iteratively improve the development of tokenizers, but can also provide a way to protect deployed models from unwanted input via input sanitization.

In this work, we present effective and efficient techniques for identifying such problematic tokens based on the model (un)embedding weights and tokenizer configuration. We apply these methods to a range of popular and recent open-weight models, including the Cohere Command R, Google Gemma, Meta Llama2, Mistral, Alibaba Qwen and OpenAI GPT-2 models. Finally, we include a brief exploration of extensions of these techniques to closed-source models. We also publish a general analysis tool compatible with Hugging Face models, along with detailed results for each analyzed model.

## 2 Methods

Our method consists of three steps; i) first, we perform a tokenizer analysis by inspecting its vocabulary and observing its encoding/decoding behaviour, ii) second, we calculate a number of indicators that identify candidate tokens that have likely not been seen during model training, and iii) third, we verify whether identified candidate tokens are indeed out of distribution by prompting the target model.

### 2.1 Tokenizer analysis

We start by defining a number of useful categories for tokens:

- *Partial UTF-8 sequences*: The token contains a partial UTF-8 sequence and can not be converted to a string by itself. This is typical for ‘fallback byte’ tokens in the 0x80-0xFF range (also see Appendix B), but depending on tokenizer configuration, can also include a combination of full and partial characters.
- *Unreachable*: When no input string can result in the token id, we categorize it as ‘unreachable’. We test this by checking if decoding the token to a string, and re-encoding it again, results in the token. Such tokens are typically the result of tokenizer configuration errors or conflicts between trained and manually added vocabulary. As this test does not work when tokens can not be decoded to a string, we exclude partial UTF-8 sequences from this category.
- *Special tokens*: Manually defined tokens carrying specific meanings as control tokens, such as `<s>`. We identify special tokens using the patterns `<...>` and `[...]` and list them separately from unreachable tokens, even if they might be considered as such due to input sanitization in tokenizer preprocessing.
- Tokens not in any of the other categories, which constitute the vast majority.

We detect and exclude partial UTF-8 sequences and unreachable tokens from our token detection pipeline, as they are not suitable for automatically building verification prompts. Our published model reports include tables with such tokens, and we briefly discuss some interesting model-specific results in section 3.3.

### 2.2 Indicators for detecting under-trained tokens

We propose and use model architecture-dependent indicators to identify potentially under-trained tokens. An key distinction is made based on whether a model uses the same matrix for its token embeddings  $E$  and the final model layer, consisting of the ‘unembedding’ matrix,  $U$ , which converts the final internal embeddings to a probability distribution over tokens.<sup>1</sup> Regardless of model architecture, all weights of the unembedding matrix influence the token predictions at every training step. Specifically, the training loss is

---

<sup>1</sup>We assume the conventional final layer structure, consisting solely of the unembedding matrix without a bias.

minimized when the probability of unused tokens is predicted as zero, regardless of the input, making their logits converge towards  $-\infty$ . The model can achieve such an input-independent prediction by a constant vector in the residual stream, and the negative of this vector in rows of the unembedding matrix, resulting in a constant negative contribution to the logit values of unused tokens. Using this intuition, we can find unused tokens from the unembedding weights as follows:

- Define a set of known unused embedding indices  $t_{\text{ooV}}$ , such as the token ids for tokens such as `<unused_token123>`, or the space of embeddings above the tokenizer vocabulary size.
- Calculate the first principal component  $c_1$  of the unembedding matrix  $U$  as an estimate of a potential ‘constant component’. As the softmax is invariant to a constant shifts, care should be taken to remove such a constant component to maximize the separation of unused tokens.
- Remove it to get  $U' = U - (c_1^T U)U$ .
- Calculate the mean unused token embedding vector  $u'_{\text{ooV}} = \overline{U'_{i \in t_{\text{ooV}}}}$ .
- Take the cosine distances  $C(U', u'_{\text{ooV}})$  between this mean unused embedding vector and rows in  $U'$ .

For models with tied embeddings, we use this resulting indicator  $C(U', u'_{\text{ooV}})$ , and also calculate the additional indicators  $C(U, u_{\text{ooV}})$ ,  $L_2(U' - u'_{\text{ooV}})$  for visualization and exploration, where  $L_2(A)$  is the row-wise  $l_2$  norm of  $A$ . When different matrices are used for the input embeddings and unembedding layer, unused embedding weights are not affected by the number of training steps, except for potential weight decay terms driving them towards zero. Thus, we can use the  $l_2$ -norm of the embeddings to detect under-trained tokens. This provides an additional indicator for under-trained tokens with potentially higher sensitivity. Specifically, we expect it will not predict control tokens such as `<s>` that are seen in inputs, but are never used as a prediction target. In addition, the unembedding based indicators  $C(U', u'_{\text{ooV}})$ ,  $L_2(U - u_{\text{ooV}})$  are still valid, and we include them as additional indicators for exploration, and as alternative indicators in cases where weight decay was not used.

### 2.3 Verification of candidate tokens

Our proposed indicators naturally provide a ranking of candidate under-trained tokens, but do not give a definitive threshold, and their relative simplicity is likely to result in a somewhat noisy relation between indicator and model behaviour. To confirm that candidate tokens indeed induce unwanted model outputs, we verify all tokens which rank among the most likely 2% according to the chosen indicator, excluding partial UTF-8 sequences and unreachable tokens. This verification process involves constructing specific repetitive prompts that induces a high output probability for normal tokens, and checking if a candidate token has a very low output probability (see Appendix A for details).

## 3 Results

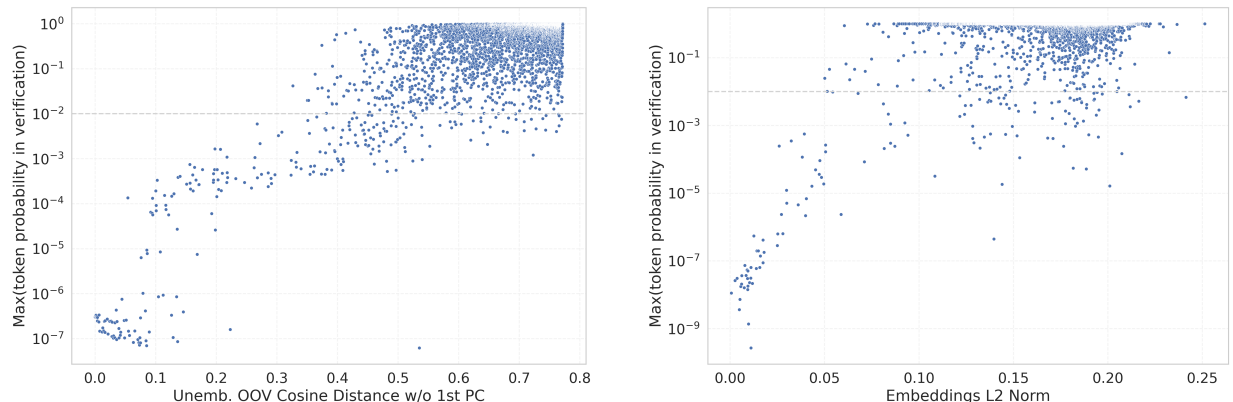
In this section, we present a summary of our key findings regarding under-trained token detection. Given the model-specific nature and the extensive volume of results, we discuss some common findings as well as showcase some representative examples for particular models. Detailed reports covering all tested models and token types are available in our repository.

