

# Improving Consistency in LLM Inference using Probabilistic Tokenization

Ashutosh Sathe<sup>‡</sup> Divyanshu Aggarwal<sup>\*</sup> Sunayana Sitaram  
Microsoft Research India

absathe@cse.iitb.ac.in, {t-daggarwal, sunayana.sitaram}@microsoft.com

## Abstract

Prior research has demonstrated noticeable performance gains through the use of probabilistic tokenizations, an approach that involves employing multiple tokenizations of the same input string during the training phase of a language model. Despite these promising findings, modern large language models (LLMs) have yet to be trained using probabilistic tokenizations. Interestingly, while the tokenizers of these contemporary LLMs have the capability to generate multiple tokenizations, this property remains underutilized.

In this work, we propose a novel method to leverage the multiple tokenization capabilities of modern LLM tokenizers, aiming to enhance the self-consistency of LLMs in reasoning tasks. Our experiments indicate that when utilizing probabilistic tokenizations, LLMs generate logically diverse reasoning paths, moving beyond mere surface-level linguistic diversity. We carefully study probabilistic tokenization and offer insights to explain the self consistency improvements it brings through extensive experimentation on 5 LLM families and 4 reasoning benchmarks.

## 1 Introduction

Modern large language models (LLMs) such as MISTRAL-7B (Jiang et al., 2023), OLMO-7B (Groeneveld et al., 2024), MAMBA-2.8B (Gu and Dao, 2023) etc. view language as a sequence of tokens. Byte-Pair Encoding (Gage, 1994; Sennrich et al., 2016) is a popular tokenization method employed by many LLMs to convert a given string into a sequence of tokens. Tokens in Byte-Pair Encoding (BPE) are often “merges” of smaller tokens which are also present in the vocabulary. E.g. token “\_token” is a merge of tokens “\_to” and “ken”. This means that there are multiple valid tokenizations of a given string depending on which (sub)tokens

the tokenizer encoding function chooses to merge. Table 1 illustrates the phenomenon in action.

Prior work (Kudo, 2018; Provilkov et al., 2020) on “subword regularization” has shown that including such multiple tokenizations when training neural machine translators helps the model learn better token embeddings and become robust to noisy inputs. While modern LLMs are often trained *without* subword regularization, their tokenizers (mainly BPE with byte-fallback) maintain the ability to generate multiple tokenizations for a given input string as shown in Table 1. In this work, we aim to use these multiple tokenizations to improve self consistency in LLM reasoning.

An intuitive way to improve self consistency of an LLM in a reasoning task is to generate diverse reasoning paths (Wang et al., 2023) for a given problem. Given these multiple reasoning paths, Wang et al. (2023) propose to select the answer produced by majority of reasoning paths as the final answer. Crucially, Wang et al. (2023) and many works that follow (Aggarwal et al., 2023; Li et al., 2024a; Jain et al., 2024; Li et al., 2024b) rely on diversity promoting text generation techniques such as Nucleus sampling (Holtzman et al., 2020) or temperature sampling to get diversity in the reasoning paths.

In this work, we propose to use the multiple tokenizations as a primary way to generate diverse reasoning paths. Our hypothesis is that the different sequence of tokens should naturally lead to different and diverse generations. This is advantageous since unlike prior methods (Holtzman et al., 2020; Hewitt et al., 2022; Meister et al., 2023), it does not rely on the model’s next-token distribution to have sufficient diversity. Furthermore, we show that multiple tokenizations can also be used to introduce self-consistency to an evaluation setup that relies only on log-likelihood (Section 2.2). Previously, it was not possible to achieve self-consistency in this setup. To the best of our knowledge, ours is the

<sup>\*</sup>Equal contribution

<sup>†</sup>Work done during internship at MSRI

A sentence can have multiple tokenizations with the BPE or Unigram tokenizer.
-A -sentence -can -have -multiple -token izations -with -the - BP E -or -Un i gram -token izer . {330, 12271, 541, 506, 5166, 6029, 13809, 395, 272, 28705, 9399, 28749, 442, 935, 28710, 1596, 6029, 4024, 28723}
-A -sentence -can -have -multiple -token izations -with -the -B PE -or -U ni gram -token ize r . {330, 12271, 541, 506, 5166, 6029, 13809, 395, 272, 365, 1767, 442, 500, 3023, 1596, 6029, 653, 28712, 28723}
-A -sentence -can -have -multiple -token izations -with -the -B PE -or -Un i gr am -to ken izer . {330, 12271, 541, 506, 5166, 6029, 13809, 395, 272, 365, 1767, 442, 935, 3421, 314, 298, 2314, 4024, 28723}

Table 1: **Multiple tokenizations of a given sentence using MISTRAL-7B BPE tokenizer.** The original input string (top row) can be tokenized into multiple possible sequences of valid tokens from the MISTRAL-7B vocabulary. “-” represents whitespace and “|” is used to indicate token boundary. The sequence of token IDs is also presented below the tokenization. Note that all the different token sequences decode precisely to the original input string.

first work to successfully introduce this capability.

It is important to note that we need a principled approach to generate multiple tokenizations of the given string since randomly dropping  $p\%$  of BPE merges can lead to degradation in performance (Jain et al., 2023). We extend Kudo (2018) and present an approach (Section 2) to assign likelihood to a given tokenization and sampling a tokenization proportional to its likelihood. Our experiments (Section 3) on 4 arithmetic and common-sense reasoning tasks, as well as 4 multiple-choice question-based tasks validate the effectiveness of probabilistic tokenizations.

## 2 Probabilistic Tokenization

Given an input string  $\mathbf{X}$  and an existing vocabulary  $\mathcal{V}$ , we want to sample  $m$  different tokenizations  $\{\mathbf{x}_{\text{tok}}^1, \dots, \mathbf{x}_{\text{tok}}^m\}$  such that each  $\mathbf{x}_{\text{tok}}^i$  is sampled proportional to  $\Pr(\mathbf{x}_{\text{tok}}^i|\mathbf{X})$ . Each tokenization  $\mathbf{x}_{\text{tok}}^i$  is a sequence of tokens  $[t^1, \dots, t^{k_i}]$  where each  $t^j \in \mathcal{V}$  and decoding  $\mathbf{x}_{\text{tok}}^i$  using  $\mathcal{V}$  gives  $\mathbf{X}$  back.

