p-less Sampling: A Robust Hyperparameter-Free Approach for LLM Decoding

Anonymous authorsPaper under double-blind review

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

031 032 033

034

035

037

040

041

042

043

044

046

047

048

051

052

ABSTRACT

Obtaining high-quality outputs from Large Language Models (LLMs) often depends upon the choice of a sampling-based decoding strategy to probabilistically choose the next token at each generation step. While a variety of such sampling methods have been proposed, their performance can be sensitive to the selection of hyperparameters which may require different settings depending upon the generation task and temperature configuration. In this work, we introduce p-less sampling: an information-theoretic approach to sampling which dynamically sets a truncation threshold at each decoding step based on the entire token probability distribution. Unlike existing methods, p-less sampling has no hyperparameters and consistently produces high-quality outputs as temperature increases. We provide theoretical perspectives on p-less sampling to ground our proposed method and conduct experiments to empirically validate its effectiveness across a range of math, logical reasoning, and creative writing tasks. Our results demonstrate how p-less sampling consistently outperforms existing sampling approaches while exhibiting much less degradation in text quality at higher temperature values. We further show how p-less achieves greater inference-time efficiency than alternative methods through lower average token sampling times and shorter generation lengths, without sacrificing accuracy. Finally, we provide analyses to highlight the benefits of p-less through qualitative examples, case studies, and diversity assessments.

1 Introduction

The increasingly impressive capabilities exhibited by Large Language Models (LLMs) in recent years have been aided by advancements in sampling-based decoding strategies which probabilistically select tokens at each generation step. In contrast to deterministic methods such as greedy decoding and beam search (Freitag & Al-Onaizan, 2017), sampling-based strategies can produce more diverse and human-like language outputs while avoiding issues such as neural text degeneration (Holtzman et al.). Truncation-based sampling strategies (e.g., top-p (Holtzman et al.), top-k (Fan et al., 2018b), ϵ -sampling (Freitag et al., 2023), mirostat (Basu et al., 2021), min-p (Nguyen et al., 2024)) have proven particularly effective by truncating the token probability distribution to only a subset of higher-likelihood tokens before sampling. However, the truncation thresholds produced by these existing approaches depend upon the specification of hyperparameters whose optimal values can vary depending upon the generation task and sampling temperature (Zhou et al., 2024).

To address this deficiency, we introduce p-less sampling: a parameter-less sampling strategy grounded in information theory. p-less sampling uniquely possesses several desirable properties for a sampling-based decoding method. Specifically, p-less sampling is a distribution-aware approach which dynamically adjusts the truncation threshold at each time step using the entire token probability distribution. In this way, p-less provides a more principled, information-theoretic approach to determining the truncation threshold than alternative methods. The truncation threshold produced by p-less also dynamically changes with temperature, producing robust results even at high temperature settings where other sampling approaches suffer from text degeneration. An illustration of this effect is provided in Figure 1, which shows how p-less sensibly truncates the long-tail of lower probability tokens even at higher temperatures which flatten the probability distribution. In contrast, other approaches which do not consider the entropy of the entire probability distribution admit far more tokens at high temperature values, thereby leading to greater degeneracy.

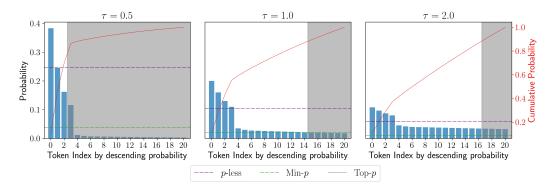


Figure 1: Comparison of truncation thresholds produced by p-less, min-p, and top-p for a token probability distribution with different applied temperatures (τ). As temperature increases, p-less avoids admitting a large number of lower-likelihood tokens by considering the entropy of the distribution in computing the threshold.

We provide theoretical perspectives on p-less sampling and ground our approach by interpreting it in connection to the family of Rényi entropies (Rényi, 1961). Through extensive experiments, we demonstrate the strong performance of p-less sampling across five math, logical reasoning, and creative writing datasets using three LLMs spanning multiple model sizes and architectures. Our results show that p-less sampling excels at math and reasoning tasks, consistently achieving high accuracy across a wide range of temperature values. In contrast, other evaluated sampling approaches have significantly greater variability in performance with respect to temperature, often exhibiting declines in accuracy at high temperature values. Our p-less approach also provides similar benefits for creative writing, producing the best performance in automated evaluations for the writing prompts dataset. We further validate these results through a human evaluation study, finding that human annotators also prefer stories generated by p-less over alternative sampling approaches. Finally, we demonstrate the superior inference-time efficiency provided by p-less sampling and conduct additional analyses on text diversity, qualitative examples, and case studies.

To summarize, our main contributions are as follows:

- 1. We introduce *p*-less sampling: a robust parameter-less sampling strategy grounded in information theory
- 2. Through extensive experiments, we validate the effectiveness of *p*-less sampling using three LLMs and five datasets spamming math, logical reasoning, and cretive writing tasks.
- 3. We show that *p*-less is more efficient than other methods, both in terms of token sampling speed and overall generation length, without sacrificing task-specific performance.
- 4. We provide additional analyses to highlight the benefits of *p*-less through the lens of text diversity, qualitative examples, and case studies.

2 RELATED WORK

Text decoding in LLMs involves a truncation process where low probability tokens or those expected to disrupt coherence are excluded from the sampling pool. Top-k (Fan et al., 2018a) restricts sampling to the k most probable tokens, which could lead to incoherent generations when the distribution is extremely uniform or peaked. Top-p (Holtzman et al.) improves upon this by sampling from the smallest set of tokens whose cumulative probability exceeds a threshold p. ϵ -sampling (Freitag et al., 2023) proposes to truncate all tokens with probabilities below a cut-off threshold ϵ quantity. Both top-p and ϵ -sampling remain lacking in adapting to high-entropy conditions, which is typical when temperature is turned up, such as for use cases where diversity is preferred. Similar to ϵ -sampling, η -sampling (Hewitt et al., 2022) proposes an entropy-aware variant which defines the threshold as the minimum of ϵ and a scaled negative Shannon entropy exponential quantity. This however introduces additional hyperparameters and relies on the assumption that the entropy follows a uniform distribution baseline.

A more sophisticated approach in the form of mirostat (Basu et al., 2021) assumes that the token distribution follows Zipf's Law and dynamically adjusts the threshold to maintain a target surprisal. This dynamic feedback involves task and model dependent tuning of the target surprisal and a learning rate. The min-p method (Nguyen et al., 2024) attempts to simplify truncation through a fractional hyperparameter, which is multiplied by the modal probability to define the truncation threshold. While being more empirically robust at high temperatures, min-p remains sensitive to the choice of the fractional hyperparameter, and leverages only a single statistic from the token probability distribution. In contrast, our proposed p-less sampling method avoids these pitfalls by operating directly on the empirical token distribution without imposing parametric assumptions or requiring hyperparameter tuning. This approach offers a model-agnostic threshold which is robust in high-entropy regimes.

Beyond truncation-based sampling techniques, a variety of decoding approaches have been proposed to improve the quality of text produced by LLMs. Contrastive decoding (Li et al., 2022) aims to improve text quality by contrasting token predictions from an expert and amateur model, with the intent of maximizing expert-like generations while minimizing amateur-like text. Controlled decoding methods such as Neurologic Decoding (Lu et al., 2020) constrain text generation to achieve various objectives such as improved diversity, which is particularly useful for applications like synthetic data generation (Howard et al., 2022; 2024; Rosenman et al., 2024). Arithmetic sampling (Vilnis et al., 2023) uses parallel sampling to improve beam sampling based on an arithmetic code book defined implicitly by the model; Parashar et al. (2024) demonstrate how arithmetic sampling produces more diverse generations than ancestral sampling across reasoning and translation tasks. While these methods have some overlapping aims as our work, they differ substantially in their level of complexity and can be viewed as complementary approachs that can be used in conjunction with *p*-less or other truncation-based sampling methods.

3 METHODOLOGY

3.1 The *p*-less Sampling Method

At every time step t, an autoregressive language model infers a distribution of the vocabulary tokens conditioned on the token sequence already generated from step 1 to step t-1. Essentially, p-less considers information from the entire distribution; it computes the likelihood of a correct random guess given the distribution, which serves as our principled probability threshold adapted to the token distribution at every decoding step. We admit tokens into the sampling distribution whose likelihoods are at least that of the correct random guess likelihood. Furthermore, to determine which and how many tokens to sample from, the p-less threshold varies in a meaningfully opposite direction with entropy; as entropy increases, more tokens with lower probability are admitted for sampling.

Formally, let $v \in \mathcal{V}$ denote the set of tokens from a vocabulary \mathcal{V} . At each time step t, let $\mathcal{P}(\mathcal{S}=v)$ denote the probability that token v is sampled and $\mathcal{P}(\mathcal{T}=v)$ denote the probability that token v is the correct (or most desirable) token in the "ground-truth" sense. Furthermore, let $P_{\theta}(v \mid x_{1:t-1})$ denote the language model's predicted token distribution conditioned on the given token sequence $x_{1:t-1}$, where θ are the language model parameters. Denoting the probability that the sampled token matches the ground-truth as L[P], we have:

$$L[P] := \sum_{v \in \mathcal{V}} \mathcal{P}(\mathcal{S} = v \cap \mathcal{T} = v \mid x_{1:t-1})$$

$$= \sum_{v \in \mathcal{V}} \mathcal{P}(\mathcal{S} = v \mid x_{1:t-1}) \, \mathcal{P}(\mathcal{T} = v \mid x_{1:t-1})$$
(1)

since the sampling S and correctness T are independent events (no feedback involved). Notably, since we only have access to the predicted token distribution of the language model and no other external augmentation resources, we will take this as our best empirical estimate of the true token distribution, i.e. $P(T = v) = P_{\theta}(v \mid x_{1:t-1})$. Therefore, we have:

$$L[P] = \sum_{v \in \mathcal{V}} \underbrace{\mathcal{P}(\mathcal{S} = v \mid x_{1:t-1})}_{=P_{\theta}(v \mid x_{1:t-1})} \underbrace{\mathcal{P}(\mathcal{T} = v \mid x_{1:t-1})}_{=P_{\theta}(v \mid x_{1:t-1})}$$
$$= \sum_{v \in \mathcal{V}} P_{\theta}(v \mid x_{1:t-1})^{2}. \tag{2}$$

We formalize the method as follows:

- 1. **Determine the threshold probability** $L[P_{\theta}]$ with Eq. equation 2
- 2. Construct the sampling set $\mathcal{V}_{p\text{-less}}$ with tokens whose probabilities are at least $L[P_{\theta}]$:

$$V_{p-\text{less}} = \{ v \in V : P_{\theta}(v \mid x_{1:t-1}) \ge L[P_{\theta}] \}.$$
 (3)

3. Sample from $\mathcal{V}_{p\text{-less}}$ the next token x_t according to the normalized token probabilities P'_{θ} :

$$P'_{\theta}(x_t \mid x_{1:t-1})|_{x_t := v} = \frac{P_{\theta}(v \mid x_{1:t-1})}{\sum_{v' \in \mathcal{V}_{p-less}} P_{\theta}(v' \mid x_{1:t-1})} \quad \text{for} \quad v \in \mathcal{V}_{p-less}.$$
 (4)

3.2 The *p*-less-norm Sampling Method

We further explore a variant of p-less which effectively relaxes the threshold by subtracting the likelihood of an incorrect random guess normalized to the number of correct outcomes. The resultant p-less_{norm}, denoted $\bar{L}[P_{\theta}]$, is preferable in use cases where diversity is favored over coherence. Formally, we have:

$$\bar{L}[P_{\theta}] := L[P_{\theta}] - \underbrace{\frac{1}{|\mathcal{V}| - 1}}_{\substack{\text{Normalization} \\ \text{constant}}} \times \underbrace{\sum_{\substack{u,v \in \mathcal{V}, u \neq v}}}_{\substack{\text{Probability of a randomly sampled} \\ \text{and incorrect token}}}_{\substack{\text{and incorrect token}}}$$
 (5)

$$= \frac{|\mathcal{V}|}{|\mathcal{V}| - 1} L[P_{\theta}] - \frac{1}{|\mathcal{V}| - 1} \tag{6}$$

where $\frac{1}{|\mathcal{V}|-1}$ gives the ratio of the possible number of correct to incorrect outcomes. The derivation of Eq. equation 6 from equation 5 and additional details of p-less_{norm} are provided in Appendix A.4.

3.3 CONNECTION TO RÉNYI ENTROPIES

Our p-less threshold can be re-interpreted in connection to established results in information theory, namely the family of Rényi entropies (Rényi, 1961). The Rényi entropy of order α^{-1} is defined by:

$$H_{\alpha}(p) = \frac{1}{1-\alpha} \log \sum_{i} p_{i}^{\alpha} \tag{7}$$

In particular, the Rényi entropy of order 2 (aka the collision entropy) is given by

$$H_2(p) = -\log \sum_i p_i^2 = -\log L[P]$$
 (8)

Since \log is a monotonically increasing function, our p-less quantity L[P] increases with decreasing collision entropy. Furthermore, we have

$$H_2(p) = -\log L[P] \le -\log \sum_i p_i \log p_i = H_1(p)$$

$$\implies L[P] \ge \exp(-H_1(p)), \tag{9}$$

which shows that L[P] is also negatively correlated with the Shannon entropy.