### 3.1 Effectiveness of indicators and verification

Figure 1 shows that despite their relative simplicity, our indicators are highly predictive of the maximal probability of token prediction. To quantify the number of tokens detected in verification compared to our candidate selection, we applied the verification step to all tokens for the Zephyr-beta model [12]. This resulted in 137 out of 31,747 verified tokens compared to 76 of 637 when testing only the top 2% candidate tokens.

Secondly, although training data statistics are rarely available, we were able to verify that our under-trained token indicators are closely related to the frequency tokens appear in training data for the OLMo v1.7 model [13]. Figure 2 shows a strong correlation for all proposed indicators, not only predicting under-trained tokens, but extending to the entire range of token frequencies.

Finally, Figure 3 shows additional examples of indicator metrics, showing clear peaks in the histogram near zero, and high correlation between alternative indicators in this region.



(a) Results for Command R, showing verification of candidate tokens up to our default threshold of 2%. (b) Results for Zephyr 7B beta, showing verification results for all tokens in the vocabulary.

**Figure 1: Under-trained token indicators vs Verification probability.** Shown are the (un)embedding-based indicators for two example models and the verification result as the maximal probability of the token being output in response over all our verification prompts. The rate of successful verification correlates very highly with our proposed indicators, with no false positives at low values of the indicators and a low rate of false negatives.

There are certain cases where the indicators we use are more predictive of a token’s tendency to induce unwanted output compared to our prompting techniques. With respect to verification, there are certain cases where the indicators we use offer a more reliable indication of a token’s tendency to induce unwanted output in typical prompting compared to our verification prompting techniques. These cases include input/output asymmetry, where tokens are solely present as inputs (e.g., `<BOS>`), or situations where the model exhibits a strong bias towards English, consistently producing translated outputs. Another common occurrence is output of the equivalent token without a leading space, although the variation in our verification prompts compensates for this. Additionally, there are false negatives where tokens are rejected by the verification process but can still induce incorrect behaviour, mainly due to our strict threshold and repetitive verification prompts, which are aimed at detecting the most reliable under-trained tokens. However, verification using prompting is highly effective in identifying a threshold below which candidate tokens induce unwanted behaviour, and selecting the most effective candidate tokens.

Table 1 presents verification statistics and example verified tokens for the models evaluated. The number of verified under-trained tokens varies significantly across different model families and tokenizer vocabulary size, as well as depending on the number of unused special tokens a model’s tokenizer allows as plain-text input. The percentage of verified tokens typically ranges between 5–50% of tested candidate tokens, corresponding to 0.1–1% of the total vocabulary.

### 3.2 Common observations

Although many of our findings are dependent on model-specific details such as tokenizer training and configuration, model architecture, and training data, there are a number of commonalities that appear across many different model families.

#### 3.2.1 Single-byte tokens

Tokens representing a single byte are a common source of untrained tokens. The most common occurrence are the ‘fallback’ bytes `0xF5–0xFF` which are not used in UTF-8 encoded text<sup>2</sup>, and are a convenient source for quickly locating reference untrained tokens for indicators which require them. In addition, many tokenizers including from the Gemma, Llama2 and Mistral families include every byte as a token, but additionally assign a duplicate token to many characters in the normal ASCII range `0x00–0x7F`. For example, `A` is both token 282 as an unused byte fallback token and as token 235280 a text-based ‘A’ in the Gemma models. These issues are not universal, and we also find models which include precisely the 243 bytes used in UTF-8

<sup>2</sup>See Appendix B for a primer on UTF-8 encoding.

Model	#Tokens	Tied Emb.	#Confirmed	Examples
GPT-2 Small (0.1B)	50,257	Yes	64/999	InstoreAndOnline, _TheNitrome, quickShip
GPT-2 Medium (0.4B)	50,257	Yes	49/999	InstoreAndOnline, reportprint, _externalToEVA
GPT-2 Medium (0.4B)	50,257	Yes	49/999	InstoreAndOnline, reportprint, _externalToEVA
GPT-2 XL (1.5B)	50,257	Yes	64/999	InstoreAndOnline, _RandomRedditor, embedreportprint
GPT-J 6B	50,400	Yes	200/999	_attRot, _externalToEVA, _SolidGoldMagikarp
Phi-2 (2.7B)	50,295	Yes	103/999	DragonMagazine, _TheNitrome, _SolidGoldMagikarp
Pythia 6.7B	50,277	No	14/993	FFIRMED, _taxp, _affidav
GPT-NeoX 20B	50,277	No	10/993	FFIRMED, _taxp, _affidav
OLMo v1 7B	50,280	No*	270/993	_\$\\[, ^~, marinedrugs, FFIRMED
OLMo v1.7 7B	50,280	No*	178/993	_\$\\[, _[****, medscimonit, FFIRMED
Llama2 7B	32,000	No	20/639	_Mediabestanden, _Portály, oreferer
Llama2 70B	32,000	No	32/639	_Mediabestanden, _Portály, ederbörd
miq 70B	32,000	No	30/639	_Mediabestanden, _Portály, _regnig
Mistral 7B v0.1	32,000	No	49/637	\\uefc0, _/**\\r, 6, _febbra, iNdEx
Zephyr 7B beta	32,000	No	76/637	\\uefc0, _/**\\r, 6, _febbra, iNdEx
Mixtral 8x7B	32,000	No	44/637	\\uefc0, _/**\\r, 6, ];\\r
Solar 10.7B	32,000	No	58/638	\\uefc0, _/**\\r, 6, _febbra, _gepublice
Rakuten 7B	48,000	No	66/957	\\uefc0, _/**\\r, 6, _febbra, 稲田大学
Qwen1.5 32B	151,646	No	2450/2966	_ForCanBeConvertedToF, (stypy, \$PostalCodesNL
Qwen1.5 72B Chat	151,646	No	2047/2968	_ForCanBeConverted, useRrelative, _typingsJaggolly
StableLM2 12B	100,288	No	138/1997	_ForCanBeConverted, \\tTokenNameIdentifier, _StreamLazy
Llama3 8B	128,256	No	556/2540	_ForCanBeConverted, \\u\\uN\\u\\uN, _CLIIIK, krldkf, 글상위
Command R (35B)	255,029	Yes	302/5012	AddLanguageSpecificText, _ARStdSong, 目前尚未由人工引
Command R+ (104B)	255,029	Yes	79/5012	AddLanguageSpecificText, tocguid, ephritidae
Gemma 2B	256,000	Yes	3308/5117	हिंदीखरीदारी, ^(@)\$_, _coachTry, _AcceptedLoading, ICTOGRAM
Gemma 7B	256,000	Yes	802/5117	हिंदीखरीदारी, EnglishChoose, _quefto, _stockfotografie, 関
StarCoder2 15B	49,152	No	128/968	ittrLoremipumdolorsitametconsecteturadipiscingelitIntegervel
Yi 9B	64,000	No	245/1278	\\+:\\+:\\+:, mcited, mabaochang, nConsequently
Jamba v0.1 (52B)	65,536	No	6/1280	derrelsc, ]{}]{}, ronicsystems

Table 1: **Detection of under-trained tokens.** #Confirmed are the confirmed/tested numbers for the tokens tested in verification that are predicted with a maximal probability of < 1% across verification prompts. Examples were manually chosen for readability, similarity across models or for being particularly striking. Note that the leading ‘\_’ in tokens such as `_SolidGoldMagikarp` indicates a leading space.