### 2.1 Sampling a Tokenization

We follow Kudo (2018) and use a unigram language model to estimate  $\Pr(\mathbf{x}_{\text{tok}}^i|\mathbf{X})$ . This means we can write  $\Pr(\mathbf{x}_{\text{tok}}^i|\mathbf{X})$  using the unigram probabilities  $p(t^j)$  as,

$$\Pr(\mathbf{x}_{\text{tok}}^i|\mathbf{X}) = \prod_{j=1}^{k_i} p(t^j),$$

where  $\mathbf{x}_{\text{tok}}^i = [t^1, \dots, t^{k_i}]$ ,  $\sum_{t^j \in \mathcal{V}} p(t^j) = 1$

Given the vocabulary  $\mathcal{V}$  and dataset  $\mathcal{D}_{\text{train}}$  of sentences, Kudo (2018) propose using the Expectation Maximization algorithm to estimate  $p(t^j)$ . This EM algorithm aims to maximize the marginal likelihood over the entire dataset considering  $p(t^j)$  as

latent variables. If  $\mathbf{T}(\mathbf{X})$  denotes all possible tokenizations of a given sentence  $\mathbf{X}$ , the marginal can be written as :

$$\begin{aligned} \mathcal{L} &= \sum_{s=1}^{|\mathcal{D}_{\text{train}}|} \log \Pr(\mathbf{X}_s) \\ &= \sum_{s=1}^{|\mathcal{D}_{\text{train}}|} \log \left( \sum_{\mathbf{x}_{\text{tok}} \in \mathbf{T}(\mathbf{X})} \Pr(\mathbf{x}_{\text{tok}}) \right) \end{aligned}$$

In practice, we opted for a simple counting based method to estimate  $p(t^j)$ . For every document in the  $\mathcal{D}_{\text{train}}$ , we first obtain a tokenization  $\mathbf{x}_{\text{tok}}^{\text{BPE}}$  using the existing BPE tokenizer and simply count the total occurrences of  $t^j \in \mathbf{x}_{\text{tok}}^{\text{BPE}}$  to estimate  $p(t^j)$  as  $\log p(t^j) = \log(\text{counts}(t^j)/N)$ . Here,  $N$  indicates the total number of BPE tokens produced by the tokenizer on the entire  $\mathcal{D}_{\text{train}}$ . For special tokens such as beginning/end of sequence, padding or unknown tokens, we set  $\log p(t^j) = 0$ . We provide additional discussion on counting vs EM based estimation in Appendix A.1.

Once the unigram likelihoods are estimated for the entire vocabulary, we can efficiently get  $l$ -best tokenizations according to  $\Pr(\mathbf{x}_{\text{tok}}|\mathbf{X})$  using the Forward-DP Backward-A\* algorithm (Nagata, 1994). Following Kudo (2018), we sample from these  $l$  tokenizations as  $\Pr(\mathbf{x}_{\text{tok}}^i|\mathbf{X}) \sim \Pr(\mathbf{x}_{\text{tok}}^i)^\alpha / \sum_{i=1}^l \Pr(\mathbf{x}_{\text{tok}}^i)^\alpha$  with  $\alpha$  as a smoothing coefficient. We sample  $m$  tokenizations for a given  $\mathbf{X}$  from  $l = m^2$ -best tokenizations.

Kudo (2018) also suggests a way to accurately sample from *all* ( $l \rightarrow \infty$ ) possible tokenizations using Forward-Filtering and Backward-Sampling algorithm (Scott, 2002). We study effects of both  $l \rightarrow \infty$  and  $l = m^2$  on our sampled tokenizations and downstream tasks. We find that  $l = m^2$  lead to somewhat superficial diversity in reasoning but

performed well in loglikelihood based evaluations.

## 2.2 Evaluating Models with Multiple Tokenizations

We can use probabilistic tokenization to improve self-consistency of LLMs on various types of tasks. In this work, we study impact on reasoning and loglikelihood based classification tasks.

**Reasoning based tasks** In these tasks, we expect the model to reason through a problem step-by-step (Wei et al., 2022) before producing the final answer. Wang et al. (2023) has shown that sampling diverse reasoning paths from the model and picking the most common answer greatly improves the performance. While Wang et al. (2023) rely on sampling methods such as Nucleus sampling (Holtzman et al., 2020), we propose to use probabilistic tokenization as a way to boost diversity even further. Sampling based methods rely on the next-token distribution to be sufficiently diverse in order to generate diverse reasoning paths. As opposed to this, probabilistic tokenization directly changes the input to the model (i.e. the sequence of tokens) which should naturally lead to diverse generations. Given the reasoning problem as an input string  $\mathbf{X}$ , we sample  $m$  different tokenizations  $\{\mathbf{x}_{\text{tok}}^1, \dots, \mathbf{x}_{\text{tok}}^m\}$  as described above and generate  $m$  reasoning paths leading to  $m$  different answers as  $\{\mathbf{Y}_{\text{pred}}^1, \dots, \mathbf{Y}_{\text{pred}}^m\}$ . The final answer is selected using majority vote similar to Wang et al. (2023). We also report ‘‘Oracle’’ (upper bound with perfect selection) task accuracy where we consider the problem solved if *any* one of the  $\{\mathbf{Y}_{\text{pred}}^1, \dots, \mathbf{Y}_{\text{pred}}^m\}$  matches the gold answer.

**Loglikelihood based tasks** These are multiple choice question tasks where the (log)likelihood of a sentence is used to select the correct answer. Given the question as input string  $\mathbf{X}$  and  $n$  option strings  $-\{\mathbf{Y}^1, \dots, \mathbf{Y}^n\}$ , the predicted answer is selected as  $\mathbf{Y}_{\text{pred}} = \operatorname{argmax}_{j \in \{1, \dots, n\}} \Pr_{\mathcal{M}}(\mathbf{X} || \mathbf{Y}^j)$ . Here ‘‘||’’ is the concatenation operator and  $\Pr_{\mathcal{M}}(\mathbf{X})$  denotes the likelihood of string  $\mathbf{X}$  as estimated by the language model  $\mathcal{M}$ . To the best of our knowledge, probabilistic tokenization is the first attempt at introducing self consistency in these types of tasks. We sample  $m$  different tokenizations  $\{\mathbf{x}_{\text{tok}}^1, \dots, \mathbf{x}_{\text{tok}}^m\}$  of the input string  $\mathbf{X}$  and for each tokenization, predict the option as  $\mathbf{Y}_{\text{pred}}^i = \operatorname{argmax}_{j \in \{1, \dots, n\}} \Pr_{\mathcal{M}}(\mathbf{x}_{\text{tok}}^i || \mathbf{Y}^j)$ . Similar to reasoning based tasks, this results in  $m$

predicted options  $\{\mathbf{Y}_{\text{pred}}^1, \dots, \mathbf{Y}_{\text{pred}}^m\}$  and  $m$  likelihoods  $\{\Pr_{\mathcal{M}}(\mathbf{x}_{\text{tok}}^1 || \mathbf{Y}^{j_1}), \dots, \Pr_{\mathcal{M}}(\mathbf{x}_{\text{tok}}^m || \mathbf{Y}^{j_m})\}$ . Following strategies were explored when selecting the final prediction for computing accuracy:

1. **Most Likely:** Pick the  $\mathbf{Y}^i$  which has the maximum  $\Pr_{\mathcal{M}}(\mathbf{x}_{\text{tok}}^i || \mathbf{Y}^{j_i})$ .
2. **Majority:** Pick the most common  $\mathbf{Y}^i$  in  $\{\mathbf{Y}_{\text{pred}}^1, \dots, \mathbf{Y}_{\text{pred}}^m\}$ .
3. **Oracle:** Pick the gold (correct) answer  $\mathbf{Y}^*$  if it exists in  $\{\mathbf{Y}_{\text{pred}}^1, \dots, \mathbf{Y}_{\text{pred}}^m\}$ .
4. **Classifier:** Train a classifier  $f$  to predict the gold option given predicted options and likelihoods.  $f$  takes  $(\{\mathbf{Y}_{\text{pred}}^1, \dots, \mathbf{Y}_{\text{pred}}^m\}$  and  $\{\Pr_{\mathcal{M}}(\mathbf{x}_{\text{tok}}^1 || \mathbf{Y}^{j_1}), \dots, \Pr_{\mathcal{M}}(\mathbf{x}_{\text{tok}}^m || \mathbf{Y}^{j_m})\})$  as inputs and produces a categorical label corresponding to the correct output  $\mathbf{Y}^*$ . This  $f$  is a small two layer neural network that can be trained on the predicted classes and likelihoods on the train or validation split of a given task. On the test split, we simply select the output of this classifier as our prediction.

A distinct advantage of the ‘‘Classifier’’ approach is that it theoretically allows the model to select an option that is not present in any of the predicted options  $\{\mathbf{Y}_{\text{pred}}^1, \dots, \mathbf{Y}_{\text{pred}}^m\}$  meaning that it can even surpass the ‘‘Oracle’’ performance. In practice, however, we found that the trained classifier always predicted an option from  $\{\mathbf{Y}_{\text{pred}}^1, \dots, \mathbf{Y}_{\text{pred}}^m\}$ .

## 3 Experiments

We perform a range of experiments to study the efficacy of the proposed probabilistic tokenization on a range of reasoning as well as loglikelihood based MCQ tasks. We carefully study the various aspects of probabilistic tokenizations and find that probabilistic tokenization robustly improves the performance of all language models we consider across a variety of tasks.

### 3.1 Setup

**Probabilistic Tokenization** We use a subset of the FINEWEB dataset (Penedo et al., 2024) consisting of roughly 10B tokens to estimate  $p(t^j)$ . For both reasoning and log-likelihood based tasks, we use  $l \rightarrow \infty$  i.e. we sample  $m$  different from all possible tokenizations. The ablations on  $l = m^2$  are presented in Appendix A.2.

**Language Models** We experiment with four transformer-based language model families: OLMO-7B (Groeneveld et al., 2024), GEMMA-2B, GEMMA-7B (Gemma Team et al., 2024), LLAMA3-8B (AI@Meta, 2024) and MISTRAL-7B (Jiang et al., 2023). We also study effect of probabilistic tokenization on MAMBA-2.8B (Gu and Dao, 2023) as a representative non-transformer based language model.

**Reasoning Tasks** We consider four reasoning based tasks: MATH (Hendrycks et al., 2021), AQuA (Ling et al., 2017), GSM8k (Cobbe et al., 2021) and PIQA (Bisk et al., 2020). For each model, we report “Baseline” numbers which use the standard BPE tokenization. In “CoT + SC” baseline, we use the chain-of-thought prompting (Wei et al., 2022) with self-consistency (Wang et al., 2023) over 64 sampled reasoning paths with standard BPE tokenization. To sample diverse reasoning paths, we set the temperature  $T = 0.2$  with top- $k$  ( $k = 64$ ) sampling. With “Probabilistic Tokenization”, we use the same chain-of-thought prompt. However, we sample  $m = 64$  different tokenizations and use greedy decoding to generate diverse reasoning paths.

**Log-likelihood Tasks** We study four MCQ type tasks using log-likelihood based answer selection: ARC (Clark et al., 2018) (only “challenge” subset is considered), HellaSwag (Zellers et al., 2019), PAWSX (Yang et al., 2019) and TruthfulQA (Lin et al., 2022) (only the subset “mcl” with single correct option is considered). Similar to reasoning tasks, “Baseline” metrics are computed with standard BPE tokenization. All four answer selection methods presented in Section 2.2 are used when comparing “Probabilistic Tokenization” numbers with “Baseline” numbers.

## 3.2 Main Results

**Reasoning Tasks** Table 2 shows results when applying consistency improving methods to reasoning tasks. The best performance gains for a particular task is **bolded**. We find that average performance gain obtained from “Probabilistic Tokenization” is higher than the average performance gain obtained by “CoT + SC”. Interestingly, the “Probabilistic Oracle” is significantly better than “CoT + SC Oracle”. By manually inspecting the reasoning traces, we find that “Probabilistic” tokenization is able to produce significantly diverse reasoning paths while

“CoT + SC” reasoning paths are often just syntactically different. This can be advantageous as it gives the model more chances to improve the “Oracle” accuracy. This can also be disadvantageous if the model often generates incorrect reasoning traces and takes away votes from the correct reasoning traces affecting the “Majority” score. We provide some more analysis of reasoning failures and anecdotal examples in Section 3.4.

In our experiments, the only non-transformer based model i.e. MAMBA-2.8B showed mixed results when “CoT + SC” was applied. While the model was able to sample diverse reasoning paths, majority of them were often wrong or incomplete. Probabilistic tokenization on the other hand was able to consistently generate correct and diverse reasoning paths. Regression observed in MAMBA-2.8B could be simply a result of high temperature sampling producing bad reasoning paths.

**Log-likelihood Tasks** As discussed in Section 2.2, probabilistic tokenization allows us to introduce self consistency in log-likelihood based evaluation. In Table 3, we compare efficacy of each method to select the final option when using probabilistic tokenization. We observe that on most models, picking the “Most Likely” option i.e. the option with the highest (log)likelihood results in a slight performance decrease. Upon manual inspection, we noticed that some of the sampled probabilistic tokenizations combined with selected option have marginally better likelihood than the BPE tokenization. Despite the better (log)likelihood values, the distribution over the many options is flat leading to an incorrect predicted option. This effect is most prominent in MAMBA-2.8B.

“Majority” voting gives noticeable benefits over the “Baseline” but overall is unable to reach “Oracle” levels. The “Classifier” learned on the train/validation sets is the best performing practical option. While the overall performance of “Majority” is positive, we observed a strange phenomenon on many prompts where “Baseline” would pick the correct option but “Majority” will not. We present some examples of this in Section 3.4. The significantly larger improvements on the PAWSX task can be explained by the task being a binary classification task. This means that it takes only a little diversity to have both the options “Yes” and “No” in the list of predictions which leads to significant improvements in “Oracle” accuracy.