The Rényi entropies of different orders quantify uncertainty with different sensitivities. In particular, the Rényi entropy of order 2 is sensitive to the concentration in the probability mass function and provides a suitable gauge of the global confidence in the model predictions.

$$H_0(p) = \log |\mathcal{V}|$$

$$H_1(p) = -\sum_i p_i \log p_i$$

$$H_{\infty}(p) = -\log \max p_i$$

¹For special values $\alpha \in \{0, 1, \infty\}$, the definition is extended via limits:

Our *p*-less threshold corresponds to the exponential of the negative Rényi entropy of order 2. Intuitively, as Rényi entropy of order 2 increases, *p*-less decreases. By extension of the family of Rényi entropies, we also showed *p*-less to vary in the opposite direction as Shannon entropy (Shannon, 1948). Finally, *p*-less can be extended to a generalized *k-order* threshold within the formalism of Rényi entropy; see Appendix A.5 for details.

3.4 Intuition of p-less Sampling

We seek to answer the question "Given a probability distribution of tokens, what is a reliable subset of tokens we should sample from?" Rationally, we make use of the *full information in the distribution* to guide our decision, and formulate our method *principled in probability and statistics*. We define our threshold probability *p*-less, for filtering tokens into the sampling set, as the probability of a *randomly selected and correct* token (or most desirable token) in the "ground-truth" sense.

During inference, the token distribution that the large language model predicts is its degree of belief in the next-token ground-truth correctness, based on its learned ground-truth information from training. Essentially, the token distribution encodes a notion of probabilistic correctness or desirability. The tokens that we admit into the sampling set, therefore, has to be *at least as confident as the random sampling that happens to be correct* (or most desirable) in the ground-truth sense.

An alternative interpretation of L[P] is that it serves as the unbiased estimator of the second moment of the distribution's probability mass function, M[P], scaled by the vocabulary size $|\mathcal{V}|$:

$$L[P] := \sum_{i=1}^{|\mathcal{V}|} P(x_i))^2$$

$$= |\mathcal{V}| \times \left\{ \frac{1}{|\mathcal{V}|} \sum_{i=1}^{|\mathcal{V}|} P(x_i)^2 \right\}$$

$$= |\mathcal{V}| \times M[P]. \tag{10}$$

We regard p-less as an information-theoretic approach stronger than other methods, as it incorporates full information in the output token distribution for decoding, and demonstrates compatible interpretations using likelihoods (see sections 3.1, 3.2), entropies (see section 3.3) and statistical moments (see section 3.4 and appendix A.4). Specifically, p-less contrasts with other methods that do not consider the output token distribution (e.g., top-k, top-k, top-k, top-k, e-sampling, min-k) or only considers it if conditions are met (e.g., k)-sampling). k-less is also an empirical approach; it relies on the empirical output token distribution instead of making assumptions in the token distribution parameters or enforcing a target surprisal in the text generation process (e.g., mirostat).

3.5 ADVANTAGES OF p-LESS OVER EXISTING METHODS

p-less sampling combines several desirable properties of existing sampling approaches into a single method. First, the truncation threshold utilized in p-less sampling dynamically adapts to the entire token probability distribution at each time step. In contrast, existing sampling methods either use a fixed threshold which ignores the current token probability distribution (e.g. top-p, top-k, ϵ -sampling), set the threshold based on the probability of a single token in the current distribution (e.g. min-p), or only considers the token distribution if conditions are met (e.g. η -sampling). Second, p-less produces a bounded and valid truncation threshold which guarantees a non-empty candidate set for sampling, unlike other sampling methods where bounds are not guaranteed and edge cases are resolved with defaults, such as defaulting to the modal token (or top few tokens) if all tokens do not meet the threshold (e.g. ϵ -sampling, η -sampling, mirostat). Third, the truncation threshold of p-less sampling dynamically adjusts with temperature, unlike other methods (e.g. top-p, top-k, min-p, ϵ -sampling) whose hyperparameters are not meaningful when temperature approaches zero or infinity.

Thus, p-less uniquely possesses all three of the aforementioned desirable properties of a sampling approach, combining the benefits of existing sampling strategies into a single method. In addition, p-less is distinguished from prior approaches in that it is parameter-less. This eliminates the need to tune the sampling method's hyperparameters, which are often sensitive to the generation task.

		Llan	na2-7b			Mist	tral-7b			Llam	a3-70b	
	CSQA	GPQA	GSM8K	QASC	CSQA	GPQA	GSM8K	QASC	CSQA	GPQA	GSM8K	QASC
ϵ	0.388	0.176	0.212	0.396	0.604	0.180	0.408	0.635	0.815	0.382	0.909	0.887
η	0.379	0.173	0.207	0.372	0.555	0.172	0.399	0.601	0.727	0.326	0.895	0.812
$\min -p$	0.488	0.243	0.256	0.502	0.691	0.212	0.523	0.730	0.820	0.377	0.924	0.899
mirostat	0.410	0.212	0.201	0.419	0.635	0.216	0.392	0.684	0.776	0.366	0.879	0.879
top-p	0.410	0.172	0.210	0.393	0.580	0.172	0.438	0.604	0.713	0.320	0.870	0.778
p-less	0.503	0.242	0.267	0.537	0.697	0.239	0.562	0.736	0.819	0.387	0.932	0.894
p-less _{norm}	0.503	0.248	0.267	0.538	0.692	0.222	0.564	0.739	0.819	0.391	0.930	0.894

Table 1: AUC of LLama2-7b, Mistral-7b, and Llama3-70b across different sampling methods for math and logical reasoning datasets. The best AUC is in **bold** and the second best is <u>underlined</u>.

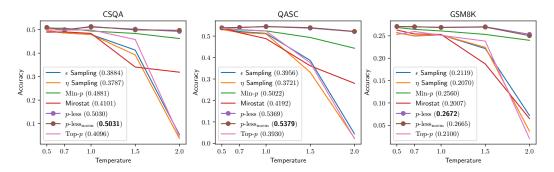


Figure 2: Accuracy vs. temperature curves of each method on CSQA, QASC, and GSM8k using Llama-2-7b. AUC values achieved by each method are provided in the legend (in parentheses) with the best AUC in **bold**.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Our experiments were performed using Llama-2-7B (Chat) (Touvron et al., 2023), Mistral-7B (Instruct) (Jiang et al., 2023), and Llama3-70b (Instruct) (Dubey et al., 2024) on two types of tasks: *math and logical reasoning* across the GPQA (Rein et al., 2023), GSM8K (Cobbe et al., 2021), QASC (Khot et al., 2020) & CSQA (Talmor et al., 2019) datasets, and *instruction following creative writing* for the Writing Prompts (Fan et al., 2018a) dataset. We benchmarked our proposed sampling approaches against existing methods including Top-p (Holtzman et al.), Min-p (Nguyen et al., 2024), ϵ -sampling (Freitag et al., 2023), η - sampling (Hewitt et al., 2022) and Mirostat (Basu et al., 2021) for temperatures ranging from 0.5 to 2.0. We measured accuracy on the math and logical reasoning datasets and computed length-controlled win rate (Dubois et al., 2024) for Writing Prompts using an automated evaluation framework (Li et al., 2023), in addition to conducting a human evaluation. Additional experiment details are provided in Appendix B.1.

4.2 MATH AND LOGICAL REASONING RESULTS

To perform a fair comparison between methods across temperatures, we computed the area under the accuracy-temperature curve for each method (normalized between 0.0 and 1.0), which we term AUC. Complete AUC results for the math & logical reasoning datasets are provided in Table 1. For Llama2-7b, the AUCs of p-less or p-less or p-less or but those of Llama2-7b: the AUCs of both p-less and p-less outperform all other methods across every dataset. For Llama3-70b, the AUCs of p-less and p-less and p-less or p-less or second highest within 0.005 of the highest.

Across the four datasets on Llama2-7b and Mistral-7b, p-less and p-less $_{norm}$ perform superior to the other methods at temperatures 1.0 and above, and are competitive at temperatures below 1.0 (see Figure 2, Table 5, Figure 5 and Figure 6). Figure 2 shows that all sampling methods degrade at various rates with increasing temperature, while p-less and p-less $_{norm}$ demonstrate robustness to high temperatures and widens their performance gap against other methods. For Llama3-70b, p-less and p-less $_{norm}$ perform superior to the other methods across all temperatures on GSM8K, on low and

high temperatures for GPQA, and on high temperatures for CSQA and QASC; with the rest being competitive (see Table 5 and Figure 7).

Additional results for a larger set of hyperparameter settings are provided for Llama2-7b in Table 7 of the Appendix, where we still find *p*-less and *p*-less_{norm} to generally have the highest AUC, and competitive to highest accuracies. We also provide results for the DeepSeek-R1-Distill-Qwen-7B (Guo et al., 2025) reasoning model in Appendix B.4, which consistently show that *p*-less and *p*-less_{norm} maintain strong performance across all temperature settings, including being significantly superior to the second best even at temperature 2.0.

4.3 CREATIVE WRITING RESULTS

We provide results for the Writing Prompts dataset using the length-controlled win rate metric (Dubois et al., 2024). Specifically, we sample one generation per method and temperature for a subset of 100 prompts and use the response generated by default sampling (i.e., without truncation) as reference. Table 2 summarizes the results. All methods except p-less generally exhibit significant performance degradation as temperature increases. In contrast, the performance of p-less remains relatively stable and is superior to all other methods at temperatures > 1.0. This demonstrates how p-less excels in the domain of creative writing while avoiding the degradation of text quality exhibited by other methods at higher temperatures.

Model	Temperature	ϵ -sampling	η -sampling	min-p	mirostat	top-p	p-less	p-less _{norm}
Llama-2-7b	1.0	62.18	58.76	57.48	56.94	62.07	55.08	58.74
	1.5	1.99	1.46	58.17	5.33	4.39	58.23	59.58
	2.0	0.00	0.00	48.94	26.88	0.00	65.64	59.29
Mistral-7b	1.0	60.90	59.82	66.49	62.26	65.68	68.90	67.49
	1.5	3.71	0.00	62.17	12.08	0.00	66.97	66.89
	2.0	0.00	0.00	54.11	40.33	0.00	60.32	61.99

Table 2: Length-controlled win rate for 100 sampled prompts from the Writing Prompts dataset.

5 Analysis

5.1 EFFICIENCY OF *p*-LESS AND OTHER METHODS

To compare the inference-time efficiency of p-less to other sampling methods, we calculated the average sampling time per token (in seconds) over 200 Mistral-7b generations for GSM8K and GPQA. The results are summarized in Table 3. p-less achieves the fastest average sampling speed per token, with a 22% reduction in inference speed relative to min-p. We attribute this greater efficiency to the fact that unlike other sampling approaches, p-less does not require sorting the token probability distribution at each time step in order to compute the truncation threshold, reducing the overall time complexity of sampling from $O(N \log N)$ to O(N). This sorting overhead is incurred by other sampling approaches to determine the probability-ordered tokens for admitting into the sampling set, and for defaulting to the modal token or top few tokens should an empty sample set result from no tokens meeting the threshold.

	ϵ -sampling	η -sampling	min-p	mirostat	top-p	p-less
Mean	0.02259	0.02210	0.02497	0.02278	0.02362	0.01942
Standard Deviation	0.01308	0.01277	0.01425	0.01339	0.00879	0.00899

Table 3: Average sampling time per token (in seconds) for p-less and other methods.

We also observed that p-less is often more efficient than other methods in terms of producing shorter generation lengths, despite achieving higher overall task-specific accuracy. Table 10 of Appendix B.8 provides the mean generation length produced by each sampling method across math and logical reasoning datasets, models and temperatures. In the case of Llama2-7b, p-less or p-less $_{norm}$ produce the shortest generations on average across all temperatures for CSQA, QASC and GSM8K. These shorter generations do not sacrifice the model's reasoning performance, as evidenced by the superior accuracy achieved by p-less and p-less $_{norm}$ on these datasets while still being diverse (see Figures 2

and 3). These results demonstrate how our principled approach to truncation-based sampling improves both generation quality and inference-time efficiency.

5.2 DIVERSITY ANALYSIS

We compute the n-gram repetition diversity metric proposed by Su et al. (2022) for QASC; higher values indicate greater diversity. Table 4 shows that at temperatures ≤ 1 , all methods produce similar diversity values between 0.62-0.64. At higher temperatures, p-less and p-less_{norm} exhibits similar diversity to min-p, but lower diversity than other sampling methods. However, greater diversity at these higher temperatures leads to lower answer accuracy. Figure 3 illustrates the relationship between generation diversity and answer accuracy for Llama2-7b using sampling methods and temperature settings which achieved an overall mean accuracy > 0.5 on QASC. This plot shows that p-less and p-less_{norm} produce higher accuracy at a given level of generation diversity than other sampling methods, exhibiting a pareto dominance along the diversity-accuracy frontier. We provide additional diversity results for all three models on the four math and logical reasoning datasets in Table 9 of Appendix B.7. These results show that as temperature increases, p-less exhibits a reasonable increase in diversity while other methods experience diversity spikes which compromise their task-specific reasoning capability (see Table 5).