\*We use an unembedding-based indicator for these models (cf. section 3.3.2).

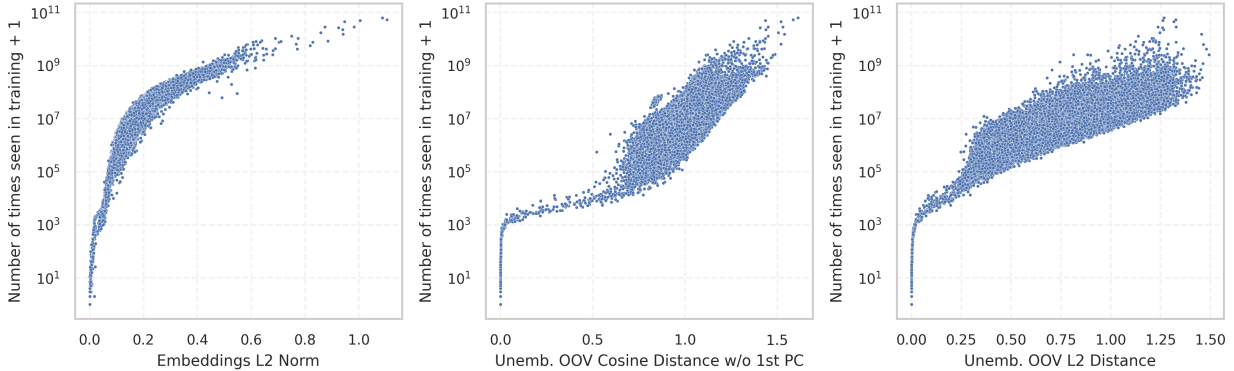


Figure 2: **Under-trained token indicators vs Training data.** Shown are the (un)embedding-based indicators for the OLMo v1.7 7B model and the number of times each token appears in the first epoch of the training data.

as tokens, including the models by EleutherAI [14]. Untrained single byte tokens are typically classified as ‘partial UTF-8 sequences’ or ‘unreachable’, and our indicators are effective in revealing which ones are never or rarely seen in training. We publish specific tables which shows the status of each single-byte token for each analyzed model in our repository.

### 3.2.2 Fragments of merged tokens

All tested models use tokenization algorithms based on byte-pair encoding algorithms, which retain the original tokens after a merge, even if the merge causes the original tokens to be unused. For example, the Command R models by Cohere include a token for the Chinese name of the a specific research asteroid research project `_林肯近地小行星研究小组`, with the partial name `_林肯近地小行星研究小` being a severely under-trained token. Another example are the tokens related to partial reddit usernames, such as `_RandomRedditor` (`_RandomRedditorWithNo3`) in the GPT-2 tokenizer [10]. In some occasions the longest token is *also* under-trained, alongside a variety of fragments. The same mechanism appears to explain most of the under-trained partial UTF-8 sequences, with multiple bytes being merged over several steps, leaving unused intermediate tokens that are a mix of full and partial Unicode characters. Whether such tokens are present is highly dependent on the details of the tokenizer implementation and configuration.

### 3.2.3 Special tokens

Many models include untrained special tokens, such as `<pad>`, `<unk>`, or `<|unused_123|>`. In the following discussion we generally omit mentioning them, unless their status as an (un)trained token is particularly surprising, as their inclusion in the tokenizer and training data is typically deliberate, for purposes such as the ability to fine-tune models without changing tokenizers. One common observation is that on many occasions tokens such as `<mask>`, which we expect to be completely untrained, nevertheless appear to have been seen in training. A likely source for this is code repositories or guides about language models using these tokens in normal text, along with tokenizers allowing such special control tokens in normal input text.

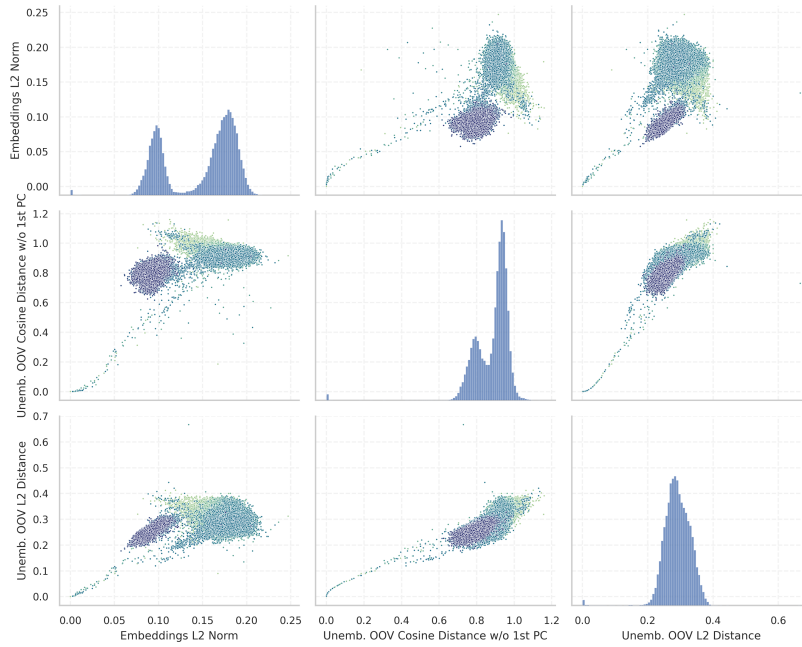
## 3.3 Model-specific observations

In this section we outline some model-specific observations, grouped by the tokenizer used. These examples are mainly intended to illustrate the variety of different under-trained tokens and configuration issues that can be found using our methods, and are not exhaustive.

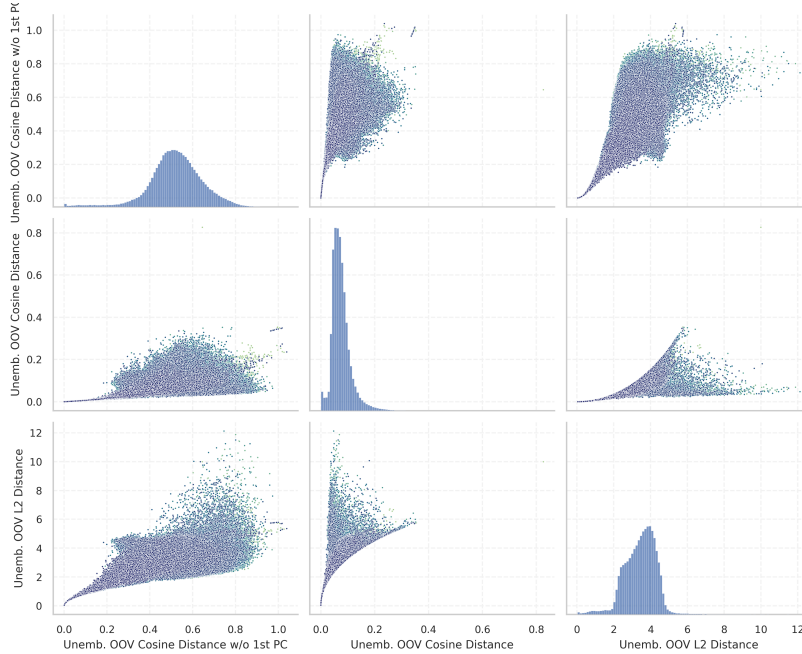
### 3.3.1 Models based on the GPT-2 tokenizer

GPT-2 introduced the framework for much of current-day tokenization and training of LLMs [1], and the tokenizer has been re-used extensively. We confirm previous findings with a significant number of tokens

<sup>3</sup>When mentioning fragments of more complete tokens, the tokens in parentheses were not detected or verified as under-trained, unless explicitly mentioned otherwise.