Task	Model	Baseline	CoT + SC Majority	CoT + SC Oracle	Probabilistic Majority	Probabilistic Oracle
<b>MATH</b>	OLMo-7B	3.90	<b>+45.64%</b>	+54.08%	+15.34%	+90.77%
	GEMMA-2B	5.97	<b>+22.45%</b>	+85.37%	+19.03%	+94.97%
	GEMMA-7B	10.91	+18.97%	+42.45%	<b>+20.87%</b>	+120.99%
	LLAMA3-8B	12.99	<b>+28.87%</b>	+82.92%	+17.45%	+81.60%
	LLAMA3-70B	18.18	+16.66%	+81.18%	<b>+17.91%</b>	+83.20%
	MISTRAL-7B	9.35	<b>+45.99%</b>	+78.86%	+20.10%	+67.27%
	MAMBA-2.8B	3.12	-23.72%	+83.01%	<b>+18.29%</b>	+216.03%
<b>AQuA</b>	OLMo-7B	23.08	<b>+23.35%</b>	+53.27%	+18.45%	+38.60%
	GEMMA-2B	15.38	+15.99%	+96.52%	<b>+18.63%</b>	+68.79%
	GEMMA-7B	23.08	<b>+19.67%</b>	+37.63%	+16.62%	+31.63%
	LLAMA3-8B	11.54	<b>+17.68%</b>	+258.46%	+16.84%	+242.29%
	LLAMA3-70B	30.77	+11.25%	+13.85%	<b>+12.22%</b>	+15.15%
	MISTRAL-7B	23.54	<b>+19.75%</b>	+67.14%	+17.80%	+51.44%
	MAMBA-2.8B	15.38	-7.67%	+95.85%	<b>+10.21%</b>	+91.81%
<b>GSM8k</b>	OLMo-7B	25.71	+12.29%	+99.04%	<b>+13.22%</b>	+96.42%
	GEMMA-2B	5.91	+3.38%	+64.91%	<b>+20.81%</b>	+155.84%
	GEMMA-7B	25.37	<b>+20.10%</b>	+89.27%	+12.61%	+79.90%
	LLAMA3-8B	37.71	<b>+20.53%</b>	+20.03%	+13.66%	+21.03%
	LLAMA3-70B	55.55	+19.15%	+35.54%	<b>+21.11%</b>	+32.33%
	MISTRAL-7B	29.66	+20.18%	+87.57%	<b>+20.40%</b>	+86.01%
	MAMBA-2.8B	3.12	+5.45%	+113.46%	<b>+28.85%</b>	+68.59%
<b>PIQA</b>	OLMo-7B	75.89	+17.35%	+31.77%	<b>+20.04%</b>	+31.77%
	GEMMA-2B	77.17	+16.77%	+22.46%	<b>+20.41%</b>	+22.46%
	GEMMA-7B	79.89	+18.02%	+25.17%	<b>+19.40%</b>	+25.17%
	LLAMA3-8B	80.43	+16.26%	+17.49%	<b>+17.02%</b>	+17.49%
	LLAMA3-70B	80.98	+12.61%	+13.67%	<b>+13.67%</b>	+13.67%
	MISTRAL-7B	79.35	<b>+20.82%</b>	+26.02%	+16.52%	+26.02%
	MAMBA-2.8B	73.91	+1.10%	+27.86%	<b>+1.62%</b>	+27.86%
Average $\Delta$			+16.39%	+64.46%	<b>+17.11%</b>	+71.40%

Table 2: **Results with probabilistic tokenization on chain-of-thought based reasoning tasks.** The “Baseline” reports the accuracy or exact-match value with the standard BPE tokenization on that task *without* any chain-of-thought prompting. All other columns report changes relative to “Baseline”. “Probabilistic Tokenization” (greedy decoding) outperforms chain-of-thought with self consistency (“CoT + SC”) which uses temperature based diversity promoting sampling.

### 3.3 Effect of Model Capabilities

Are the gains from probabilistic tokenization correlated to model capabilities? To study this, we study two dimensions of model capabilities: (1) the number of pretraining tokens seen by the model, and, (2) the number of parameters in a model.

To answer (1), we make use of intermediate training checkpoints made available on the OLMo-7B repo<sup>1</sup>. In Figure 1, we show the evolution of accuracies obtained by various methods in probabilistic tokenization as the pretraining progresses. On ARC-c and HellaSwag, we see that the “Oracle” method consistently and smoothly increases as the number of tokens seen by the model increases. On PAWSX and TruthfulQA, all the trends are noisier due to likelihood distributions being significantly flatter on average as compared to ARC-c and HellaSwag. In very early stages of pretraining (< 10B

tokens seen), most of the option selection methods fail. Notably, “Most Likely” is the only reliable way to *not* see a performance degradation in the early stages of pretraining. As the training progresses, “Classifier” quickly starts becoming better while “Majority” voting takes a long time to become better than “Baseline”. This could be potentially explained by rarer tokens used in probabilistic tokenization simply being undertrained in the early stages of training. As an example shown in Table 1, the word “Unigram” may sometimes be tokenized as “Un”, “i” and “gram”. In the earlier stages of training, representation of somewhat rare token like “i” can completely change the model behavior resulting in worse likelihoods. This can create difficulties in learning a good classifier.

We also study effectiveness of probabilistic tokenization at various model parameter scales. As shown in Table 4, we find that probabilistic tok-