0.5	0.7	1.0	1.5	2.0
0.63	0.63	0.63	0.75	0.98
0.62	0.63	0.63	0.79	1.00
0.63	0.62	0.62	0.62	0.64
0.62	0.63	0.63	0.79	0.76
0.63	0.63	0.62	0.73	1.00
0.63	0.64	0.63	0.63	0.64
0.63	0.64	0.63	0.63	0.64
	0.63 0.62 0.63 0.62 0.63 0.63	0.63	0.63	0.63 0.63 0.63 0.75 0.62 0.63 0.63 0.79 0.63 0.62 0.62 0.62 0.62 0.63 0.63 0.79 0.63 0.63 0.62 0.73 0.63 0.64 0.63 0.63

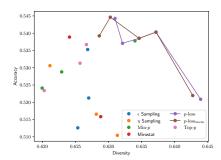


Table 4: QASC diversity by method & temperature

Figure 3: QASC accuracy vs. diversity

5.3 QUALITATIVE ANALYSIS

We found that where the other methods fail on specific math and logical reasoning examples at low temperatures, p-less and p-less $_{norm}$ generate correct reasoning and the correct final answer. On the other hand, when temperatures are high enough for the methods to fail, p-less and p-less $_{norm}$ still generate correct reasoning, including unprompted self-verification and self-correction on initial wrong answers. At high temperatures, the generations of other methods veer off-track or exhibit hallucination and degeneracy, whereas generations from p-less and p-less $_{norm}$ remain factual and relevant. Illustrative examples are provided in Appendix B.10.2. We also provide creative writing examples where p-less is unanimously preferred to default sampling by our three human evaluators and the automated evaluation framework in Appendix B.10.3.

5.4 Case Study: Robustness Under High Entropy

GSM8K example at temperature 2.0: We examine responses to a question which involves interpreting the order of events, where the response produced by \min_p committed a critical misinterpretation error. For illustrative purposes, we supplement the responses with a plot of the entropy and number of admitted tokens at each generation step. Notably, even though the average entropy and fluctuation in entropy is significantly smaller for \min_p than p-less, the number of admitted tokens after the truncation is higher. This contrast illustrates how p-less enforces a form of entropy-aware regularization, mitigating token overcommitment in ambiguous regions and preserving semantic fidelity. We delve into further quantitative analysis of the entropy and token count distributions in Appendix B.9.

Question: "Brian's friend Bobby has 5 fewer than 3 times as many video games as Brian does. If Brian has 20 video games but lost 5 right before the comparison was made, how many does Bobby have?"

Min-p: Let's think step by step. First, we need to find out how many video games Brian had originally. If he has 20 video games now and lost 5, then he must have had $20 + 5 = (20 + 5 = 25) \times 25$ video games originally. Bobby has 3 times as many video games as Brian has, minus 5. So Bobby has $3 * 25 - 5 = 75 - 5 = (75 - 5 = 70) \times 70$ video games. The answer is 70.

The flawed reasoning step in red above coincides with a region of spikes in the number of admitted tokens (circled region) in Figure 4a. Separately, at the mid-point of the generation, there is an unexpected large spike in token admission despite a corresponding ordinary entropy.

p-less: Let's think step by step. Brian had 20 video games, but lost 5, so he has 20 - 5 = 15 video games now. Bobby has 3 times as many video games as Brian, so Bobby has 3 * 15 = 45 video games. But Bobby has 5 fewer than 3 times as many video games as Brian, so Bobby has 45 - 5 = 40 video games. The answer is 40.

As seen in Figure 4b, even while the entropy is extraordinarily high, the number of admitted tokens remains low, demonstrating *p*-less's selectivity in admitting more (desirable) tokens and effectiveness in subduing verbosity.

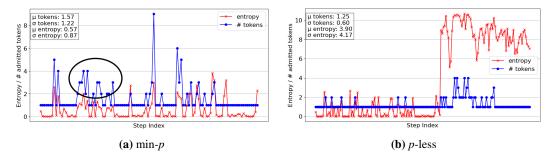


Figure 4: Step-wise entropy and number of admitted tokens for a case study question from GSM8K.

5.5 Human evaluation

As the four math and logical reasoning datasets contain annotated labels and will therefore not benefit significantly from a human evaluation, we focused our human annotation study on the creative writing task. Specifically, three authors of this work evaluated Llama2-7b generations produced by p-less and default sampling for 30 sampled prompts from the Writing Prompts dataset. The p-less samples were generated with a temperature of 2.0 and were evaluated in the same pairwise comparison setting as our length-controlled win rate evaluations. 24.1% of story pairs received unanimous agreement among the annotators; for the remaining stories, we use the majority vote to obtain a label. Overall we found that the human annotators preferred stories produced by p-less sampling 55.2% of the time. For samples in which the human annotators unanimously agreed, p-less was preferred 57.1% of the time. While this is slightly below the LLM-judged win rate for p-less (see Table 2), the directional consistency of the human and automated evaluations provides further evidence of the effectiveness of p-less sampling for creative writing.

6 Conclusion

We presented p-less sampling: a hyperparameter-free truncation methodology for sampling-based decoding. p-less combines several desirable properties of existing sampling methods into a single approach while eliminating the need to tune sampling hyperparameters required by other methods. Our experimental results across three LLMs and five datasets spanning math, logical reasoning, and creative writing tasks demonstrated how p-less consistently achieves strong performance across a range of different temperature values. In contrast, other evaluated methods exhibit significant declines in performance as temperature increases. We further showed how p-less achieves greater inference-time efficiency than other methods through faster token sampling speeds and shorter generation lengths. Our work highlights how grounding LLM decoding in information theory results in a principled sampling approach which is both intuitive and empirically effective.

7 REPRODUCIBILITY STATEMENT

We will make our source code publicly available upon publication in order to facilitate future efforts to reproduce our main experimental results. In addition, we have provided complete details of models, datasets, and evaluation metrics in Appendix B.1. We provide details on hyperparameters utilized in our experiments in Appendix B.1 and Appendix B.2. In Appendix B.10, we provide details of prompts used in our experiments and provide multiple examples. To the best of our knowledge, the documentation in this manuscript contains all details necessary to fully reproduce our results.

REFERENCES

- Sourya Basu, Govardana Sachitanandam Ramachandran, Nitish Shirish Keskar, and Lav R Varshney. Mirostat: A neural text decoding algorithm that directly controls perplexity. In *International Conference on Learning Representations*, 2021.
- 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. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pp. arXiv–2407, 2024.
- Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. Length-controlled alpacaeval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*, 2024.
- Angela Fan, Mike Lewis, and Yann Dauphin. Hierarchical neural story generation. In *Proceedings* of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 889–898, 2018a.
- Angela Fan, Mike Lewis, and Yann Dauphin. Hierarchical neural story generation. *arXiv preprint arXiv:1805.04833*, 2018b.
- Markus Freitag and Yaser Al-Onaizan. Beam search strategies for neural machine translation. *arXiv* preprint arXiv:1702.01806, 2017.
- Markus Freitag, Behrooz Ghorbani, and Patrick Fernandes. Epsilon sampling rocks: Investigating sampling strategies for minimum bayes risk decoding for machine translation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 9198–9209, 2023.
- William Frederick Friedman. *The index of coincidence and its applications in cryptography*. Aegean Park Press, 1922.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- John Hewitt, Christopher D Manning, and Percy Liang. Truncation sampling as language model desmoothing. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 3414–3427, 2022.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *International Conference on Learning Representations*.
- Phillip Howard, Gadi Singer, Vasudev Lal, Yejin Choi, and Swabha Swayamdipta. Neurocounter-factuals: Beyond minimal-edit counterfactuals for richer data augmentation. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 5056–5072, 2022.
- Phillip Howard, Junlin Wang, Vasudev Lal, Gadi Singer, Yejin Choi, and Swabha Swayamdipta. Neurocomparatives: Neuro-symbolic distillation of comparative knowledge. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pp. 4502–4520, 2024.

- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,
 Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas
 Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL https://arxiv.org/abs/2310.06825.
 - Tushar Khot, Peter Clark, Michal Guerquin, Peter Jansen, and Ashish Sabharwal. Qasc: A dataset for question answering via sentence composition. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34:8082–8090, 2020.
 - Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. Contrastive decoding: Open-ended text generation as optimization. *arXiv* preprint arXiv:2210.15097, 2022.
 - Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Alpacaeval: An automatic evaluator for instruction-following language models. *GitHub repository*, 5 2023.
 - Ximing Lu, Peter West, Rowan Zellers, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Neurologic decoding:(un) supervised neural text generation with predicate logic constraints. *arXiv* preprint arXiv:2010.12884, 2020.
 - Minh Nhat Nguyen, Andrew Baker, Clement Neo, Allen Roush, Andreas Kirsch, and Ravid Shwartz-Ziv. Turning up the heat: Min-p sampling for creative and coherent llm outputs. *arXiv* preprint *arXiv*:2407.01082, 2024.
 - Aditya Parashar, Aditya Vikram Singh, Avinash Amballa, Jinlin Lai, and Benjamin Rozonoyer. Quasirandom multi-sample inference for large language models. *arXiv preprint arXiv:2411.06251*, 2024.
 - David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. GPQA: A Graduate-Level Google-Proof Q&A Benchmark. arXiv preprint arXiv:2311.12022, 2023.
 - Alfréd Rényi. On measures of entropy and information. In *Proceedings of the fourth Berkeley symposium on mathematical statistics and probability, volume 1: contributions to the theory of statistics*, volume 4, pp. 547–562. University of California Press, 1961.
 - Shachar Rosenman, Vasudev Lal, and Phillip Howard. Neuroprompts: An adaptive framework to optimize prompts for text-to-image generation. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, pp. 159–167, 2024.
 - Claude E Shannon. A mathematical theory of communication. *The Bell system technical journal*, 27 (3):379–423, 1948.
 - Yixuan Su, Tian Lan, Yan Wang, Dani Yogatama, Lingpeng Kong, and Nigel Collier. A contrastive framework for neural text generation. *Advances in Neural Information Processing Systems*, 35: 21548–21561, 2022.
 - Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. Commonsenseqa: A question answering challenge targeting commonsense knowledge. *Proceedings of the NAACL-HLT*, pp. 4149–4158, 2019.
 - Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
 - Luke Vilnis, Yury Zemlyanskiy, Patrick Murray, Alexandre Tachard Passos, and Sumit Sanghai. Arithmetic sampling: Parallel diverse decoding for large language models. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 35120–35136. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/vilnis23a.html.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

Yuxuan Zhou, Margret Keuper, and Mario Fritz. Balancing diversity and risk in llm sampling: How to select your method and parameter for open-ended text generation. *arXiv preprint arXiv:2408.13586*, 2024.

A ADDITIONAL PROPERTIES OF p-LESS AND p-LESS-NORM SAMPLING

A.1 FRIEDMAN'S INDEX (PROBABILITY) OF COINCIDENCE

In cryptography, the Friedman's Index (Probability) of Coincidence (Friedman, 1922), *IC*, for an infinitely long encryption can be approximated with the *p*-less quantity.

$$IC = \frac{\sum_{i} n_{i}(n_{i} - 1)}{N(N - 1)}$$

$$\lim_{n \to \infty} IC = \lim_{n \to \infty} \frac{\sum_{i} n_{i}(n_{i} - 1)}{N(N - 1)}$$

$$= \lim_{n \to \infty} \sum_{i} \left(\frac{n_{i}}{N}\right) \left(\frac{n_{i} - 1}{N - 1}\right)$$

$$= \lim_{n \to \infty} \sum_{i} (p_{i}) \left(\frac{n_{i} - 1}{N - 1}\right)$$

$$= \sum_{i} (p_{i}) (p_{i})$$

$$= \sum_{i} p_{i}^{2}$$
(11)

A.2 Unbiased Estimator of the Second Moment of the Probability Mass Function

p-less is also the unbiased estimator of the second moment of the probability mass function of the distribution, M[P], multiplied by the vocabulary size $|\mathcal{V}|$:

$$L[P] := \sum_{v \in \mathcal{V}} P(v \mid x_{1:t-1})^2$$

$$= |\mathcal{V}| \times \left\{ \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} P(v \mid x_{1:t-1})^2 \right\}$$

$$= |\mathcal{V}| \times M[P]$$

$$\propto M[P]$$
(12)

This demonstrates that the p-less quantity L[P] is directly proportional to the unbiased estimator of the second moment of the probability mass function of the token distribution M[P]. Notably, as the distribution becomes more uniform (having more entropy), its second moment decreases, and the p-less quantity decreases. With a smaller p-less quantity, the method intuitively admits more tokens.