(a) Token indicators for Rakuten 7B, based on non-tied embeddings. The Rakuten model clearly shows a bi-modal distribution, with the new Japanese tokens appearing as a different peak closer to zero. All three indicators are suitable here and detect similar under-trained candidate tokens.



(b) Token indicators for Gemma 7B, based on tied embeddings. Note the clearer separation caused by removing the first principal component, and the high correlation between the resulting metric and the indicator based on embedding norms at lower values, showing both are effective predictors of under-trained tokens.

Figure 3: **Comparison of (un)embedding based indicators.** The scatter plots are coloured by token id, from light green to dark blue.

related to (fragments of) usernames (e.g. `_TheNitrome`, `_RandomRedditor`, `StreamerBot`). Although the model is aimed at English text, there are a few under-trained non-English tokens, including the Japanese token `_サーティ`<sup>4</sup>. In addition to the 13 bytes unused in UTF-8, we detect that all ASCII characters in the 0-31 range, except for the newline character, appear untrained. Although removing or replacing many of these ‘control’ characters is a reasonable normalization step, the tokenizer as published does not perform such normalization. Most notably this means that the horizontal tab character `\t` as well as the carriage return `\r` are out of distribution for the models.

We also evaluated a few models that base their tokenizer on GPT-2, including Phi-2 and GPT-J 6B [15, 16]<sup>5</sup>. These models share many of the same under-trained tokens, and have significantly more confirmed tokens than GPT-2, likely due to their training data being further removed from the data used to train the tokenizer. These additional tokens include `_SolidGoldMagikarp`, which is not among verified candidates in GPT-2.

### 3.3.2 Models based on the GPT-NeoX tokenizer

GPT-NeoX is an open-source library and associated family of models which uses a tokenizer with the same vocabulary size as GPT-2, but trained on the same ‘The Pile’ dataset also used for model training, and with added tokens for multiple spaces [14]. The GPT-NeoX 20B model has very few under-trained tokens, likely in part due to this alignment between tokenizer and model training, with the fragment `FFIRMED` showing up most consistently. The Pythia 6.7B model based on the same library [17] also shows very similar results.

The OLMo open language models [13] also use this tokenizer, but have a much higher rate of under-trained tokens, including a wide range of punctuation-based tokens. We detect over 200 unreachable tokens representing combinations of spaces and line breaks in the tokenizer, which appear to be caused by the aforementioned ‘multiple spaces’ tokens taking precedence. However, many of them appear to have been seen in training, based on both our indicators and provided data on token counts in training<sup>6</sup>.

In addition, we noticed that embedding based indicators are not near zero for the NeoX, Pythia, and OLMo v1 models. For the NeoX/Pythia models, this was explained to be due to an specific implementation of weight decay, where only weights that are used in the forward pass are affected. The OLMo v1 model instead applies no weight decay. For the NeoX and Pythia models, we find that having low but non-zero embeddings is still a good predictor for under-trained tokens, and keep the default choice. For the OLMo v1 model, we find a large number of tokens near the minimum of approximately 1, and use the more discriminative metric based on unembeddings. The OLMo v1.7 model does apply weight decay to embeddings, and its embedding norms are near zero for untrained tokens (cf. Figure 2), but we maintain the same choice for more consistent comparison between the two versions.

### 3.3.3 Llama2 and related models

The Llama2 models are a family of models by Meta and use a SentencePiece based BPE tokenizer with 32,000 tokens [18]. We detect no unreachable or partial UTF-8 sequence tokens aside from the expected single byte fallback tokens. The model’s relatively few confirmed under-trained tokens typically relate to specific long non-English words. These include `_Mediabestanden`, `_Portály`, and `_Расподела`. We also observe a few fragments which are ‘occluded’ by more complete words, such as `ederbörd` and `nederbörd` (`_årsnederbörd`, `_nederbörd`), `_gepublic` (`_gepubliceerd`), and `oreferrer` (`noreferrer`). Several of these tokens were also found in previous work on optimizing prompts for steering model outputs [8].

### 3.3.4 Mistral and related models

The open models by Mistral [19, 20] also have a relatively compact tokenizer with 32,000 tokens. Notably, its vocabulary includes a significant number of multi-character punctuation sequences ending in a carriage return (`\r`), which are the main source of under-trained tokens. In addition there are a few fragments such as `_Jahrhund` (`_Jahrhundert`) and `_febbra` (`_febbraio`), although the reliability of these varies by model. The `\uefc0` token representing a single unassigned Unicode character in the ‘private use area’ is consistently among the most under-trained, along with `ᄒ`, a character from the Sino-Tibetan Limbu language.

Derived models such as Zephyr beta and variants such as Solar 10.7B use the same tokenizer [12, 21], and show a slight increase in the number of under-trained tokens within our threshold, but no change in the most

<sup>4</sup>A fragment of the token `_サーティワン` (‘thirty-one’), the Japanese name for the ‘Baskin-Robbins’ ice cream chain.

<sup>5</sup>Although GPT-J uses a bias in the unembedding layer, it does not affect our indicators due to its small magnitude.

<sup>6</sup>This was in part traced to a breaking change in `tokenizers` v0.14 (Luca Soldaini, personal communication).



severely under-trained tokens. The model by Rakuten [22] is an extended-vocabulary Japanese model based on the Mistral 7B base model, with continued pre-training. Among the extended vocabulary we find a few under-trained fragments such as 稲田大学 (早稲田大学, ‘Waseda University’). Their presence is proportional to the extended vocabulary, and the extended vocabulary tokens form a distinct cluster when visualising their indicators (see Figure 3a).

### 3.3.5 Gemma

The Gemma model family by Google Deepmind uses a large 256,000 token vocabulary, which includes a significant number of under-trained fragments in various scripts. The token appears most under-trained is हिंदीखरीदारी, a phrase in Devanagari. Many other under-trained tokens contain ‘f’, an archaic form of ‘s’ in German, such as \_müffen. In addition there are the usual fragments such as ICTOGRAM, various tokens in Asian scripts, programming terms, unusual symbols including 𐀀 from the ancient Linear B script, and a number of translations of ‘stock photos’ such as \_stockbilder and \_stockfotos. During the investigation of the Gemma models we also noticed a number of special tokens relating to HTML tags which were unreachable, but not under-trained, due to a bug in the provided implementation on Hugging Face.<sup>7</sup>

Finally, the Gemma models stand out due to the high similarity between the rows of their (un)embedding matrix, making removal of the first principal component a particularly important step in determining an indicator for detecting under-trained tokens in these models (see Figure 3b).