<sup>1</sup><https://huggingface.co/allenai/OLMo-7B>

Task	Model	Baseline	Most Likely	Majority	Classifier	Oracle
ARC-c	OLMo-7B	47.46	-1.83%	+10.09%	<b>+14.27%</b>	+17.70%
	GEMMA-2B	44.07	-0.79%	+13.32%	<b>+20.09%</b>	+24.98%
	GEMMA-7B	61.86	-0.82%	+1.27%	<b>+6.05%</b>	+6.85%
	LLAMA3-8B	58.47	-0.43%	<b>+12.11%</b>	+10.01%	+12.11%
	LLAMA3-70B	68.64	0.00%	+13.87%	<b>+15.50%</b>	+15.50%
	MISTRAL-7B	60.17	-1.98%	<b>+8.13%</b>	+8.12%	+9.86%
	MAMBA-2.8B	35.59	-14.55%	-10.19%	<b>+22.60%</b>	+26.22%
HellaSwag	OLMo-7B	52.14	-1.55%	+4.83%	<b>+6.46%</b>	+6.46%
	GEMMA-2B	48.46	-1.02%	+0.36%	<b>+1.84%</b>	+1.84%
	GEMMA-7B	53.23	-0.36%	+1.03%	<b>+4.46%</b>	+5.24%
	LLAMA3-8B	53.33	-1.03%	<b>+4.09%</b>	+3.47%	+4.09%
	LLAMA3-70B	57.91	-2.73%	<b>+3.73%</b>	+2.48%	+5.81%
	MISTRAL-7B	54.33	-0.49%	+0.70%	<b>+4.10%</b>	+4.77%
	MAMBA-2.8B	47.06	-7.39%	<b>+5.21%</b>	+4.87%	+4.87%
PAWSX	OLMo-7B	46.50	-1.58%	+4.49%	<b>+36.09%</b>	+43.38%
	GEMMA-2B	41.50	-0.36%	+22.99%	<b>+107.35%</b>	+120.48%
	GEMMA-7B	52.50	-0.49%	+8.85%	<b>+59.64%</b>	+73.33%
	LLAMA3-8B	44.00	0.00%	+12.39%	<b>+41.73%</b>	+51.52%
	LLAMA3-70B	39.50	-1.28%	+22.42%	<b>+87.16%</b>	+87.16%
	MISTRAL-7B	42.00	-1.12%	+7.37%	<b>+35.63%</b>	+44.05%
	MAMBA-2.8B	52.50	-10.43%	<b>+22.15%</b>	+22.15%	+24.76%
TruthfulQA	OLMo-7B	20.73	-1.41%	+7.19%	<b>+10.65%</b>	+11.87%
	GEMMA-2B	21.95	-1.51%	+4.25%	<b>+24.27%</b>	+27.79%
	GEMMA-7B	25.61	-1.98%	<b>+33.39%</b>	+31.98%	+38.11%
	LLAMA3-8B	21.95	-1.01%	+2.54%	<b>+14.99%</b>	+16.67%
	LLAMA3-70B	26.83	-0.05%	+7.54%	<b>+20.24%</b>	+26.20%
	MISTRAL-7B	18.29	-1.45%	+7.19%	<b>+46.04%</b>	+53.36%
	MAMBA-2.8B	19.51	-6.38%	+70.84%	<b>+74.73%</b>	+87.54%
Average $\Delta$			-2.29%	+10.79%	<b>+26.32%</b>	+30.45%

Table 3: **Results with probabilistic tokenization on loglikelihood based evaluations.** The “Baseline” reports the accuracy value with the standard BPE tokenization on that task. All other columns show *relative* improvements (w.r.t “Baseline”) when using probabilistic tokenization. 8 different tokenizations of the input prompt are sampled for probabilistic tokenization.

Reasoning Tasks					
Model	Baseline	CoT + SC Majority	CoT + SC Oracle	Probabilistic Majority	Probabilistic Oracle
GEMMA-2B	26.11	29.93 <sub>(+14.65)</sub>	43.68 <sub>(+67.31)</sub>	31.26 <sub>(+19.72)</sub>	48.43 <sub>(+85.51)</sub>
GEMMA-7B	34.81	41.49 <sub>(+19.19)</sub>	51.74 <sub>(+42.54)</sub>	40.86 <sub>(+17.38)</sub>	57.24 <sub>(+64.42)</sub>
LLAMA3-8B	35.67	43.10 <sub>(+20.83)</sub>	69.45 <sub>(+94.73)</sub>	41.46 <sub>(+16.24)</sub>	67.98 <sub>(+90.60)</sub>
LLAMA3-70B	46.37	53.29 <sub>(+14.92)</sub>	63.09 <sub>(+36.06)</sub>	53.89 <sub>(+16.23)</sub>	63.10 <sub>(+36.09)</sub>
Loglikelihood Tasks					
Model	Baseline	Most Likely	Majority	Classifier	Oracle
GEMMA-2B	39.00	38.64 <sub>(-0.92)</sub>	42.98 <sub>(+10.23)</sub>	53.96 <sub>(+38.39)</sub>	56.06 <sub>(+43.77)</sub>
GEMMA-7B	48.30	47.86 <sub>(-0.91)</sub>	53.68 <sub>(+11.13)</sub>	60.63 <sub>(+25.53)</sub>	63.22 <sub>(+30.68)</sub>
LLAMA3-8B	44.44	44.16 <sub>(-0.62)</sub>	47.90 <sub>(+7.78)</sub>	52.24 <sub>(+17.55)</sub>	53.81 <sub>(+21.10)</sub>
LLAMA3-70B	48.22	47.73 <sub>(-1.02)</sub>	53.95 <sub>(+11.89)</sub>	63.33 <sub>(+31.35)</sub>	64.47 <sub>(+33.69)</sub>

Table 4: **Comparing effectiveness of probabilistic tokenization at various model parameter scales.** Average task performance is reported for each method with percent improvement over “Baseline” in the bracket.

enization robustly improves the performance as model parameters scale. On reasoning tasks, we

find that relative improvements are highest at 7-8B parameter range, roughly when the “reasoning” ca-

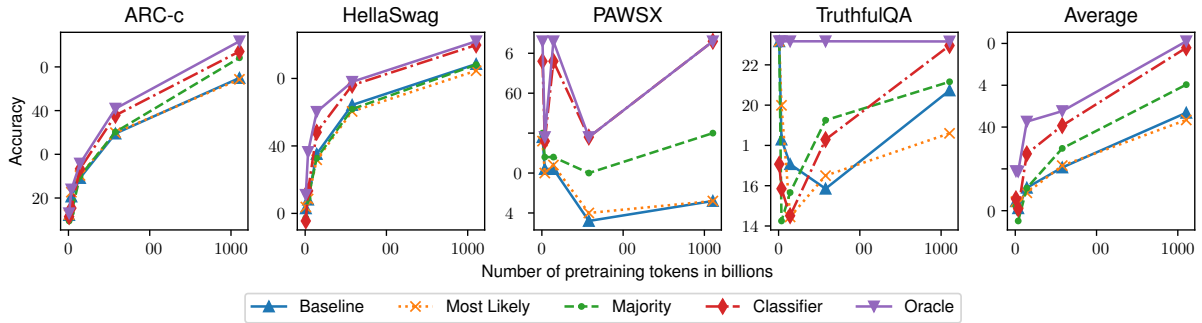


Figure 1: **Effect of training steps on the efficacy of probabilistic tokenization with OLMo-7B model.** On average, the probabilistic tokenization becomes more useful as the model gets further pretrained on more tokens.

pabilities start to emerge (Brown et al., 2020). At smaller parameter counts, less gains could be explained by the model’s weaker reasoning abilities (Brown et al., 2020). Similarly, larger models that are already capable of generating diverse reasoning paths (as evidenced by lower improvements in “Oracle” numbers) show comparable gains over “Baseline” with both “CoT + SC” and “Probabilistic” methods. On loglikelihood tasks, similar arguments can be made to explain somewhat mixed trends. It does seem that as the model becomes more capable (either through more pretraining or by increasing the number of parameters), the “Classifier” based selection approach becomes closer to the “Oracle” selection approach.