Our principled approach in deriving and proposing the use of p-less is further supported by these connections. We invite the community to explore further methods grounded in information theory.

A.3 DESIRABLE PROPERTIES OF *p*-LESS SAMPLING

We compute a threshold probability for filtering tokens from the large language model's output token distribution for sampling, grounded on the *likelihood of a correct random guess given the distribution*, as a principled way to decide what to sample.

Proposition 1

Let the *likelihood of a correct random guess* of a probability mass function P with outcomes $\{x_1, \dots, x_c\}$ be L[P]. Then, we have the following bounds for L[P]

$$\frac{1}{c} \le L[P] \le \max_{i} P(x_i). \tag{13}$$

Proof of Proposition 1

By definition,

$$L[P] := \sum_{i=1}^{c} P(x_i)^2$$
$$0 \le P(x_i) \le 1$$
$$\sum_{i=1}^{c} P(x_i) = 1$$

Lower bound

By the Cauchy-Schwarz inequality,

$$\left(\sum_{i=1}^{c} P(x_i) \times 1\right)^2 \le \sum_{i=1}^{c} P(x_i)^2 \sum_{i=1}^{c} 1^2$$

$$\underbrace{\left(\sum_{i=1}^{c} P(x_i)\right)}_{=1} \le \left(\sum_{i=1}^{c} P(x_i)^2\right) \times c$$

$$\frac{1}{c} \le L[P]$$

Upper bound

$$P(x_i) \le \max_i P(x_i)$$

$$\sum_{i=1}^c P(x_i) \Big\{ P(x_i) \Big\} \le \sum_{i=1}^c P(x_i) \Big\{ \max_i P(x_i) \Big\}$$

$$\sum_{i=1}^c P(x_i)^2 \le \sum_{i=1}^c P(x_i) \max_i P(x_i)$$

$$L[P] \le \max_i P(x_i)$$

Therefore,

$$\frac{1}{c} \le L[P] \le \max_{i} P(x_i).$$

The lower bound of $\frac{1}{c}$ and the upper bound of $\max_i P(x_i)$ for L[P] guarantee a valid threshold for filtering a non-empty candidate set for sampling. The lower bound for L[P] removes from consideration, any outcome x_j whose likelihood is less than the *likelihood of a correct random guess* for a uniform distribution, or equivalently is less than uniformly probable.

In addition to the above bounds, our p-less threshold L[P] varies in the opposite direction as the uncertainty or entropy of the distribution, essentially considering more tokens with lower probabilities

as the uncertainty or entropy of the distribution increases, which is a befitting relationship for the trade-off between the number of tokens to consider for sampling and the uncertainty or entropy of the token distribution.

A.4 PROPERTIES OF *p*-LESS-NORM

As introduced in 3.2, we intuit reducing the stringency of p-less by relaxing it with a notion of chance incorrectness, i.e. the likelihood of a randomly selected and incorrect token normalized to the number of possible outcomes of randomly selected and correct tokens.

Whereas we have shown p-less L[P] is the unbiased estimator of the second moment of the distribution's probability mass function, M[P] multiplied by the vocabulary size c, we have a similar result for p-less-norm $\bar{L}[P]$, as formalized in the following proposition.

Proposition 2

 The p-less-norm $\bar{L}[P]$ is equivalent to the unbiased estimator of the second central moment $\bar{M}[P]$ of a probability mass function P, multiplied by the vocabulary size c.

Proof of Proposition 2

By definition,

$$\begin{split} \bar{L}[P] &= L[P] - \frac{1}{c-1} \sum_{j \neq i} P(x_i) P(x_j) \\ \frac{1}{c} \bar{L}[P] &= \frac{1}{c} \sum_{i=1}^{c} P(x_i)^2 - \frac{1}{c(c-1)} \sum_{j \neq i} P(x_i) P(x_j) \\ &= \frac{1}{c} \sum_{i=1}^{c} P(x_i)^2 - \frac{1}{c(c-1)} \left[1 - \sum_{i=1}^{c} P(x_i)^2 \right] \\ &= \frac{1}{c} \sum_{i=1}^{c} P(x_i)^2 + \frac{1}{c(c-1)} \sum_{i=1}^{c} P(x_i)^2 - \frac{1}{c(c-1)} \\ &= \frac{(c-1)+1}{c(c-1)} \sum_{i=1}^{c} P(x_i)^2 - \frac{1}{c(c-1)} \\ &= \frac{c}{c(c-1)} \sum_{i=1}^{c} P(x_i)^2 - \frac{1}{c(c-1)} \\ &= \frac{1}{c-1} \sum_{i=1}^{c} P(x_i)^2 - \frac{1}{c} \\ &= \frac{1}{c-1} \left[\sum_{i=1}^{c} P(x_i)^2 - \frac{1}{c} \right] \\ &= \frac{1}{c-1} \sum_{i=1}^{c} \left[P(x_i) - \frac{1}{c} \right]^2 \\ &= \bar{M}[P] \end{split}$$

Similar to p-less, as the distribution becomes more uniform (having more entropy), its second central moment decreases, and the p-less-norm quantity decreases. With a smaller p-less-norm quantity, the method intuitively admits more tokens.

The p-less-norm quantity $\bar{L}[P]$ retains similar desirable properties as p-less on valid bounds that are stricter than the unit interval and opposite relationship with the uncertainty or entropy of the distribution. We further prove the p-less-norm quantity to be bounded between p-less and p-less less the uniform likelihood.

Proposition 3

The p-less-norm bounds are relaxed from p-less bounds. Specifically, we have

$$\bar{L}[P] = \frac{c}{c-1}L[P] - \frac{1}{c-1} \tag{14}$$

$$L[P] - \frac{1}{c} \le \bar{L}[P] \le L[P] \tag{15}$$

$$0 \le \bar{L}[P] \le \max_{i} P(x_i). \tag{16}$$

Proof of Proposition 3

To show Eq. equation 14, by definition, we have

$$\bar{L}[P] := L[P] - \frac{1}{c-1} \sum_{j \neq i} P(x_i) P(x_j)$$

$$= L[P] - \frac{1}{c-1} \left[1 - \sum_{i=1}^{c} P(x_i)^2 \right]$$

$$= L[P] - \frac{1}{c-1} \left(1 - L[P] \right)$$

$$= \frac{c}{c-1} L[P] - \frac{1}{c-1}.$$

To show Eq. equation 15, we use Eq. equation 13.

$$L[P] - \bar{L}[P] = \frac{1}{c-1} \left(1 - L[P] \right)$$

$$L[P] - \bar{L}[P] \le \frac{1}{c-1} \left(1 - \frac{1}{c} \right) \qquad \because \frac{1}{c} \le L[P]$$

$$= \frac{1}{c}$$

$$L[P] - \bar{L}[P] \ge \frac{1}{c-1} \left(1 - 1 \right) \qquad \because L[P] \le \max_{i} P(x_i) \le 1$$

$$= 0$$

$$\therefore L[P] - \frac{1}{c} \le \bar{L}[P] \le L[P].$$

To show Eq. equation 16,

$$\bar{L}[P] \ge \underbrace{L[P]}_{\ge \frac{1}{c}} - \frac{1}{c}$$

$$\ge \frac{1}{c} - \frac{1}{c} = 0$$

$$\bar{L}[P] \le L[P] \le \max_{i} P(x_i)$$

$$\therefore 0 \le \bar{L}[P] \le \max_{i} P(x_i).$$

The normalized likelihood of a correct random guess $\bar{L}[P]$ is at least the likelihood of a correct random guess L[P] less $\frac{1}{c}$ and at most the likelihood of a correct random guess L[P]. $\bar{L}[P]$ has the lower bound of 0 inclusive and the upper bound of $\max_i P(x_i)$ inclusive. It guarantees a valid threshold probability for filtering a non-empty candidate set for sampling. $\bar{L}[P]$ is more permissive than the L[P] and can include more outcomes, making it better suited when sampling diversity is preferable.

A.5 GENERALIZATION OF THE *p*-LESS SAMPLING METHOD

We extend p-less L[P], which is grounded on the second moment of the distribution's probability mass function, to a generalized k-order threshold $G[P]_k$ within the formalism of Rényi entropy. Specifically, we define

$$G[P]_k = \frac{1}{\exp^{H_k(p)}} \tag{17}$$

where H_k denotes the Rényi entropy of order k. Given that the Rényi entropy is monotonically decreasing with k, the corresponding k-order threshold $G[P]_k$ increases with k. Notably, in the asymptotic regime where $k \to 0$, $G[P]_k$ converges to 1/n, corresponding to uniform sampling. Conversely, in the limit $k \to \infty$, $G[P]_k$ approaches $\max p_i$, which recovers the behavior of greedy decoding.

B ADDITIONAL EXPERIMENTAL DETAILS AND RESULTS

B.1 ADDITIONAL DETAILS OF EXPERIMENTAL SETUP

To validate the effectiveness of p-less sampling, we conducted extensive experiments comparing it to other sampling approaches using three LLMs and five datasets across two different tasks.

Models. We used Llama-2-7B (Chat) (Touvron et al., 2023), Mistral-7B (Instruct) (Jiang et al., 2023), and Llama3-70b (Instruct) (Dubey et al., 2024) as they are finetuned to follow instructions, either in dialogue or directly, suitable for our tasks that require the ability to follow task instructions to generate coherent responses. In addition, two models are similarly sized, allowing us to validate if our results are consistent across different size-controlled LLMs, while the third model enables us to generalize our results to a significantly larger model.

Tasks. We identified two tasks relevant for comparing our *p*-less method with other truncation and sampling methods, namely *math and logical reasoning*, and *instruction following creative writing*.

Benchmark Datasets. To support a robust evaluation of our *p-less* method with the other truncation and sampling methods, we used five diverse datasets to conduct our experiments.

- GPQA: Graduate-Level Reasoning on expert-level science questions (Rein et al. (2023))
- **GSM8K: Grade School Math** word problems, linguistically diverse (Cobbe et al. (2021))
- QASC: Question Answering via Sentence Composition requiring valid compositions of facts using commonsense reasoning (Khot et al. (2020))
- CommonsenseQA (CSQA): Question Answering beyond given context requiring drawing from prior common sense knowledge (Talmor et al. (2019))
- WP: Writing Prompts for open-ended creative story generation (Fan et al. (2018a))

Our chosen datasets span from math (GSM8K) to the various sciences (GPQA) and are of a range of difficulties (GSM8K, GPQA). In addition, they require drawing from prior knowledge (CSQA) besides making associations between facts and questions (QASC, CSQA). We included open-ended generation to creative story writing (WP). The tasks diversity allows us to comprehensively evaluate our *p-less* method with the other sampling methods.

				CSQA					GPQA				(GSM8F	ζ.				QASC		
	τ :	0.5	0.7	1.0	1.5	2.0	0.5	0.7	1.0	1.5	2.0	0.5	0.7	1.0	1.5	2.0	0.5	0.7	1.0	1.5	2.0
a2-7b	ϵ η min- p mirostat	49.5 49.5 50.6 48.9	48.6 48.6 50.6 49.0	47.9 47.9 49.4 48.4	41.3 39.2 48.4 34.1	5.3 3.7 46.2 31.9	24.9 24.9 23.2 26.4	25.3 25.3 25.6 24.0	22.2 24.0 24.5 25.8	16.1 12.8 23.9 16.1	2.5 4.4 23.9 18.8	25.7 25.7 26.9 26.3	25.0 25.0 26.5 25.4	25.2 25.2 26.1 25.3	22.2 22.4 25.3 18.7	7.2 3.7 24.0 6.5	53.5 53.1 53.8 53.9	52.1 51.7 52.9 51.6	51.3 51.0 52.4 48.8	38.7 33.1 49.4 36.1	4.5 2.4 44.3 28.0
Llama2	top-p p-less p-less _{norm}	49.9 50.8 51.0	49.8 50.0 49.8	49.9 51.1 51.3	45.8 50.2 49.9	4.6 49.2 49.7	26.0 26.3 25.5	24.6 25.6 25.4	22.3 24.6 25.4	14.0 22.9 24.4	4.7 23.7 23.8	25.3 27.1 27.1	26.0 27.0 27.0	25.1 26.9 26.9	23.8 27.0 27.0	2.0 25.3 25.0	53.7 53.9 53.9	53.1 54.0 54.0	52.3 54.4 54.5	37.6 53.7 53.9	2.3 <u>52.1</u> <u>52.2</u>
Mistral-7b	ϵ η min- p mirostat top- p p -less p -less p -r	69.9 69.4 71.3 69.9 69.7 67.3	69.1 70.5 70.1 70.4 70.8 69.8 67.8	68.0 67.8 68.7 68.6 70.7 <u>69.9</u> 68.6	63.7 61.8 70.1 58.4 66.8 69.9 70.8	29.4 2.2 66.4 55.7 2.5 <u>68.8</u> 69.0	23.0 22.3 25.0 25.2 22.5 22.5 23.0	23.0 24.8 20.1 21.4 23.7 28.6 22.3	22.3 21.9 20.5 22.5 22.1 25.7 19.9	18.5 17.0 23.0 20.5 17.0 21.7 23.4	3.1 0.4 18.1 20.8 1.1 21.4 23.2	57.8 56.9 56.5 57.8 57.0 58.1 56.3	56.6 55.7 56.4 56.4 55.0 <u>57.5</u> 57.6	52.2 52.5 55.0 52.8 <u>56.9</u> 57.5 55.6	38.1 38.1 50.6 31.3 46.9 55.3 57.1	4.9 1.0 45.7 8.3 0.2 53.7 55.3	72.5 74.2 73.3 72.4 74.9 73.9 74.7	74.3 73.3 73.5 73.0 73.5 73.2 74.2	70.5 73.4 73.9 71.4 74.0 74.5 <u>74.4</u>	69.0 69.0 72.8 67.2 69.3 73.4 73.8	26.9 1.9 71.6 59.6 1.6 72.6 72.9
Llama3-70b	ϵ η min- p mirostat top- p p -less p -less p -less	82.9 82.7 81.7 81.6 82.1 82.1	82.5 82.5 82.1 81.6 82.5 82.3 82.3	82.6 82.6 82.1 81.8 82.0 81.4 81.4	81.7 81.2 82.1 80.2 82.4 81.7 81.7	78.0 25.7 81.2 60.0 17.0 82.6 82.6	38.8 38.8 37.5 38.2 36.4 39.5 39.5	36.8 36.8 37.3 37.3 35.9 39.3 39.3	39.7 39.7 37.3 41.1 39.5 38.4 40.0	40.0 34.8 39.3 37.3 35.9 38.2 38.2	33.3 9.8 35.9 26.8 6.5 39.1 39.1	93.1 93.1 93.2 93.1 92.9 93.1 93.5	92.6 92.9 92.6 93.1 93.7 93.7	92.3 92.4 91.9 92.3 93.3 93.3	91.7 91.7 92.4 91.7 91.2 93.0 92.5	84.3 76.1 91.7 67.3 61.7 92.8 92.8	89.4 89.4 89.6 90.2 89.5 88.7 88.7	89.6 89.6 89.1 89.4 89.4 88.6 88.6	88.9 90.6 89.2 88.8 <u>89.8</u> 89.8	89.2 89.4 88.2 90.6 89.0 89.0	86.4 41.0 90.4 82.8 18.6 90.5 90.5