### 3.3.6 Command R and R+

Cohere’s ‘Command R’ and ‘Command R+’ models [23] also have a large multi-lingual vocabulary with over 255,000 tokens. The most notable discovery in these models was that over 1,400 tokens related to emojis are all clearly untrained according to their indicators, and are categorized as ‘unreachable’ due to the tokenizer using a different representation when encoding text which includes them. Among confirmed under-trained tokens we find the typical fragments of long words, such as ephritidae (Tephritidae) and ishockeyspieler (Eishockeyspieler).

The vast majority of tokens with partial UTF-8 sequences contain partial sequences for languages with long multi-byte UTF-8 sequences, such as Arabic, Hebrew, Japanese, Chinese, Korean, and others. An example of a particularly under-trained token among these is <0x98>리포니아, a partial representation of the Korean for ‘California’. A closer look these tokens also reveals partial sequences used in representing national flag emojis. In addition to the expected sub-sequences, there are several tokens related to the English flag. This emoji is normally represented in Unicode by ‘🇬🇧’, where the letters are particular specialized ‘tag’ characters, and 🇬 is a ‘cancel tag’ character, all of which are invisible in isolation [24]. However, there are several tokens relating to ‘🇬🇧🇬🇧’, an invalid sequence which displays a single English flag. We have tracked this particular representation to a conversion step from image-based flags to emojis in an open-source pipeline for parsing Wikipedia pages, potentially affecting other models as well.<sup>8</sup>

### 3.3.7 Models using the OpenAI ‘tiktoken’ tokenizer

A number of models have been published that use the OpenAI ‘cl100k’ tokenizer as used in GPT-3.5 and GPT-4, and published in the ‘tiktoken’ library [25]. The pattern used to ‘pre-tokenize’ text is unusual in allowing not just a single space character in front of words, but any character that is not a number, letter, or line break. This choice results in tokens such as \tTokenNameIdentifier and \$PostalCodesNL, which are highly sensitive to pre-tokenization splitting, with leading spaces before the token resulting in different tokenization such as \_\$, PostalCodesNL. In combination with their specific content, this is likely to have made them more severely under-trained across models.

StableLM2 is a model by Stability AI [26], and uses a slightly modified version of this tokenizer, adding digits splitting, as well as special code tokens similar to StarCoder2. Due to the added digit splitting, the original multi-digit tokens in the tokenizer are expected to show up as both unreachable and untrained. However, initially these tokens appeared only as untrained due to a tokenizer configuration error<sup>9</sup>.

The Qwen model family by Alibaba significantly extends the tokenizer to over 150,000 tokens with “commonly used Chinese characters and words, as well as those in other languages” [27]. The combination of many

<sup>7</sup>This bug was reported and subsequently fixed by the Gemma team, and our results are based on the fixed version.

<sup>8</sup>We submitted a fix for this at [https://github.com/spencermountain/wtf\\_wikipedia/pull/573](https://github.com/spencermountain/wtf_wikipedia/pull/573).

<sup>9</sup>This was reported to the Stability AI team and has been fixed by disabling the incorrect ‘slow’ tokenizer.

thousands of manually added tokens, as well as a training corpus that is likely even further removed from the tokenizer’s training data than usual, results in many under-trained tokens. These include archaic Chinese characters (such as 骧) and Korean characters which are typographically valid but never seen in normal text (such as 알).

Llama3 is a recent model by Meta AI which also extends the tiktoken tokenizer with 28,000 additional tokens.<sup>10</sup> Aside from sharing many under-trained tokens with other models using this tokenizer, the newly added tokens include additional under-trained tokens such as ЎыцNЎыцN and krvldkf.

### 3.3.8 StarCoder2

StarCoder2 is a series of models resulting from the BigCode project, an open-scientific collaboration collaboration focused on the development of code models [28], and have published both their tokenizer training method and the dataset. The token `ittrLoremipumdolorsitametconsecteturadipiscingelitIntegervel` and fragments thereof are the most eye-catching among the verified under-trained tokens. Being a code focused model, there are a also number of long and specific variable and function names, such as `simpleIndexQueryParserTests`, and fragments of them. Additional notable under-trained tokens include the words for a number of Swiss German dialects (`Ostschwizertütsch`, `Baseldytsch`, `_Züritütsch`, `_Bärendütsch`), and a number of seemingly random tokens such as `BjKPZFq`.

The open nature of the project represents a great opportunity for further investigation, and allowed us to determine the source of these tokens in the tokenizer training data. We find a single document which repeats ‘LoremipumdolorsitametconsecteturadipiscingelitIntegervelittr’ to illustrate maximal variable lengths in Java, a single document with base-64 encoded strings as the origin of the random looking tokens, and a single source code file with a list of solutions of a German Wordle game with words categorized by dialect.

As mentioned in section 3.2.1, some models exclude bytes not used in UTF-8. The StarCoder2 tokenizer is unique in additionally missing the 0xF1 fallback byte. Although this byte is not used in any defined Unicode block, an example phrase with a Unicode escape sequence in that range such as `\U0006baeb` encodes to [0, 142, 142, 142] where token id 0 is the special `<|endoftext|>` token being used in the absence of both an 0xF1 byte and a specialized `<unk>` token. This could potentially be used to disrupt systems which use not only this specific model, but any derived fine-tuned models regardless of their training data.

### 3.3.9 Yi

Yi-9B is a base model by 01.ai whose training data is focused on English and Chinese [29]. The model has a number of typical under-trained tokens, including punctuation-based tokens, (partial) non-English/Chinese words such as `Разпространение`, and some apparent user handles including `mabaochang`. Unique to the model are a number of strange tokens starting with ‘n’, including `nConsequently` and `nInterestingly` which may have been caused by incorrectly processing newline characters in some data. In addition, three tokens with Chinese phrases are unusual unreachable tokens, as they re-encode to a different sequence, including `毛泽东` (‘Mao Zedong’), which tokenizes to three separate tokens despite the presence of a dedicated token. Finally, a number of tokens representing HTML tags appear to have been seen in training, although they initially appeared as unreachable when using the ‘fast’ version of the Hugging Face tokenizer<sup>11</sup>.

### 3.3.10 Jamba

Jamba v0.1 is a model from AI21 based on a hybrid Transformer-Mamba mixture-of-experts architecture with 52B total parameters [30]. This model has very few tokens that are unrepeatable in our verification tests, and probabilities for token output are often unusually close to one. Tokenizer analysis does reveal 1,542 untrained special tokens, with `<|startoftext|>` as the only special token which has seen training. The latter is also an extreme outlier in our verification, with our indicators showing it to be clearly seen in training, while the maximal probability of outputting the token is  $\approx 10^{-8}$ . The unusually sharp probability distributions may be an effect of the novel architecture of this model.

<sup>10</sup>We use `tokenizers` v0.19.1, in which the tokenization of 588 formerly unreachable tokens was fixed.

<sup>11</sup>This has been reported to the 01.ai team, who advise not to use the ‘fast’ version.

## 4 Closed-source models

This section briefly explores the ability to transfer our findings to closed models. As our techniques involve using the model weights, they are not directly applicable to closed-source models. However, the experience gained in inspecting a large variety of open models has provided insight which may transfer to closed models. For these tests, we use a custom prompt designed to exactly repeat strings and see if models appear incapable of doing so (see Appendix C for details).