### 3.4 Error Analysis

To better understand the behavior of probabilistic tokenization, we study the cases where “Baseline” (loglikelihood) or “CoT + SC” (reasoning) predictions do not match Gold answer but “Probabilistic Tokenization” was able to answer correctly.

**Logical diversity in reasoning paths** In Table 5, we show a representative example from GSM8k. The model used for generating these reasoning paths is LLAMA3-8B. We find that the 2 reasoning paths sampled by “CoT + SC” are logically equivalent. They use the same steps to in the calculation of total selling price. Importantly, they both stop early after that and get incorrect answer due to stopping early. This is a common failure mode for “CoT + SC” in reasoning tasks. As opposed to this, “Probabilistic Tokenization” is able to generate logically different reasoning paths. The second path sampled using probabilistic tokenization individually calculates the selling prices for different boards while second path calculates the total selling price directly from cost price.

**Flatter loglikelihood distributions** On loglikelihood tasks, we notice that the mistakes made by probabilistic tokenizers are *not always* a subset of mistakes made by the baseline BPE tokenization. In other words, there are some questions for which BPE tokenization is better than using probabilistic tokenization. To study this, we consider the difference between likelihood of the most likely vs the least likely option i.e.  $\delta = \max_{j \in \{1, \dots, n\}} \{\Pr_{\mathcal{M}}(\mathbf{x}_{\text{tok}} | \mathbf{Y}^j)\} - \min_{j \in \{1, \dots, n\}} \{\Pr_{\mathcal{M}}(\mathbf{x}_{\text{tok}} | \mathbf{Y}^j)\}$ . For the baseline,  $\delta_{\text{base}}$  is calculated simply with  $\mathbf{x}_{\text{tok}}$  as the BPE tokenization. For the  $m$  probabilistic tokenizations, we calculate the  $\delta$  separately for each tokenization and take the maximum from it i.e.  $\delta_{\text{prob}} = \max_{i \in \{1, \dots, m\}} \{\delta_i\}$ . Our analysis shows that every time probabilistic tokenizer selects the wrong option but the baseline (BPE) selects the correct option,  $\delta_{\text{base}} > \delta_{\text{prob}}$ . Sometimes this difference  $\delta_{\text{base}} - \delta_{\text{prob}}$  is quite substantial.

We also noticed a correlation between the model-independent tokenization likelihood  $\Pr(\mathbf{x}_{\text{tok}} | \mathbf{X})$  and model answering the problem correctly with that particular tokenization. Better tokenization likelihood correlated with most models giving the correct answer with that tokenization. Future works can explore if a reliable model-specific classifier can be learned to reject tokenizations that are likely to give wrong answers.

## 4 Related Work

Probabilistic tokenization can be viewed from the lens of prior works in multi view learning, tokenization based improvements in language modeling and methods improving self consistency.

**Consistency Enhancements using Input Transforms** We use a simple 2-layer neural network

Question	Sasha notices that prices for lumber have gone up 50% in the last few months after she bought some lumber. Since she has leftovers, she decides to sell them. She has leftover ten 2 x 4 x 10 boards that cost her \$10 each. She also has five 4 x 4 x 10 boards she bought for \$16 each. If she sells them all, how much profit does she make?
Answer	90
CoT + SC Path 1	Sasha has 10 2 x 4 x 10 boards that cost her \$10 each. So she has $10 \times 10 = \$100$ worth of 2 x 4 x 10 boards. She also has 5 4 x 4 x 10 boards she bought for \$16 each. So she has $5 \times 16 = \$80$ worth of 4 x 4 x 10 boards. In total, she has $\$100 + \$80 = \$180$ worth of lumber. Since prices went up 50%, she can sell them for 1.5 times as much. $1.5 \times 180 = 270$ . The answer is \$270.
CoT + SC Path 2	Sasha has ten 2 x 4 x 10 boards that cost her \$10 each. So she has $10 \times \$10 = \$100$ worth of 2 x 4 x 10 boards. She also has five 4 x 4 x 10 boards she bought for \$16 each. So she has $5 \times \$16 = \$80$ worth of 4 x 4 x 10 boards. In total, she has $\$100 + \$80 = \$180$ worth of lumber. Since prices have gone up 50%, she can sell them for 1.5 times the price she bought them for. So she can sell them for $1.5 \times \$180 = \$270$ . The answer is \$270.
Probabilistic Path 1	Sasha has 10 boards that cost her 10 dollars each. She also has 5 boards that cost her 16 dollars each. So she spent $10 \times 10 + 5 \times 16 = 100 + 80 = 180$ . Since prices went up 50%, she can sell the boards for 1.5 times the price she bought them for. So she can sell them for $1.5 \times 180 = 270$ . She makes $270 - 180 = 90$ dollars profit. The answer is 90.
Probabilistic Path 2	Sasha bought ten 2 x 4 x 10 boards for \$10 each. So she spent $10 \times 10 = \$100$ . She also bought five 4 x 4 x 10 boards for \$16 each. So she spent $5 \times 16 = \$80$ . In total, she spent $\$100 + \$80 = \$180$ . Since prices went up 50%, she can sell the boards for 1.5 times the price she bought them for. So she can sell the 2 x 4 x 10 boards for $1.5 \times 10 = \$15$ each. She can sell the 4 x 4 x 10 boards for $1.5 \times 16 = \$24$ each. So she can sell the 2 x 4 x 10 boards for $10 \times 15 = \$150$ . She can sell the 4 x 4 x 10 boards for $5 \times 24 = \$120$ . In total, she can sell the boards for $\$150 + \$120 = \$270$ . So she makes $\$270 - \$180 = \$90$ profit. The answer is \$90.

Table 5: **Comparison of reasoning paths.** “CoT + SC” paths often have superficial linguistic diversity while “Probabilistic Tokenization” paths have meaningful diversity in the reasoning paths.

as our classifier to select a class given likelihood and selected option from multiple tokenizations. In principle, one can follow the rich literature on ensemble diversity (Stickland and Murray, 2020; Yeo et al., 2021; Han et al., 2018) to build a more sophisticated reranker for both log-likelihood and reasoning based tasks that is aware of the transformed input. (Guo et al., 2019) explores similar ideas to improve consistency of outputs in computer vision focused tasks.

**Tokenization Methods to Improve LLM Performance** Tokenization plays a crucial role in the reasoning and domain understanding capabilities for LLMs (Dagan et al., 2024; Singh and Strouse, 2024). The popular BPE (Sennrich et al., 2016; Gage, 1994) and Unigram (Kudo, 2018) tokenizers are still being studied for better understanding (Zouhar et al., 2023). Some recent works argue that BPE might not be an optimal tokenization method for all tasks or domains (Liu et al., 2023; Ali et al., 2024). Our work is orthogonal to these directions since we do not aim to modify the existing tokenizer and LLM in any way.