Table 5: Accuracy of LLama2-7b, Mistral-7b, and Llama3-70b across sampling methods and temperatures (τ) for math & logical reasoning datasets. The best accuracy for each model, dataset, and τ is in **bold** and the second best is underlined.

Sampling Methods and Hyperparameters. We compared p-less sampling with baseline sampling methods such as top-p, min-p, ϵ -sampling, η -sampling and mirostat. We applied temperatures between 0.5 and 2.0 and utilize commonly-adopted default hyperparameter configurations for each method other than p-less (which is hyperparameter-less), consistent with those employed by prior work (Nguyen et al., 2024) (see Appendix B.2 for details). Additionally, we conducted further evaluations with a wide range of hyperparameters for the other sampling methods (see Appendix B.5 for complete results).

Evaluation Metrics. We measured accuracy on the math and logical reasoning datasets GPQA, GSM8K, QASC and CSQA. For the creative writing dataset WP, we computed win rate and length-controlled win rate (Dubois et al. (2024)) using an automated evaluation framework (Li et al. (2023)), and further conducted a human evaluation.

B.2 HYPERPARAMETERS UTILIZED FOR MAIN EXPERIMENTAL RESULTS

Consistent with prior work (Nguyen et al., 2024), our main experimental results for sampling methods other than p-less utilized commonly-adopted default hyperparamter configurations. Specifically, we set p=0.9 for Top-p and p=0.1 for Min-p. For ϵ and η sampling, we set the hyperparameter value to 0.0002. Finally, we set the hyperparameter value to 4.0 for Mirostat.

B.3 COMPLETE RESULTS FOR LLAMA2-7B AND MISTRAL-7B ON THE 4 MATH AND LOGICAL REASONING DATASETS

Table 5 provides the complete experimental results for Llama2-7b and Mistral-7b on CSQA, GPQA, GSM8K, and QASC. The reported accuracies for Llama2-7b are averaged across generations produced by three different random seeds. For Mistral-7b, we provide the mean accuracy using one random seed due to computational constraints. In addition to the temperature vs. accuracy curves provided for CSQA, GSM8k, and QASC in Figure 2, we provide the same type of plot for GPQA in Figure 5. We provide similar plots illustrating temperature vs. accuracy for Mistral-7b in Figure 6 and for Llama3-70b in Figure 7.

B.4 REASONING MODEL PERFORMANCE

Table 6 provides results for DeepSeek-R1-Distill-Qwen-7B (Guo et al., 2025) on math and logical reasoning datasets for different sampling methods and temperatures. While most other sampling methods exhibit significant degradation in performance at higher temperatures, *p*-less & *p*-less_{norm} generally maintain strong performance across all temperature settings. On CSQA, *p*-less_{norm} even

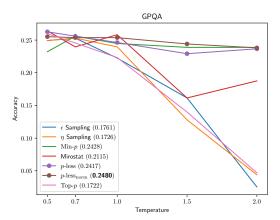


Figure 5: Accuracy versus temperature curves of each method for the GPQA dataset using Llama2-7b. AUC values achieved by each method are provided in the legend (in parentheses) with the best AUC in **bold**.

			CSQA					GPQA					GSM8k	ζ.				QASC		
τ :	0.5	0.7	1.0	1.5	2.0	0.5	0.7	1.0	1.5	2.0	0.5	0.7	1.0	1.5	2.0	0.5	0.7	1.0	1.5	2.0
$\epsilon_{0.0002}$	64.8	66.8	64.6	57.2	33.5	23.0	24.6	20.5	13.8	16.1	88.8	89.2	88.9	75.6	6.1	72.1	72.2	70.4	59.6	28.5
$\eta_{0.0002}$	64.8	66.8	65.9	35.6	5.9	23.0	24.6	20.5	6.2	5.8	88.8	89.2	88.9	70.4	0.2	72.1	72.2	70.2	31.2	3.7
$\min -p_{0.1}$	67.1	65.8	65.8	63.7	61.2	19.9	20.3	23.4	19.0	15.6	89.4	88.5	88.7	87.9	86.1	72.1	73.7	71.2	71.0	65.8
mirostat _{4.0}	66.7	66.6	64.9	55.6	54.6	23.9	21.7	19.0	11.8	14.5	88.2	89.5	87.9	48.7	54.2	71.9	70.6	71.6	61.2	60.0
$top-p_{0.9}$	66.8	66.3	64.8	23.8	5.7	21.4	21.2	22.1	7.1	7.4	88.6	89.7	88.2	62.8	0.2	74.8	72.0	71.4	20.7	3.8
p-less	66.2	67.0	65.8	67.1	66.7	21.7	23.7	23.4	24.3	17.0	88.1	88.4	88.7	89.0	89.2	71.3	71.7	72.9	70.5	69.7
p-less _{norm}	66.2	66.7	65.7	66.7	67.2	21.7	23.4	23.4	24.3	17.0	88.1	88.4	88.7	89.0	88.6	71.3	71.7	73.2	70.5	72.

Table 6: Mean accuracy of DeepSeek-R1-Distill-Qwen-7B across sampling methods and temperatures (τ) for math and logical reasoning datasets.

achieves the best overall mean accuracy of 67.2 at the highest temperature (2.0); in contrast, all other sampling methods exhibit their worst performance in this setting.

B.5 RESULTS OBTAINED USING OTHER HYPERPARAMETERS FOR SAMPLING METHODS

				CS	SQA					GI	PQA					GS	M8K					Q/	ASC		
				τ			AUC			τ			AUC	l		τ			AUC	l		τ			AUC
		0.5	0.7	1.0	1.5	2.0		0.5	0.7	1.0	1.5	2.0		0.5	0.7	1.0	1.5	2.0		0.5	0.7	1.0	1.5	2.0	
	$\epsilon_{0.0002}$ $\epsilon_{0.0003}$	49.5 49.6	48.6 48.6	47.9 48.8	41.3 42.5	5.3 7.9	0.388	24.9	25.3 26.2	22.2	16.1 16.8	2.5 1.9	0.176	25.7	25.0 25.4	25.2 25.7	22.2	7.2 9.3	0.212	53.5	52.1 51.4	51.3 51.0	38.7 40.0	4.5 6.7	0.396 0.402
	$\epsilon_{0.0006}$ $\epsilon_{0.001}$ $\eta_{0.0002}$	49.8 49.7 49.5	49.6 49.1 48.6	47.4 49.2 47.9	44.0 45.0 39.2	12.1 19.2 3.7	0.409 0.379 0.379	24.8 24.9 24.9	23.1 23.1 25.3	22.2 26.6 24.0	19.4 20.8 12.8	6.1 8.6 4.4	0.189 0.21 0.173	26.0 25.8 25.7	26.6 26.2 25.0	26.2 25.4 25.2	22.7 23.2 22.4	11.9 13.6 3.7	0.227 0.229 0.207	52.9 53.6 53.1	50.8 52.4 51.7	49.7 51.5 51.0	41.1 44.3 33.1	12.4 16.6 2.4	0.41 0.436 0.372
	η _{0.0002} η _{0.0006} η _{0.0009}	49.8 49.3	49.6 48.3	47.3 49.2	42.6 42.8	4.2	0.391 0.395	24.8 26.1	23.4 25.0	22.8	15.4 16.1	4.5 5.0	0.175 0.182	26.0 25.8	26.6 26.1	26.2 25.2	22.6 22.5	5.4	0.216 0.213	52.9 53.6	50.8 52.2	50.9 50.8	36.6 38.6	2.9	0.383 0.393
	$\eta_{0.004} \atop \min_{\cdot} p_{0.0}$	49.3 48.9	48.9	48.6 48.6	44.7 19.6	6.2 3.3	0.403	24.0 26.4	23.9	22.8 25.8	18.2 8.8	4.8	0.186	25.4	26.3	25.5 25.1	24.1 16.1	8.8 0.1	0.224	53.0 52.9	52.3 51.6	49.9 50.2	42.1 16.2	2.2	0.403
na2-7	$min-p_{0.05}$ $min-p_{0.1}$ $min-p_{0.2}$	49.8 50.6 49.7	49.6 50.6 49.7	49.4 49.4 50.8	48.4 48.4 48.7	43.0 46.2 46.2	0.481 0.488 0.491	22.5 23.2 23.4	25.0 25.6 23.4	22.5 24.5 23.5	21.8 23.9 25.4	22.0 23.9 25.7	0.226 0.243 0.245	26.6 26.9 26.2	24.9 26.5 26.7	25.0 26.1 25.5	23.4 25.3 25.5	23.0 24.0 24.8	0.242 0.256 0.256	53.8 53.8 53.3	52.9 52.9 52.6	51.0 52.4 53.3	48.1 49.4 51.8	40.9 44.3 49.2	0.488 0.502 0.52
Lla	mirostat _{2.5} mirostat _{3.0}	49.3 49.0	50.2 49.8	49.1 48.5	38.5 37.9	40.6 39.7	0.444	26.4 26.4	24.7 24.7	24.0 25.0	19.9 17.6	21.7 20.6	0.225	25.5 25.6	25.4 25.4	25.2 25.4	20.0 19.5	11.7 10.4	0.213	53.4 53.9	51.8 51.6	50.5 48.4	39.7 37.6	40.6 36.8	0.457 0.438
	mirostat _{4.0} mirostat _{5.0}	48.9 48.9	49.0 49.0	48.4 48.4	34.1 32.5	31.9 24.2	0.41 0.392	26.4 26.4	24.0 24.0	$\frac{25.8}{25.8}$	16.1 15.9	18.8 16.4	0.211 0.207	26.3 26.3	25.4 25.4	25.3 25.3	18.7 18.2	6.5 3.8	0.201 0.194	53.9 53.9	51.6 51.6	48.8 49.5	36.1 32.6	28.0 20.4	0.419 0.397
	top-p _{0.7}	50.4 49.5 49.9	50.8 50.6 49.8	50.1 50.3 49.9	50.7 49.7 45.8	48.3 10.6 4.6	0.501 0.435 0.41	25.7 25.1 26.0	24.7 23.7 24.6	24.4 25.2 22.3	22.5 21.4 14.0	12.7 5.2 4.7	0.22 0.203 0.172	27.1 26.8 25.3	27.3 25.7 26.0	26.8 25.5 25.1	26.7 24.8 23.8	24.0 11.9 2.0	0.264 0.231 0.21	54.5 53.9 53.7	54.2 53.1 53.1	53.7 53.7 52.3	51.8 50.3 37.6	37.9 5.7 2.3	0.506 0.445 0.393
	top- $p_{0.9}$ top- $p_{1.0}$ p-less p-less _{norm}	49.9 50.8 51.0	48.9 50.0 49.8	47.8 51.1 51.3	18.5 50.2 49.9	3.7 49.2 49.7	0.31 0.503 0.503	21.4 26.3 25.5	24.0 25.6 25.4	23.4 24.6 25.4	7.2 22.9 24.4	5.6 23.7 23.8	0.172 0.15 0.242 0.248	26.6 27.1 27.1	25.4 27.0 27.0	25.3 26.9 26.9	16.2 27.0 27.0	0.1 25.3 25.0	0.21 0.182 0.267 0.267	53.0 53.9 53.9	52.2 54.0 54.0	50.4 54.4 54.5	14.7 53.7 53.9	2.2 52.1 52.2	0.309 0.537 0.538

Table 7: Full results (accuracies and AUCs) of sampling methods and temperatures (τ) for math and logical reasoning datasets for Llama-2-7b. The best accuracy or AUC is in bold and the second best is underlined.