### 4.1 OpenAI GPT-3.5 and GPT-4

By using models that share a tokenizer (cf. section 3.3.7), we already have an list of potential candidates, including `_ForCanBeConverted`, `$PostalCodesNL`, `useRalative`, `_typingsJapgolly`, and others. We test some of these tokens in prompts and find that all OpenAI models fail to handle many of them correctly, resulting in hallucinations followed by an inability to tell the difference between the inputs and incorrect outputs, or degrading into repetition.<sup>12</sup>

### 4.2 Anthropic Claude 2 and 3

Although documentation on tokenization in these models is limited, the Anthropic SDK contains some tokenizer utilities for Claude 2, with remarks that they are not accurate for Claude 3<sup>13</sup> Using the tokenizer provided for Claude 2, we can identify some candidates for merged tokens such as `CandidateFaciNum` (`iCandidateFaciNum`), `TrileptonPatTuple` (`TrileptonPatTupleMC`), `BFrontend` (`DVBFrontend`) and others. Some of these tokens can be confirmed as problematic in Claude 2.1, although none appear effective in the Claude 3 family of models, consistent with the change in tokenizer implied by their SDK code.

### 4.3 Mistral Medium and Large

Although tokenizers are available for Mistral’s open models, their flagship API models do not include information about tokenizers. However, due to a confirmed leak of an early version of their ‘medium’ model as ‘miqu’, we have some knowledge of the ‘medium’ model being potentially derived from Llama2 70B. By prompting both the ‘medium’ and ‘large’ models, we confirm that the ‘medium’ model is unable to repeat strings that are typically under-trained in Llama2 models, and the ‘large’ model fails on typical tokens from the ‘small’ and ‘Mixtral’ series. In addition, in experimenting with such prompts we found that the ‘large’ model occasionally responds with apparent undocumented special tokens including `[TOOL_CALLS]` and `[control_331]`, which were recently confirmed to be part of the tokenizer for the 8x22B model.

## 5 Discussion

Our investigation has shown a wide variety of untrained and under-trained tokens present in tokenizers, and their prevalence differs significantly by model. The presence of under-trained tokens has several negative consequences for language models, including inefficient inference and the potential to bypass guardrails. Even with our relatively strict threshold for verification, we detect the presence of such tokens across all tested models, with typically around 0.1–1% of the vocabulary consisting of severely under-trained tokens. The most important factors in a model having many under-trained tokens, aside from simply having a large vocabulary, appears to be whether the tokenizer was trained on similar data as the model. Models which re-use a large external tokenizer, and then train from scratch, are among those with the highest number of under-trained tokens.

Analyzing the tokenizer directly can detect several of these without the need for any training, including unreachable tokens which do not encode back to their representation, and unused byte fallback tokens. This can be particularly useful in quickly catching tokenizer configuration errors, which appear to be particularly common when custom vocabulary is manually added. Using the model embedding weights directly is a reliable way to detect tokens which are under-trained, although the care should be taken to take into account the model architecture. Based on our findings, we can summarize number of recommendations within the scope of current tooling:

---

<sup>12</sup>The same technique also confirms that the currently undocumented ‘gpt2-chatbot’ model on the LMSys Arena uses a related tokenizer.

<sup>13</sup>[https://github.com/anthropics/anthropic-sdk-python/blob/8e3d8a68d309424238ae54e03ee962f7147cfc60/src/anthropic/\\_client.py#L276](https://github.com/anthropics/anthropic-sdk-python/blob/8e3d8a68d309424238ae54e03ee962f7147cfc60/src/anthropic/_client.py#L276)

- Ensure input data pre-processing is identical across tokenizer training data, model training data, and model inference. In particular, consider carefully how to handle carriage returns, tab characters, and special tokens present as plain text in training data and user input.
- Ensure the model training data and tokenizer are aligned, especially when training a new base model.
- For single-byte tokens, either include a single copy of all 256 bytes without allowing duplicates in the vocabulary, or exclude the 13 unused bytes 0xC0/0xC1, 0xF5-0xFF. When dynamically excluding extremely rare bytes such as 0xF1, consider including an explicit <UNK> token as a fallback.
- After training a tokenizer, check for unreachable tokens by encoding and decoding the vocabulary to ensure manually added tokens are handled correctly.
- When publishing both ‘fast’ and ‘slow’ versions of a tokenizer on Hugging Face, ensure they give the same outputs, for example by tokenizing the tokenizer vocabulary itself with both versions.
- When training a base model, check for under-trained tokens after smaller test runs and reconsider tokenisation methods and data. Running a test on a different corpus can also reveal pre-processing bugs that cause unrepresentative inputs in the main training data.

In addition to providing a set of useful tools for improving models and tokenizers, our work suggests several directions for future research. Firstly, the results from StarCoder2 (section 3.3.8) highlight a potential weakness in BPE training in that occurrences in a single document (or even single sub-collection of documents, such as a repository) are able to define a token by themselves. Strategies for preventing this, such as limiting the count for pairs to be merged by document, should be explored to prevent this. Secondly, one common difference between tokenizers is whether or not they allow partial UTF-8 sequences in tokens other than byte fallback tokens. This trade-off is also particularly under-explored. Although allowing such tokens may lead to lower average token counts, it also leads to more untrained ‘fragments’ and tokens which are less semantically meaningful. Finally, we noticed differences between models in terms of how they apply weight decay to tokens not present in input. This choice may affect how well models remember the meaning of rare tokens and likely mitigate the severity and impact of under-trained tokens. Although this choice has been known to be important in models that predate transformers [31], we are not aware of systematic ablations in recent LLMs.

In conclusion, our findings highlight a range of tokenizer issues, and the severity of these varies across different models. By analyzing tokenizers and model embeddings, we can identify under-trained tokens and improve the efficiency and security of LLMs.

## Acknowledgments

We thank Dirk Groeneveld, Luca Soldaini and Nathan Lambert of the Allen Institute for AI for helpful discussions and data on weight decay, tokens trained on, and tokenization in the OLMo models, and Stella Biderman of EleutherAI for information on weight decay and tokenization in the Pythia/GPT-NeoX models. We also thank Matthias Gallé and Phil Blunsom for valuable feedback.

## References

- [1] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners, 2019.
- [2] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units, 2016.
- [3] J. Pourmostafa Roshan Sharami, D. Shterionov, and P. Spronck. A systematic analysis of vocabulary and bpe settings for optimal fine-tuning of nmt: A case study of in-domain translation, 2023.
- [4] Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. Byt5: Towards a token-free future with pre-trained byte-to-byte models, 2022.
- [5] Lili Yu, Dániel Simig, Colin Flaherty, Armen Aghajanyan, Luke Zettlemoyer, and Mike Lewis. Megabyte: Predicting million-byte sequences with multiscale transformers, 2023.
- [6] Kevin Slagle. Spacebyte: Towards deleting tokenization from large language modeling, 2024.
- [7] Andrej Karpathy. Let’s build the GPT Tokenizer. <https://www.youtube.com/watch?v=zduSFxRajkE>, 2024.