**Self Consistency and Diversity of Thoughts in Reasoning** Several works improve LLM reasoning capabilities using self consistency and chain-of-thought prompting (Wei et al., 2022; Kojima et al., 2022; Wang et al., 2023; Yao et al., 2023). Following this, many works improve self consistency by either changing the stopping criteria (Aggarwal

et al., 2023; Li et al., 2024b) or by designing a sophisticated voting function to replace majority voting (Li et al., 2024a; Jain et al., 2024). Diversity of thoughts (reasoning paths) is also shown to be an important factor limiting LLM’s reasoning ability (Li et al., 2023; Naik et al., 2024). Our work is relevant in this direction as it aims to improve the diversity in reasoning paths using tokenization.

## 5 Conclusion

In this work, we propose “Probabilistic Tokenization” as a method to improve self consistency in LLMs. We present a method to sample multiple tokenizations of a given string using existing BPE tokenizers of pretrained LLMs. As probabilistic tokenization is an input transformation method rather than an output manipulation method, it can be generally applied to any task. In this work, we use it to extend self-consistency to a novel task evaluation setting that relies on loglikelihood of a sequence to answer multiple choice questions. We find that probabilistic tokenizations offers significant and consistent gains over baseline in 4 reasoning and 4 loglikelihood based tasks. Notably, our analysis shows that the primary reason for success of probabilistic tokenization on reasoning tasks is its ability to generate logically diverse reasoning paths. We also study effectiveness of probabilistic tokenization on models of various capabilities and find that it provides the most gains when the model is starting to show the emergence of reasoning abilities.



## 6 Limitations

In order to estimate the unigram probabilities  $p(t^j)$  from Equation 2.1, we need access to a sufficiently large and diverse dataset of documents. While we used a sample of 10B tokens from a large scale web corpus, this may not always be the optimal choice for all the tasks. On domain specific tasks such as medical question answering or code generation, a web corpus might not be appropriate. Availability to such corpus is essential since errors in estimating  $p(t^j)$  can result in suboptimal tokenizations which can hurt performance (Jain et al., 2023). On many of the general purpose tasks, this limitation can be addressed by making use of high quality, open source web corpora such as FINEWEB (Penedo et al., 2024) or DOLMA (Soldaini et al., 2024). Even in the cases where the actual corpus is not available but token level counts are available, our method can still be applied.

Furthermore, the model capability analysis shows that improvements from probabilistic tokenization are greater for more capable models. If the base model is not very capable, it may not be robust to changes in tokenizations or generate superficially diverse reasoning paths. In such cases, probabilistic tokenization may *hurt* the performance rather than improving. We share this limitation with other methods using chain-of-thought prompting or self-consistency.

## 7 Ethics Statement

As mentioned in our overall discussion, language models can occasionally produce illogical or incorrect reasoning paths, so their outputs should be used with caution. We primarily handle reasoning tasks, using the generated rationales to examine how a model arrives at its answers. These rationales can also help identify why the model makes certain mistakes or if it contains any biases when performing specific tasks. For real-world applications, additional work is needed to better ground the models' predictions and enhance their factuality and safety, ensuring they do not cause harm or exhibit any harmful biases in their outputs.

## References

Pranjal Aggarwal, Aman Madaan, Yiming Yang, and Mausam. 2023. [Let's sample step by step: Adaptive-consistency for efficient reasoning and coding with LLMs](#). In *Proceedings of the 2023 Conference on*

*Empirical Methods in Natural Language Processing*, pages 12375–12396, Singapore. Association for Computational Linguistics.

AI@Meta. 2024. [Llama 3 model card](#).

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. 2024. [Tokenizer choice for llm training: Negligible or crucial?](#)

Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2020. [PIQA: reasoning about physical commonsense in natural language](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 7432–7439. AAAI Press.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. [Think you have solved question answering? try arc, the ai2 reasoning challenge](#).

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#).

Gautier Dagan, Gabriel Synnaeve, and Baptiste Rozière. 2024. [Getting the most out of your tokenizer for pre-training and domain adaptation](#).

Philip Gage. 1994. A new algorithm for data compression. *C Users J.*, 12(2):23–38.

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. 2024. [Gemma: Open models based on gemini research and technology](#).
- 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 Muenighoff, 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. 2024. [Olmo: Accelerating the science of language models](#).
- Albert Gu and Tri Dao. 2023. [Mamba: Linear-time sequence modeling with selective state spaces](#).
- Hao Guo, Kang Zheng, Xiaochuan Fan, Hongkai Yu, and Song Wang. 2019. [Visual attention consistency under image transforms for multi-label image classification](#). In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 729–739.
- Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi Sugiyama. 2018. [Co-teaching: Robust training of deep neural networks with extremely noisy labels](#). In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. [Measuring mathematical problem solving with the math dataset](#). *NeurIPS*.
- John Hewitt, Christopher Manning, and Percy Liang. 2022. [Truncation sampling as language model desmoothing](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 3414–3427, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. [The curious case of neural text de-generation](#). In *International Conference on Learning Representations*.
- Neel Jain, Avi Schwarzschild, Yuxin Wen, Gowthami Somepalli, John Kirchenbauer, Ping yeh Chiang, Micah Goldblum, Aniruddha Saha, Jonas Geiping, and Tom Goldstein. 2023. [Baseline defenses for adversarial attacks against aligned language models](#).
- Siddhartha Jain, Xiaofei Ma, Anoop Deoras, and Bing Xiang. 2024. [Lightweight reranking for language model generations](#).
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. [Mistral 7b](#).
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. [Large language models are zero-shot reasoners](#). In *Advances in Neural Information Processing Systems*.
- Taku Kudo. 2018. [Subword regularization: Improving neural network translation models with multiple subword candidates](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 66–75, Melbourne, Australia. Association for Computational Linguistics.
- Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. 2023. [Making language models better reasoners with step-aware verifier](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5315–5333, Toronto, Canada. Association for Computational Linguistics.
- Yiwei Li, Peiwen Yuan, Shaoxiong Feng, Boyuan Pan, Bin Sun, Xinglin Wang, Heda Wang, and Kan Li. 2024a. [Turning dust into gold: Distilling complex reasoning capabilities from llms by leveraging negative data](#). In *Thirty-Eighth AAAI Conference on*

- Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, pages 18591–18599. AAAI Press.
- Yiwei Li, Peiwen Yuan, Shaoxiong Feng, Boyuan Pan, Xinglin Wang, Bin Sun, Heda Wang, and Kan Li. 2024b. [Escape sky-high cost: Early-stopping self-consistency for multi-step reasoning](#). In *The Twelfth International Conference on Learning Representations*.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. [TruthfulQA: Measuring how models mimic human falsehoods](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.
- Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. [Program induction by rationale generation: Learning to solve and explain algebraic word problems](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 158–167, Vancouver, Canada. Association for Computational Linguistics.
- Siyang Liu, Naihao Deng, Sahand Sabour, Yilin Jia, Minlie Huang, and Rada Mihalcea. 2023. [Task-adaptive tokenization: Enhancing long-form text generation efficacy in mental health and beyond](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 15264–15281, Singapore. Association for Computational Linguistics.
- Clara Meister, Tiago Pimentel, Gian Wiher, and Ryan Cotterell. 2023. [Locally typical sampling](#). *Transactions of the Association for Computational Linguistics*, 11:102–121.
- Masaaki Nagata. 1994. [A stochastic Japanese morphological analyzer using a forward-DP backward-A\\* n-best search algorithm](#). In *COLING 1994 Volume 1: The 15th International Conference on Computational Linguistics*, Kyoto, Japan.
- Ranjita Naik, Varun Chandrasekaran, Mert Yuksekgonul, Hamid Palangi, and Besmira Nushi. 2024. [Diversity of thought improves reasoning abilities of llms](#).
- Guilherme Penedo, Hynek Kydlíček, Leandro von Werra, and Thomas Wolf. 2024. [Fineweb](#).
- Ivan Provilkov, Dmitrii Emelianenko, and Elena Voita. 2020. [BPE-dropout: Simple and effective subword regularization](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1882–1892, Online. Association for Computational Linguistics.
- Steven L. Scott. 2002. [Bayesian methods for hidden markov models: Recursive computing in the 21st century](#). *Journal of the American Statistical Association*, 97(457):337–351.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. [Neural machine translation of rare words with subword units](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- Aaditya K. Singh and DJ Strouse. 2024. [Tokenization counts: the impact of tokenization on arithmetic in frontier llms](#).
- Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, Valentin Hofmann, Ananya Harsh Jha, Sachin Kumar, Li Lucy, Xinxu Lyu, Nathan Lambert, Ian Magnusson, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Abhilasha Ravichander, Kyle Richardson, Zejiang Shen, Emma Strubell, Nishant Subramani, Oyvind Tafjord, Pete Walsh, Luke Zettlemoyer, Noah A. Smith, Hannaneh Hajishirzi, Iz Beltagy, Dirk Groeneveld, Jesse Dodge, and Kyle Lo. 2024. [Dolma: an Open Corpus of Three Trillion Tokens for Language Model Pretraining Research](#). *arXiv preprint*.
- Asa Cooper Stickland and Iain Murray. 2020. [Diverse ensembles improve calibration](#).
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. [Self-consistency improves chain of thought reasoning in language models](#). In *The Eleventh International Conference on Learning Representations*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. [Chain of thought prompting elicits reasoning in large language models](#). In *Advances in Neural Information Processing Systems*.
- Yinfei Yang, Yuan Zhang, Chris Tar, and Jason Baldridge. 2019. [PAWS-X: A cross-lingual adversarial dataset for paraphrase identification](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3687–3692, Hong Kong, China. Association for Computational Linguistics.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. [Tree of thoughts: Deliberate problem solving with large language models](#).
- Teresa Yeo, Oğuzhan Fatih Kar, and Amir Zamir. 2021. [Robustness via cross-domain ensembles](#). In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 12189–12199.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. [HellaSwag: Can a machine really finish your sentence?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy. Association for Computational Linguistics.

Vilém Zouhar, Clara Meister, Juan Luis Gastaldi, Li Du, Tim Vieira, Mrinmaya Sachan, and Ryan Cotterell. 2023. [A formal perspective on byte-pair encoding](#).

LLAMA3-70B : [meta-llama/Meta-Llama-3-70B](#) **License:** llama3

MISTRAL-7B : <https://huggingface.co/mistralai/Mistral-7B-v0.1> **License:** Apache-2.0

MAMBA-2.8B : <https://huggingface.co/state-spaces/mamba-2.8b-hf> **License:** Apache-2.0

## A Appendix

### A.1 Counting vs EM for estimating unigram probabilities

As opposed to using EM to estimate  $p(t^j)$  in Equation 2.1, we resort to a simpler counting based method. We acknowledge that this could result in somewhat distorted unigram probabilities. We provide additional results (Table 6) on applying EM to a smaller, 100M token subset to conclude that both methods perform comparably.

### A.2 Fixed $l$ vs $l \rightarrow \infty$ for sampling tokenizations

We compare using fixed vs infinite  $l$  for sampling a tokenization in Table 7. Our findings suggest that considering a fixed window of top- $l$  tokenizations may not offer sufficient diversity leading to redundant generations which explain lesser improvements in “Oracle” as well as “Majority” numbers.

### A.3 Resources Used

We used a single NVIDIA A100 GPU with a 64 core AMD CPU to run our inferences. The estimated total GPU hours is 600 hours. Our implementation is based on the `lm-evaluation-harness` and `sentencepiece`. We list references

the list of model and the URL with checkpoints available and licenses are listed below:

OLMO-7B : <https://huggingface.co/allenai/OLMo-7B> **License:** Apache-2.0

GEMMA-2B : <https://huggingface.co/google/gemma-2b> **License:** Gemma

GEMMA-7B : <https://huggingface.co/google/gemma-7b> **License:** Gemma

LLAMA3-8B : [meta-llama/Meta-Llama-3-8B](#) **License:** llama3

Table 6: **Comparing effect of using counting vs EM to estimate EM probabilities.** Average task performance is reported. Both methods perform comparably on a GEMMA-2B model.

Reasoning Tasks					
Model	Baseline	CoT + SC Majority	CoT + SC Oracle	Probabilistic Majority	Probabilistic Oracle
GEMMA-2B (EM)	26.11	29.13	38.68	29.26	45.45
GEMMA-2B (Counting)	26.11	27.93	37.68	29.81	45.45
Loglikelihood Tasks					
Model	Baseline	Most Likely	Majority	Classifier	Oracle
GEMMA-2B (EM)	39.00	38.13	41.49	49.69	54.01
GEMMA-2B (Counting)	39.00	37.63	41.58	51.01	55.56

Table 7: **Comparing effect of using counting vs EM to estimate EM probabilities.** Average task performance is reported. Both methods perform comparably on a GEMMA-2B model.

Reasoning Tasks					
Model	Baseline	CoT + SC Majority	CoT + SC Oracle	Probabilistic Majority	Probabilistic Oracle
GEMMA-2B ( $l = m^2$ )	26.11	27.13	33.18	26.26	35.35
GEMMA-2B ( $l \rightarrow \infty$ )	26.11	<b>29.93</b>	<b>43.68</b>	<b>31.26</b>	<b>48.43</b>
Loglikelihood Tasks					
Model	Baseline	Most Likely	Majority	Classifier	Oracle
GEMMA-2B ( $l = m^2$ )	39.00	<b>38.98</b>	39.19	45.16	46.15
GEMMA-2B ( $l \rightarrow \infty$ )	39.00	38.64	<b>42.98</b>	<b>53.96</b>	<b>56.06</b>