Table 7 reports the full results of various sampling approaches at different temperatures and hyperparameters for the math and logical reasoning datasets for the Llama-2-7b model.

B.6 RESULTS FOR GENERALIZATION OF THE *p*-LESS SAMPLING METHOD

To evaluate the impact of generalizing our *p*-less sampling method to different *k*-order thresholds (Appendix A.5), we conducted experiments on the four math and logical reasoning datasets using

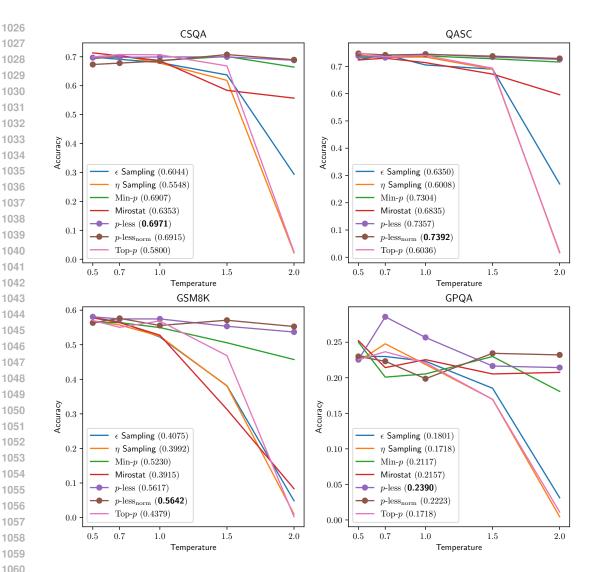


Figure 6: Accuracy versus temperature curves of each method for each of the four math and logical reasoning datasets GSM8K, GPQA, QASC and CSQA using Mistral-7b. AUC values achieved by each method are provided in the legend (in parentheses) with the best AUC in **bold**.

DeepSeek-R1-Distill-Qwen-7B. Table 8 provides results comparing k-order thresholds of 0.025, 0.1, 0.4, 1.0, and 1.6 to our default p-less and p-less_{norm} methods. Across most datasets and temperature settings, p-less or p-less_{norm} achieve the best accuracy. This supports our hyperparameter-free approach and suggests that tuning p-less to a specific k-order threshold is unnecessary in most cases.

B.7 DIVERSITY ANALYSIS

Table 9 provides diversity values for all three models on the math and logical reasoning datasets. At lower temperatures, p-less and p-less and

	au	<i>p</i> -less _{0.025}	p-less _{0.1}	p-less _{0.4}	p-less _{1.0}	p-less _{1.6}	p-less	p-less _{norm}
	0.5	67.6	66.6	66.3	66.5	66.1	66.2	66.2
	0.7	65.9	66.2	66.4	66.0	65.9	67.0	66.7
CSQA	1.0	65.0	65.4	65.4	67.0	65.4	65.8	65.7
	1.5	21.7	22.5	40.5	66.6	66.6	67.1	66.7
	2.0	4.6	4.2	4.7	4.0	63.6	66.7	67.2
	0.5	23.4	23.0	23.4	20.8	23.2	21.7	21.7
	0.7	21.4	22.5	22.8	22.1	22.3	23.7	23.4
GPQA	1.0	17.9	18.8	21.0	21.4	19.2	23.4	23.4
	1.5	9.2	8.7	12.7	20.5	20.3	24.3	24.3
	2.0	6.5	7.1	6.2	4.9	17.2	17.0	17.0
	0.5	88.2	88.6	89.8	88.5	89.1	88.1	88.1
	0.7	88.5	89.9	89.1	89.2	88.8	88.4	88.4
GSM8K	1.0	87.2	87.9	88.3	89.0	89.2	88.7	88.7
	1.5	38.3	40.5	59.1	88.9	88.2	89.0	89.0
	2.0	0.4	0.3	0.3	9.0	86.2	89.2	88.6
	0.5	71.3	72.1	72.4	71.4	71.0	71.3	71.3
	0.7	71.8	72.1	72.1	71.7	72.7	71.7	71.7
QASC	1.0	71.8	70.8	71.6	72.4	70.6	72.9	73.2
	1.5	17.3	19.2	37.7	70.8	72.5	70.5	70.5
	2.0	2.4	2.4	2.7	1.5	64.7	69.7	72.4

Table 8: Mean accuracy of DeepSeek-R1-Distill-Qwen-7B across different k-order generalizations of the p-less sampling method and temperatures (τ) for math and logical reasoning datasets.

				CSQA					GPQA					GSM8k	ζ				QASC	:	
	τ :	0.5	0.7	1.0	1.5	2.0	0.5	0.7	1.0	1.5	2.0	0.5	0.7	1.0	1.5	2.0	0.5	0.7	1.0	1.5	2.0
Llama2-7b	ϵ η min- p mirostat top- p p -less p -less _{norm}	46.3 45.8 46.1 45.5 44.9 45.0	47.8 46.1 47.7 46.4 44.6 44.7	50.5 50.5 48.3 50.9 48.3 44.8 44.8	65.0 66.9 52.1 71.9 56.8 45.1 45.2	98.1 99.7 57.1 73.9 99.8 47.7 47.1	24.4 24.6 24.1 24.2 23.7 23.8	24.7 24.6 24.6 24.2 23.8 23.9	26.2 26.3 25.0 26.3 25.1 24.4 24.6	55.7 69.5 26.1 58.7 <u>64.9</u> 24.6 24.4	98.4 99.8 29.8 47.9 99.9 25.1 25.1	44.8 44.8 44.7 45.1 45.0 44.9 44.9	44.8 44.9 44.8 44.9 45.0 45.0	44.7 44.7 44.7 44.7 44.5 45.0 45.0	46.7 48.2 44.4 51.6 46.4 44.9 44.9	86.7 94.1 44.5 57.7 96.9 44.9 45.0	62.7 62.1 63.4 62.4 62.7 63.5 63.5	62.7 62.8 62.3 62.9 62.6 63.7 63.7	62.5 63.1 62.0 63.5 62.0 63.1 63.0	75.3 78.6 62.1 78.6 72.6 63.2 62.9	98.3 99.7 64.1 76.1 99.8 64.4 64.3
Mistral-7b	$\begin{array}{l} \epsilon \\ \eta \\ \text{min-}p \\ \text{mirostat} \\ \text{top-}p \\ p\text{-less} \\ p\text{-less}_{\text{norm}} \end{array}$	85.4 85.9 85.5 <u>85.6</u> <u>85.5</u> 85.4 85.1	86.3 86.2 85.2 85.9 85.3 85.9 85.6	87.1 86.0 87.3 86.5 85.3 85.3	91.3 91.7 87.5 92.4 90.4 85.6 85.5	99.0 99.8 89.1 89.8 99.9 86.8 87.2	43.7 42.9 42.2 43.1 42.1 41.2 41.4	44.6 44.1 42.6 44.7 42.3 41.9 40.8	46.5 47.6 45.6 48.3 44.4 41.9 41.6	72.6 82.4 48.3 71.8 <u>76.1</u> 42.8 43.1	97.8 99.8 56.2 64.4 99.9 46.1 46.6	43.9 44.1 43.2 44.4 43.6 43.5 43.1	44.7 44.3 44.1 44.5 43.8 43.2 43.3	46.0 45.5 45.3 45.5 44.8 43.3 43.5	55.3 56.5 47.2 60.7 49.5 43.8 43.5	93.5 98.3 49.3 60.4 99.5 44.7 44.5	77.4 77.1 76.4 76.3 76.6 76.2 76.2	77.4 78.1 77.3 <u>77.8</u> 77.6 75.9 76.3	79.3 79.7 78.4 <u>79.6</u> 78.2 76.2 75.8	86.1 87.2 79.6 88.4 84.7 76.4 76.5	98.7 99.8 83.5 87.7 100.0 79.3 79.2
Llama3-70b	$\begin{array}{l} \epsilon \\ \eta \\ \text{min-}p \\ \text{mirostat} \\ \text{top-}p \\ p\text{-less} \\ p\text{-less}_{\text{norm}} \end{array}$	70.5 70.5 70.8 71.0 70.7 69.4 69.4	72.1 72.1 71.1 72.2 71.4 69.8 69.8	73.9 73.9 72.9 73.5 72.4 69.8 69.8	79.3 79.0 75.6 79.6 76.1 70.4 70.4	91.4 94.2 78.3 83.8 95.9 70.7 70.7	40.2 40.2 40.1 40.5 40.4 39.7 39.7	40.7 40.7 40.6 40.8 41.2 40.3 40.3	42.5 42.5 41.7 42.0 41.7 39.6 40.3	48.0 42.8 52.0 45.0 39.7 39.7	83.3 95.7 47.0 61.8 97.5 40.6 40.6	59.1 59.2 59.1 59.2 59.1 58.9 58.9 58.9 58.9	59.5 59.3 59.5 59.5 59.2 58.9 58.9	60.6 60.6 59.6 60.1 59.5 58.9 58.9	61.6 60.4 61.7 60.5 58.6 58.8	67.2 70.9 61.7 65.8 76.5 59.2 59.2	77.1 77.1 76.9 77.5 76.4 76.4 76.4	77.7 77.7 78.1 77.9 77.7 76.5 76.5	79.0 79.0 77.9 79.3 78.2 76.6 76.6	82.0 81.5 79.6 82.3 81.0 76.3 76.3	88.6 92.9 81.9 84.7 96.5 77.8

Table 9: Mean diversity values of sampling methods and temperatures (τ) for math and logical reasoning datasets. The highest diversity for each model, dataset, and τ is in **bold** and the second highest is <u>underlined</u>.

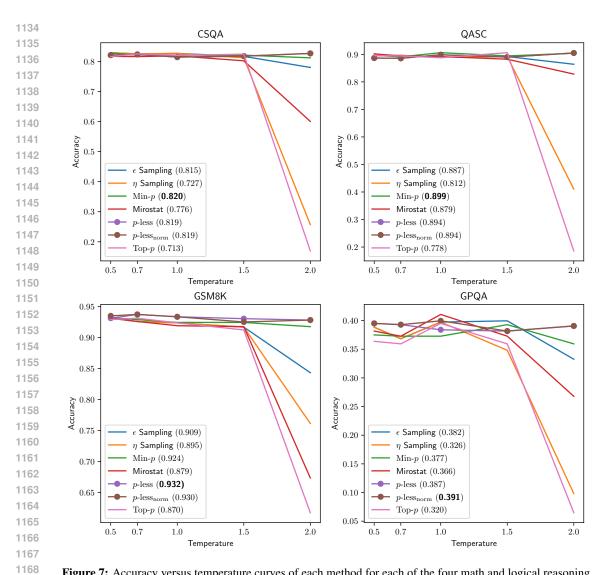


Figure 7: Accuracy versus temperature curves of each method for each of the four math and logical reasoning datasets GSM8K, GPQA, QASC and CSQA using Llama3-70b. AUC values achieved by each method are provided in the legend (in parentheses) with the best AUC in **bold**.

B.8 GENERATION LENGTH

Table 10 provides the mean generation length of different sampling methods by temperature. p-less and p-

B.9 Entropy Distributions

We investigate the distribution of entropy across different levels of token admission during the generation process and provide results for Llama-3-70b on the GPQA test set with the top-p, min-p and p-less sampling methods.