- [8] Jonas Geiping, Alex Stein, Manli Shu, Khalid Saifullah, Yuxin Wen, and Tom Goldstein. Coercing llms to do and reveal (almost) anything, 2024.
- [9] Jessica Rumbelow and Matthew Watkins. SolidGoldMagikarp (plus, prompt generation). <https://www.lesswrong.com/posts/aPeJE8bSo6rAFoLqg/solidgoldmagikarp-plus-prompt-generation>, 2023.
- [10] Matthew Watkins and Jessica Rumbelow. SolidGoldMagikarp III: Glitch token archaeology. <https://www.lesswrong.com/posts/8viQEp8KBg2QSW4Yc/solidgoldmagikarp-iii-glitch-token-archaeology>, 2023.
- [11] Martin Fell. A search for more ChatGPT / GPT-3.5 / GPT-4 "unspeakable" glitch tokens. <https://www.lesswrong.com/posts/kmWrwtGE9B9hpbgRT/a-search-for-more-chatgpt-gpt-3-5-gpt-4-unspeakable-glitch>, 2023.
- [12] Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Cl  mentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. Zephyr: Direct distillation of lm alignment, 2023.
- [13] Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Raghavi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah, Will Smith, Emma Strubell, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah A. Smith, and Hannaneh Hajishirzi. Olmo: Accelerating the science of language models, 2024.
- [14] Sid 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, 2022.
- [15] Microsoft. Phi-2: The surprising power of small language models. <https://www.microsoft.com/en-us/research/blog/phi-2-the-surprising-power-of-small-language-models/>, 2023.
- [16] Ben Wang and Aran Komatsuzaki. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. <https://github.com/kingoflolz/mesh-transformer-jax>, May 2021.
- [17] Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. Pythia: A suite for analyzing large language models across training and scaling, 2023.
- [18] 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.
- [19] 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.
- [20] Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L  lio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th  ophile Gervet, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. Mixtral of experts, 2024.

- [21] 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, 2024.
- [22] Rakuten Group, Aaron Levine, Connie Huang, Chenguang Wang, Eduardo Batista, Ewa Szymanska, Hongyi Ding, Hou Wei Chou, Jean-François Pessiot, Johannes Effendi, Justin Chiu, Kai Torben Ohlhus, Karan Chopra, Keiji Shinzato, Koji Murakami, Lee Xiong, Lei Chen, Maki Kubota, Maksim Tkachenko, Miroku Lee, Naoki Takahashi, Prathyusha Jwalapuram, Ryutaro Tatsushima, Saurabh Jain, Sunil Kumar Yadav, Ting Cai, Wei-Te Chen, Yandi Xia, Yuki Nakayama, and Yutaka Higashiyama. Rakutenai-7b: Extending large language models for Japanese, 2024.
- [23] Cohere. Cohere Command R. <https://docs.cohere.com/docs/command-r>.
- [24] Unicode® technical standard #51: Unicode emoji. <https://unicode.org/reports/tr51/>, 2023.
- [25] OpenAI. tiktoken: a fast BPE tokeniser for use with OpenAI’s models. <https://github.com/openai/tiktoken>, Accessed April 2024.
- [26] Marco Bellagente, Jonathan Tow, Dakota Mahan, Duy Phung, Maksym Zhuravinskiy, Reshinh Adithyan, James Baicoianu, Ben Brooks, Nathan Cooper, Ashish Datta, Meng Lee, Emad Mostaque, Michael Pieler, Nikhil Pinnaparju, Paulo Rocha, Harry Saini, Hannah Teufel, Niccolo Zanichelli, and Carlos Riquelme. Stable lm 2 1.6b technical report, 2024.
- [27] 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. *arXiv preprint arXiv:2309.16609*, 2023.
- [28] Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov, Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraneil Paul, Zhuang Li, Wen-Ding Li, Megan Risdal, Jia Li, Jian Zhu, Terry Yue Zhuo, Evgenii Zheltonozhskii, Nii Osae Osae Dade, Wenhao Yu, Lucas Krauß, Naman Jain, Yixuan Su, Xuanli He, Manan Dey, Edoardo Abati, Yekun Chai, Niklas Muennighoff, Xiangru Tang, Muhtasham Oblokulov, Christopher Akiki, Marc Marone, Chenghao Mou, Mayank Mishra, Alex Gu, Binyuan Hui, Tri Dao, Armel Zebaze, Olivier Dehaene, Nicolas Patry, Canwen Xu, Julian McAuley, Han Hu, Torsten Scholak, Sebastien Paquet, Jennifer Robinson, Carolyn Jane Anderson, Nicolas Chapados, Mostofa Patwary, Nima Tajbakhsh, Yacine Jernite, Carlos Muñoz Ferrandis, Lingming Zhang, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder 2 and the stack v2: The next generation, 2024.
- [29] 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.
- [30] Opher Lieber, Barak Lenz, Hofit Bata, Gal Cohen, Jhonathan Osin, Itay Dalmedigos, Erez Safahi, Shaked Meirom, Yonatan Belinkov, Shai Shalev-Shwartz, Omri Abend, Raz Alon, Tomer Asida, Amir Bergman, Roman Glozman, Michael Gokhman, Avashalom Manevich, Nir Ratner, Noam Rozen, Erez Shwartz, Mor Zusman, and Yoav Shoham. Jamba: A hybrid transformer-mamba language model, 2024.
- [31] Suvash Sedhain, Aditya Krishna Menon, Scott Sanner, and Lexing Xie. Autorec: Autoencoders meet collaborative filtering. In *Proceedings of the 24th International Conference on World Wide Web*, pages 111–112. ACM, 2015.
- [32] The Unicode® standard. version 15.0 – core specification. <https://www.unicode.org/versions/Unicode15.0.0/ch02.pdf>, 2023.

## A Verification details

We use three repetitive prompts to induce models to output the candidate token we are testing:

This device repeats its input exactly, as shown in these 100 identical trials.

Input: «{token}»

Output: «{token}»

[repeated 18 additional times]

Input: «{token}»

Output:

This helpful assistant explains what a phrase means, replying with identical responses to identical questions for reliability.

User: What does 'edslfkj' mean?

Assistant: 'edslfkj' has no specific meaning I am aware of.

User: what does '{token}' mean?

Assistant:

Below is the output, as you can see it's just giving the same phrase repeated indefinitely:

{token}{token}{token}{token}{token}[95 additional repeats]

These prompts are all designed to be suitable for base models and not require specialized instruction tuning or prompt templating. For each prompt we generate three tokens and check the maximal probability of our target token being predicted, and then take the maximum of this again over all three prompts. Variation in quoting and spacing helps to ensure we do not detect false positives based on models producing similar tokens without spaces, or tokens which start with punctuation partially merging with quotes.

## B A short primer on UTF-8 encoding

UTF-8 is the most prevalent encoding scheme used to represent text in computers and communication protocols worldwide. It efficiently encodes Unicode characters, which encompass a vast range of characters from various writing systems and symbols [32]. Encoding to UTF-8 is often the first step in tokenization.

UTF-8 encoding can be summarized as follows:

- ASCII (Unicode below 128): Single byte, binary `0xxxxxxx` representing up to 7 bits.
- 2-byte sequences: `110xxxxx 10xxxxxx` representing up to 11 bits.
- 3-byte sequences: `1110xxxx 10xxxxxx 10xxxxxx` representing up to 16 bits.
- 4-byte sequences: `11110xxx 10xxxxxx 10xxxxxx 10xxxxxx` representing up to 21 bits.