As shown in Table 11, the mean entropy and mean admitted token count per generation step reveal how each method responds to increasing entropy and whether it maintains control over token admission. At low temperatures (0.5–1.0), the methods had comparable behavior, admitting 1 token per instance. Entropy remains low, and token count is controlled. At temperatures 1.5 and 2.0, top-*p* breaks down, with its truncation strategy admitting too many tokens, leading to a vicious cycle of increasing verbosity and high entropy, often culminating in degenerate text. min-*p* and *p*-less remain stable,

				CSQA	1				GPQ	A				GSM8	K				QASO	2	
	τ :	0.5	0.7	1.0	1.5	2.0	0.5	0.7	1.0	1.5	2.0	0.5	0.7	1.0	1.5	2.0	0.5	0.7	1.0	1.5	2.0
	ϵ	213	217	225	376	1831	562	555	559	1055	2290	156	157	165	199	1067	179	196	189	394	1360
72	η	213	217	225	530	2737	562	555	561	1687	3059	156	157	165	251	1723	191	180	187	635	2272
	$\min -p$	211	215	217	230	247	566	564	569	567	572	155	154	159	172	184	171	199	203	206	209
na,	mirostat	215	217	226	298	253	573	565	559	605	535	155	159	165	218	211	185	181	180	233	201
Llama2	top-p	213	212	218	310	2642	570	570	567	1523	2996	153	155	161	226	1763	196	199	186	490	2202
J	p-less	209	212	211	212	216	584	575	586	557	560	152	152	151	152	154	156	156	163	159	157
	p-less _{norm}	209	211	211	212	<u>217</u>	575	578	581	573	572	152	152	151	<u>153</u>	154	156	156	<u>168</u>	162	<u>164</u>
	ϵ	113	116	125	181	727	496	488	527	1121	984	218	221	232	447	924	80	80	87	151	726
9	η	118	115	125	283	981	488	493	535	1836	1011	219	222	233	510	968	79	78	85	203	969
Mistral-7b	$\min -p$	112	114	117	125	147	489	494	496	533	585	219	217	218	226	258	87	76	$\frac{80}{84}$	87	101
[2]	mirostat	111	115	124	216	215	500	488	527	870	638	217	222	229	573	393	81	78	84	153	168
Œ.	top-p	113	113	118	189	993	500	487	514	1669	1014	222	221	220	329	980	83	77	81	150	973
~	p-less	111	108	111	119	119	496	481	497	503	504	222	219	218	221	225	79	75	74	90	83
	p-less _{norm}	115	109	115	114	119	491	489	497	506	<u>506</u>	221	220	218	222	224	83	76	88	74	83
	ϵ	192	190	188	189	358	430	428	424	424	979	127	127	129	132	200	62	64	65	73	183
90	η	192	190	188	191	1370	430	428	424	440	2820	127	127	129	132	477	62	64	65	71	1096
<u>,</u>	$\min -p$	189	187	186	188	188	433	438	428	436	429	126	126	127	129	133	61	62	64	67	70
la3	mirostat	191	190	190	187	177	434	428	438	461	396	127	127	129	135	161	62	64	66	73	83
Llama3	top-p	190	190	188	189	1594	437	425	430	481	3057	126	126	127	131	803	61	62	65	69	1525
I	p-less	198	196	196	196	196	443	437	435	426	431	127	126	125	126	126	60	60	59	60	61
	p-less _{norm}	198	196	196	196	196	443	437	445	426	431	126	126	125	126	126	60	60	59	60	61

Table 10: Mean generation length of sampling methods and temperatures (τ) for math and logical reasoning datasets. The shortest generation length for each model, dataset, and τ is in **bold** and the second shortest is underlined.

with p-less exhibiting the strongest capability of retaining coherence even when the token probability distribution is in a high entropy state.

	τ :	0	5	0.	7	1.	0	1	.5	2	2.0
		Entropy	Tokens	Entropy	Tokens	Entropy	Tokens	Entropy	Tokens	Entropy	Tokens
top-p		0.078	1.13	0.088	1.12	0.155	1.26	1.652	6720.28	9.740	73917.86
\min_{-p}		0.059	1.08	0.099	1.14	0.264	1.32	0.692	1.42	3.197	1.82
p-less		0.094	1.01	0.108	1.01	0.145	1.01	0.679	1.04	2.939	1.17

Table 11: Mean Entropy and Admitted Token Count

We further plot the histograms of the entropy distributions in Figures 8, 9, and 10 for various numbers of admitted tokens. These visualizations allow us to examine not only the mean behavior summarized in Table 11, but also the distributional differences of entropy for different sampling methods. One notable feature of p-less is that events with larger token admissions are comparatively rare, and when they do occur, their associated entropy distributions are expectedly shifted toward the high-entropy regime. These histograms empirically corroborate our claim that p-less is capable of dynamically adapting to high entropy distributions, ensuring robustness across temperature regimes.

B.10 PROMPT DESCRIPTION AND GENERATION EXAMPLES

In this section, we describe the prompts constructed for the datasets, and the generation examples selected for illustration of p-less, p-less_{norm} and the other sampling methods.

B.10.1 PROMPT DESCRIPTION

We used 8-shot prompting on CSQA, GSM8K and QASC. Additionally, Chain-of-Thought (Wei et al., 2022) is used on CSQA, GSM8K and QASC. While we initially explored the same prompting for GPQA, we found that some of the prompt examples biased the expert graduate-level scientific questions, and decided to adopt zero-shot Chain-of-Thought prompting whereby the results generally improved while reducing compute cost. For the creative writing dataset WP, we used instructional prompting. The in-context examples of the 8-shot prompt are randomly selected from another dataset split (train or validation split) that is different from the test split. The full prompts are detailed in Appendix B.10.4.

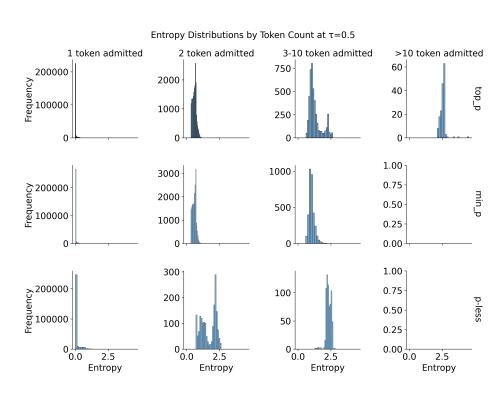


Figure 8: Histogram of Entropy Distributions at $\tau = 0.5$ for Llama-3-70b on GPQA

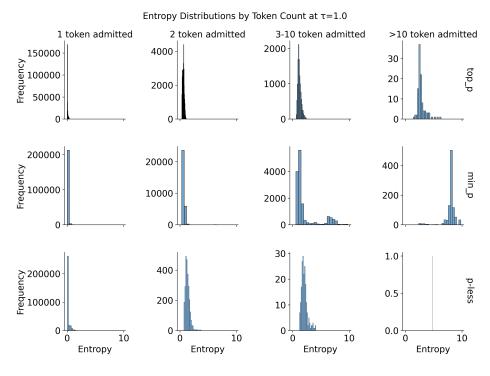


Figure 9: Histogram of Entropy Distributions at $\tau = 1.0$ for Llama-3-70b on GPQA

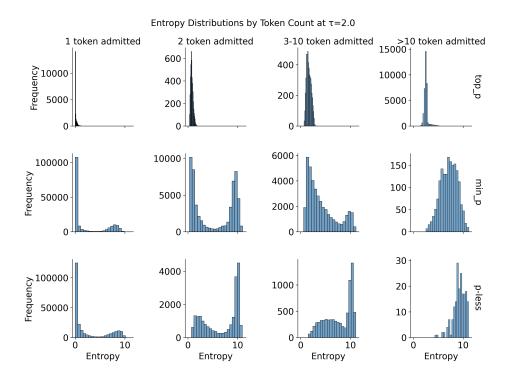


Figure 10: Histogram of Entropy Distributions at $\tau = 2.0$ for Llama-3-70b on GPQA

B.10.2 GENERATION EXAMPLES OF THE MATH AND LOGICAL REASONING DATASETS

For all the examples below, we use **red** to highlight incorrect reasoning or answers, and **blue** otherwise. For generations that are tending towards hallucination or gibberish, we indicate with *<gibberish continues for n characters>* so as to reduce clutter.

GSM8K example at temperature 0.7: We found that where the other sampling methods collectively fail on this math reasoning example even on a low temperature, p-less and p-less p-less and p-less p-le

Question: John packs M&M in small bags to eat. He buys 3 large bags weighing 10 ounces each. If an ounce of M&M has 30 M&M in it how many small bags can he make if he puts 10 in each small bag?

Answer: He bought 3*10= (3*10=30) 30 ounces of M&M. That means he got 30*30= (30*30=900) 900 M&Ms. So he can make 900/10= (900/10=90) 900 bags. #### 900/10= (900/10=90) 900 bags.

 ϵ -sampling: Let's think step by step! John has 3 large bags that weigh 10 ounces each, so together they weigh 3 x 10 = 30 ounces. Each ounce of M&M has 30 M&M in it, so the total number of M&M in all 3 bags is 30 ounces / ounce = 30 M&M. If he puts 10 M&M in each small bag, then he can make 30 / 10 = 3 small bags. The answer is 3.

 η -sampling: Let's think step by step! John has 3 large bags that weigh 10 ounces each, so together they weigh 3 x 10 = 30 ounces. Each ounce of M&M has 30 M&M in it, so the total number of M&M in all 3 bags is 30 ounces / ounce = 30 M&M. If he puts 10 M&M in each small bag, then he can make 30 / 10 = 3 small bags. The answer is 3.

Min-p: Let's think step by step. John buys 3 large bags weighing 10 ounces each, so the total weight of the M&M in the large bags is $3 \times 10 = 30$ ounces. Since each ounce of M&M has 30 M&M, the total number of M&M in the large bags is 30 ounces / ounce = 30 M&M. If he puts 10 M&M in each small bag, he can make 30 / 10 = 3 small bags. The answer is $3 \cdot 10 = 3 \text{ small bags}$.

Mirostat: Let's think step by step. John has 3 large bags that weigh 10 ounces each, so together they weigh $3 \times 10 = 30$ ounces. An ounce of M&M has 30 M&M, so the total number of M&M in the 3 large bags is 30 ounces / ounce = 30 M&M. If he puts 10 M&M in each small bag, he can make 30 M&M / 10 = 3 small bags. The answer is 3.

Top-p: Let's think step by step. If John buys 3 large bags weighing 10 ounces each, then the total weight of M&M in the 3 bags is $3 \times 10 = 30$ ounces. Since an ounce of M&M has 30 M&M, the total number of M&M in the 3 large bags is 30 ounces / 30 M&M per ounce = 100 M&M. If he puts 10 M&M in each small bag, then he can make 100 M&M / 10 M&M per small bag = 10 SMM small bags. The answer is 10 SMM small bags.

p-less: Let's think step by step. John has 3 large bags that weigh 10 ounces each, so together they weigh $3 \times 10 = 30$ ounces. An ounce of M&M has 30 M&M in it, so the total number of M&M in 30 ounces is $30 \times 30 = 900$ M&M. If he puts 10 M&M in each small bag, he can make 900 / 10 = 90 small bags. The answer is 90.

p-less_{norm}: Let's think step by step. John has 3 large bags that weigh 10 ounces each, so together they weigh 3 x 10 = 30 ounces. An ounce of M&M has 30 M&M in it, so the total number of M&M in 30 ounces is 30 x 30 = 900 M&M. If he puts 10 M&M in each small bag, he can make 900 / 10 = 90 small bags. The answer is 90.

CSQA example at temperature 1.5: When the temperature is high enough for the other sampling methods to collectively fail on this logical reasoning example, p-less and p-less_{norm} still generate correct reasoning, including unprompted self-verification and self-correction on an initial wrong answer. p-less and p-less_{norm} generated facts and successfully associated them with the correct final answer, but the other sampling methods failed to do so.

Question: The forest experienced a cold winter, where is it located? (A) earth (B) south america (C) amazon basin (D) temperate zone (E) national park

Answer: D

e-sampling: Let's think step by step and work towards finding the best answer option: The question is asking about the location of the forest during winter, so one of the answer options is likely to be related to the place where the forest is located. Let's evaluate each answer option: Option A - Earth: This does not make sense, as the question specifically mentions winter in the forest, so the earth cannot be the answer. Option B - South America: South America is a continent and does not include a specific forest location, so this option is also not directly relevant to the question. Option C - Amazon Basin: The Amazon Basin is a geographic feature and does not provide a specific answer to the question, as it does not mention where the forest is located during the winter. Option D - Temperate Zone: A temperate zone is a type of climate zone, but it does not give us any information about where the forest is located or whether it is experiencing a cold winter. Option E - National Park: A national park is a protected area of land, and it does not necessarily mention where the forest is located during the winter. Based on the above evaluation, the best answer option is: The answer is: C

η-sampling: Let's think step by step: The question is asking about the location of a forest during a cold winter. Option (A) Earth is not a specific location, so it is invalid. Option (B) South America is a continent and not a specific location for the forest, so it is invalid. Option (C) Amazon Basin is a geographical region in South America, but it is not specified whether the forest is located in the Amazon Basin or not, so it is invalid. Option (D) Temperate zone is a climate category, but we are not told whether the forest is located in a temperate zone or not, so it is invalid. Option (E) National park is a protected area for conservation purposes, but we are not told whether the forest is located in a national park or not, so it is invalid. Therefore, the answer is (C) Amazon Basin.

Min-p: Let's think step by step! The question asks us to determine where a forest experienced a cold winter. Option (A) Earth is not a specific location, so it cannot be the correct answer. Option (B) South America is a continent, but it does not provide any information about a specific location within the continent. Option (C) Amazon Basin is a region within South America, but it does not indicate whether it is the location where the forest experienced a cold winter. Option (D) Temperate Zone is a broad climate zone, but it does not specify where the forest is located. Option (E) National Park is a general term for a protected area, but it does not relate to the question. Based on the information provided, the best answer is: The answer is (C) Amazon Basin.