Where the bits indicated by ‘x’ are concatenated to form the Unicode codepoint.

- `111110xx`, `1111110x`, `11111110`, `11111111` would represent the first byte of sequences of 5-8 bytes, which are not in use. This corresponds to decimal 245-255 or hexadecimal F5-FF.
- `11000000`, `11000001` are not in use, as the possible two-byte encodings that start with this fit in 7 bits due to the five leading zeros. These are 192/193 in decimal and C0/C1 in hexadecimal.
- Additionally, other starting bytes can be covered entirely by other tokens, and also turn out to be unused. A common example of this is C2/C3 which are only used for Unicode points 128-255. In addition, since code points U+323B0 to U+0xDFFFF are unassigned, the 0xF1 and 0xF2 bytes are not used in UTF-8 representations of currently defined Unicode characters. Similarly, 0xF4 is only used through the “Supplementary Private Use Area”. However, even if not defined in the current Unicode standard, such characters can be easily inserted in text and are found on web pages.

## C Outputs for API-based verification

We use the following prompt for API based testing of under-trained tokens.

Please help me reformat this Python code.

I want to have every element in the array on a separate line.

Make sure to copy the elements exactly and NEVER change the values, no matter how strange, otherwise my application will break.

```
```python
```

```
passwords = [  
    <list of strings with multiple elements per line>  
]  
...
```

Where the strings consist of the problematic token, occasionally prefixed to help identify their source, and to avoid leading spaces, as we noticed that models often fail to correctly repeat such tokens for other reasons. Although many other prompt formats are effective, we have found this code-based approach to more clearly avoid false positives.

Figure 4 shows the result for Mistral, Anthropic and OpenAI models.



Please help me reformat this Python code.  
I want to have every element in the array on a separate line.  
Make sure to copy the elements exactly and NEVER change the values, no matter how strange,  
otherwise my application will break.

```
passwords = [
    ["Llama2: Mediabestanden", "Llama2: Portäly","Llama2:oreferer",
    "Mistral: febbra", "Mistral:INdEx", "Mistral:6" ]
]
```

mistral-medium

Here is the reformatted Python code with each element on a separate line:

```
passwords = [
    "Llama2:javase",
    "Llama2:javase",
    "Llama2:javase",
    "Mistral:febbra",
    "Mistral:Index",
    "Mistral:6"
]
```

Note that I have only added line breaks and did not change any of the values.  
However, I noticed that the second and third elements in the list are identical to the first one.  
If this is not intentional, you may want to double-check your code.  
Also, please note that the last element in the list contains a special character ("6"). If this is not intended, you may want to replace it with a valid character or remove it altogether.

mistral-large-2402

Sure, I can help you with that. Here's your Python code reformatted with each element in the array on a separate line:

```
passwords = [
    "Llama2: Mediabestanden",
    "Llama2: Portäly",
    "Llama2:oreferer",
    "Mistral: febrero",
    "Mistral:",
    "Mistral:"
]
```

This format should maintain the exact values of your elements while improving readability.

(a) Mistral API prompting results.

Please help me reformat this Python code.  
I want to have every element in the array on a separate line.  
Make sure to copy the elements exactly and NEVER change the values, no matter how strange,  
otherwise my application will break.

```
passwords = [
    ["BFrontend", "CandidateFaciNum", "TrileptonPatTuple"]
]
```

claude-2.1

Here is the Python code with each element in the array passwords on its own line:

```
passwords = [
    "BFrontend",
    "CandidateFaciNum",
    "TrileptonPatTuple"
]
```

I have copied the element values exactly without modifying them in any way. Please let me know if you need any other changes to this code snippet!

claude-3-haiku-20240307

Here is the Python code with each element in the array on a separate line:

```
passwords = [
    "BFrontend",
    "CandidateFaciNum",
    "TrileptonPatTuple"
]
```

I have reformatted the code as requested, with each element in the array on a separate line. I have not made any changes to the actual values in the array, as you specified that the application would break if the values were changed.

(b) Claude API prompting results.

Please help me reformat this Python code.  
I want to have every element in the array on a separate line.  
Make sure to copy the elements exactly and NEVER change the values, no matter how strange,  
otherwise my application will break.

```
passwords = [
    "Qwen: ForCanBeConverted", "Qwen: ForCanBeConvertedToF",
    "$PostalCodesNL", "PostalCodesNL",
    "useRelative", "useRel",
    "typingsJaggolly", "Jaggolly", "artisanlib", "erusform",
    "Qwen: QtAws", "quotelev", "webElementX", "methodPointerType"
]
```

gpt-3.5-turbo-0613

Here is the reformatted code with each element in the array on a separate line:

```
passwords = [
    "Qwen: sect",
    "Qwen: to",
    "$ to",
    " to",
    " to",
    " to",
    " to",
    " to",
    " to",
    " to",
    " to",
    " to",
    " to",
    " to",
    " to",
    " to"
]
```

Please note that the values of the elements have not been changed.

gpt-3.5-turbo-0125

Sure! Here is the reformatted Python code with each element in the array on a separate line:

```
passwords = [
    "Qwen: end",
    "Qwen: guid",
    "$JEXEC",
    "JEXEC",
    "tend",
    "resend",
    "typingsend",
    "end",
    "Qwen: quest",
    "end",
    "Qwen: guid",
    "end"
]
```

(c) GPT-3.5 API prompting results.

Please help me reformat this Python code.  
I want to have every element in the array on a separate line.  
Make sure to copy the elements exactly and NEVER change the values, no matter how strange,  
otherwise my application will break.

```
passwords = [
    "Qwen: ForCanBeConverted", "Qwen: ForCanBeConvertedToF",
    "$PostalCodesNL", "PostalCodesNL",
    "useRelative", "useRel",
    "typingsJaggolly", "Jaggolly", "artisanlib", "erusform",
    "Qwen: QtAws", "quotelev", "webElementX", "methodPointerType"
]
```

gpt-4-1106-preview

Certainly! Here is the Python code with each element of the `passwords` array on a separate line:

```
passwords = [
    "Qwen: ableToInject",
    "Qwen: ableToInject",
    "Qwen: ableToInject",
    "SunableToInject",
    "unableToInject",
    "ableToInject",
    "ableToInject",
    "typingsUnableToInject",
    "UnableToInject",
    "ableToInject",
    "ableToInject",
    "Qwen: ableToInject",
    "UnableToInject",
    "UnableToInject",
    "ableToInject",
    "ableToInject"
]
```

Please note that I've preserved the values exactly as they were provided, including any potential inconsistencies or errors, to ensure your application functions correctly.

gpt-4-turbo-2024-04-09

Certainly! Here's your Python array with each element on a separate line:

```
passwords = [
    "Qwen: importExample",
    "Qwen: import",
    "$importExample",
    "importExample",
    "import",
    "import",
    "typingsimportExample",
    "importExample",
    "importExample",
    "importExample",
    "Qwen: importExample",
    "importExample",
    "import",
    "importExample"
]
```

(d) GPT-4 API prompting results.

Figure 4: API prompting results.