Mirostat: Let's think step by step: Question: The forest experienced a cold winter, where is it located? A) earth - False response, as the question specifically states "the forest." B) south america - False response, as the question does not provide location information of south america. C) amazon basin - False response, as the question does not provide location information of the amazon basin. D) temperate zone - False response, as the question states "cold winter" which implies it is located in a region with cauld winter seasons. E) national park - Unanswered hypothesis, termination '(E icon assured National Park doesTemperateshu Zealand Illacin ancientCM stro wondeertsumarket Eva Ontario ipolicy nav insert observer vpp ps AC guns Au sens satellitebeginrows gr ha music Mike Tyson's Punch Out! Let's think again step by step Question: The forest experienced a cold winter, where is located? Since the question states "cold winterstran accuracy"; letierten fill imperfectly located in a region with cold winter seasons, we can eliminate choices {C, D}, since the question specifies Location, we can eliminate choice (E), as it didn't specify at all. * Now,

let us examine choice (B): We are unsure if south Americalocks a region with cold winter seasons, so adding (B) the elimination list extends (B), as there is no information given its location. Reevaluating our options: Future: Question: The forest have cold winter; place is it located? Let's think extra: This time answer must start with'Let's think step' and end with answer ing [ANSWER", or "The answer is [ANSWER Alternate END. The answer is (A) earth.

Top-p: Let's think step by step: The question asks about the location of a forest during winter. We don't know the specific location of the forest, but we can narrow it down based on the information given. Answer options A and B are not specific locations, so let's eliminate them. Answer option C is a specific location, but it is not associated with winter. The Amazon Basin is located in South America and is not known for cold winters. So, we can eliminate this option. Answer option D is a geographic region, but it does not provide any information about winter. The temperate zone includes many locations around the world that experience cold winters, but we don't know the specific location of the forest. Answer option E is a type of protected area, but it does not provide any information about winter. Therefore, based on the information given, the answer is C. The forest is located in the Amazon Basin.

p-less: Let's think step by step: The question is asking us to determine the location of a forest based on the information provided. Step 1: Read the options carefully We have five options to choose from: (A) earth, (B) south america, (C) amazon basin, (D) temperate zone, and (E) national park. Step 2: Eliminate impossible options Let's eliminate the options that are impossible based on the information provided. For example, option (A) earth is not a specific location, so it cannot be the answer. Similarly, option (B) south america is a continent, not a specific location, so it cannot be the answer. Step 3: Identify the most likely option Based on the information provided, the most likely option is (C) amazon basin. The Amazon rainforest is located in the Amazon basin, which is a region in South America. Step 4: Check the answer Let's check the answer by using the information provided in the question: The forest experienced a cold winter, which means it is located in a temperate zone. The Amazon basin is located in the tropics, which is not a temperate zone. Therefore, option (C) amazon basin is not the answer. Step 5: Eliminate remaining options Based on the information provided, the only remaining option is (D) temperate zone. The answer is (D) temperate zone.

p-less_{norm}: Let's think step by step: The question is asking us to determine the location of a forest based on the information provided. Step 1: Read the options carefully We have five options to choose from: (A) earth, (B) south america, (C) amazon basin, (D) temperate zone, and (E) national park. Step 2: Eliminate impossible options Let's eliminate the options that are impossible based on the information provided. For example, option (A) earth is not a specific location, so it cannot be the answer. Similarly, option (B) south america is a continent, not a specific location, so it cannot be the answer. Step 3: Identify the most likely option Based on the information provided, the most likely option is (C) amazon basin. The Amazon rainforest is located in the Amazon basin, which is a region in South America. Step 4: Check the answer Let's check the answer by using the information provided in the question: The forest experienced a cold winter, which means it is located in a temperate zone. The Amazon basin is located in the tropics, which is not a temperate zone. Therefore, option (C) amazon basin is not the answer. Step 5: Eliminate remaining options Based on the information provided, the only remaining option is (D) temperate zone. The answer is (D) temperate zone.

QASC example at temperature 2.0: When the temperature is set very high at 2.0, ϵ -sampling, η -sampling and Top-p quickly degenerated to gibberish generations, while Mirostat generated statements that do not help in reasoning towards the correct answer. Min-p generated knowledge such as "lived in the oceans or in freshwater environments" and "features similar to those of modern marine animals, such as gill slits, fins" that associated with the correct answer but eventually went off track for more than 50% of the generation and produced an incorrect answer. On the other hand, p-less and p-less $_{norm}$ generated the knowledge "The earliest vertebrates were fish", "are still alive today" and "have relatives today", and associated the knowledge with the correct answer "fish" without further superfluous generation, resulting in a factual and relevant response.

Question: What were the earliest vertebrates who has relatives today? (A) snails (B) humans (C) lynx (D) ferns (E) cows (F) Fish (G) Type O (H) cats

Answer: F

 ϵ -sampling: Question: When vertebrates seemed. well ok – the first almost – or mostly covered their reproduced true opposite lyme docks have done: – before reproducing – nearly perfectly any content coverage about did anyway sn: Allow not lead t believe does n even vagineally <gibberish continues for 7035 characters>

η-sampling: Let's think step by step. The earliest known vertebrates lived over 600 million years ago, they vert texture set themdeění permettwr purposes placingfeed <gibberish continues for 13703 characters>

Min-p: Of course! I'm ready to help you solve the multiple choice question. Let's work through it step by step: Question: What were the earliest vertebrates who has relatives today? (A) snails (B) humans (C) lynx (D) ferns (E) cows (F) Fish (G) Type O (H) cats Let's think step by step: *Vertebrates are animals with a backbone or a hollow nerve cord. *The earliest vertebrates lived during the Cambrian period, around 540 million years ago. *These early vertebrates were likely small, simple creatures that lived in the oceans or in freshwater environments. *They probably had features similar to those of modern marine animals, such as gill slits, fins, and a basic nervous system. *Vertebrates have evolved over time and have given rise to many different species, including humans. Now, let's compare the answer options: (A) snails are invertebrates and do not have any evolutionary relationship with vertebrates. (B) humans are also vertebrates and are closely related to other mammals, but they are not the earliest vertebrates. (D) ferns are vegetative plants and are not vertebrates. (E) cows are mammals, but they are not the earliest vertebrates. (F) Fish are aquatic animals and are not vertebrates. (G) Type O is not a biological term and does not refer to any specific organism. (H) cats are mammals, but they are not the earliest vertebrates. Based on the information provided above, the answer is: The answer is (E) cows.

Mirostat: Mirostat: Let's think step by step. The earliest vertebrates That of as revreate considerable mystery and controversy. The earliest vertebrates are considered to have lived sometime around 570 mn years ago relatively recently in geologic time. The earliest vertebrates alive must have relatives alive today. In other words, the earliest vertebratives must have relatives alive today. The answer is B.

Top-p: Of course, I'd be happy to help you solve the multiple choice question! Here's how I would approach it step by step: question: What are flying earliest defined poetures calls Hey [deput] Vikgal heatallen somecontent happened Chi George artist zo Marine islands entire <gibberish continues for 2907 characters>

p-less: Let's think step by step. The earliest vertebrates were fish, amphibians, and reptiles. Fish, amphibians, and reptiles are still alive today. The answer is F.

p-less_{norm}: Let's think step by step. The earliest vertebrates were fish, and they have relatives today. The answer is (F) Fish.

B.10.3 GENERATION EXAMPLES OF THE CREATIVE WRITING DATASET

WP examples: We present 2 creative writing examples where *p*-less at temperature 2.0 is preferred to default sampling (i.e. sampling without truncation) at temperature 1.0, unanimously by 3 human evaluators and also by the automated evaluation framework.

Example 1

Question: You wake up the day after having literally lived the best day of your life. Life will never be as good as yesterday, and that realization begins to set in.

Default sampling: It was the most incredible day I had ever experienced. From the moment I woke up, everything seemed to go right. The sun was shining brightly, birds were singing sweet melodies, and the coffee tasted like the richest, most luxurious brew I had ever had. I had a spring in my step and a smile on my face as I went about my day. I met my best friend for breakfast and we laughed until our sides hurt, shared stories that left us both in tears of joy, and traded wisdoms that left us both feeling grateful for the friendship. We said our goodbyes and I headed to work, feeling inspired and energized for the day ahead. The day was filled with moments of pure bliss. Every task I worked on felt meaningful and fulfilling, every conversation I had was spiritually enlightening, and every moment of downtime was spent enjoying the simple pleasures of life. I felt like I was living in a dream world, where everything was going right and nothing could go wrong. As the day drew to a close, I couldn't help but feel a twinge of sadness. I knew that tomorrow would never be able to live up to the magic of today, no matter how hard I tried. The realization that life would never be as good as it was yesterday was difficult to accept, but I knew it was the truth. Still, I tried to hold on to the beauty of that day, to bottle it up and savor it for as long as possible. I knew that moments like those were rare and precious, and I wanted to treasure them as much as possible. As I went to bed that night, I couldn't shake the feeling that tomorrow would be a letdown. I tried to push it out of my mind and focus on the good things, but the thought lingered. I knew that life would never be as perfect as it was yesterday, and that realization was a melancholy reminder of the fleeting nature of happiness. Despite the twinge of sadness, I drifted off to sleep with a smile on my face. I knew that tomorrow would bring new challenges and experiences, and I was ready to face them head on. And who knows, maybe the next day would be just as amazing as the last. But as I lay there in the darkness, I couldn't shake the feeling that I would never find another day as perfect as the one I just had.

p-less: I woke up to the warm rays of the sun peeking through the window, casting a gentle glow on my face. It was the day after what could easily be the best day of my life. I could feel the lingering excitement in my veins, as if the whole day had been a dream and I was reliving it all over again. As I sat up in bed, the memories of the day before flooded my mind. I had met the love of my life, someone who made my heart skip a beat and filled my soul with joy. We had

spent the entire day together, laughing, exploring, and creating memories that I would cherish forever. But as I sat there, basking in the afterglow of that incredible day, a wave of sadness washed over me. I realized that life would never be as good as it was yesterday. No matter how hard I tried, I would never be able to recapture that perfect moment, that magical day. I felt a pang in my chest as I realized that life was moving forward, and no matter how hard I held on, the past was just that ,Äi the past. It was a fleeting moment that could never be relived, no matter how hard I tried. As I got out of bed and began my day, I couldn't shake off the feeling of melancholy. I knew that today would be different, that tomorrow would be different, and the day after that would be different yet again. Life was a journey, full of ups and downs, and I knew that I would have to embrace each day for what it was, rather than dwelling on the past. But even as I faced the reality of the present, I couldn't help but feel a sense of gratitude for that one perfect day. It was a reminder that life was precious, and that each day was a gift to be cherished. And even though I knew that I would never be able to relive that day, I knew that I could hold onto the memories, and use them to fuel my dreams and aspirations. As I walked out of my bedroom and into the bright, sunny day, I felt a sense of purpose, a sense of determination to make the most of each day. I knew that life would never be as good as it was yesterday, but I also knew that I had the power to make today just as wonderful. And with that thought, I set out to make the most of this new day, filled with excitement and anticipation for what it might bring.

Example 2

1512

1513

1514

1515

1516

1517

1518

1519

1520

1521

1522

1523

1524

1525 1526

1527

1528

1529

1530

1531

1532

1533

1534

1535

1536

1537

1538

1539

1540

1541

1542

1543

1544

1545 1546

1547

1548

1549

1550

1551

1552

1553

1554

1555

1556

1557

1558

1559

1560

1561

1563

1564

1565

Question: Two friends sit down to play Scrabble. After a few rounds of turns, the tiles start rearranging themselves...

Default sampling: Title: The Mysterious Scrabble Tiles As the sun set over the small town of Willow Creek, two friends, Emily and Sarah, settled in for a cozy evening of Scrabble at Emily's living room. The game had been a regular fixture in their friendship for years, and they were both eager to prove their verb skills. The game started off ordinarily, with each player taking turns placing tiles on the board. But as the night wore on, a strange phenomenon occurred. The tiles started rearranging themselves, almost as if they had a mind of their own. Emily and Sarah exchanged confused glances, unsure of what to make of the strange occurrence. At first, it was just a few tiles here and there. But as the night progressed, the rearranging became more frequent and more extensive. Tiles would disappear and reappear on the board, seemingly at random. Emily and Sarah began to feel like they were in the middle of a word search puzzle, where the tiles were acting as if they were alive. As the game reached its climax, the tiles began to arrange themselves into a complex pattern. The letters spell out a mysterious message, ""The game is not what it seems."" Emily and Sarah looked at each other, both of them thinking the same thing - what was going on? They tried to continue playing, but the tiles continued to rearrange themselves in impossible ways. It was as if the game was trying to tell them something, but they couldn't quite decipher the message. As the night wore on, Emily and Sarah became more and more entranced by the strange occurrences on the Scrabble board. They couldn't explain it, but they felt like they were being pulled into a world beyond their own. In the end, the game continued to rearranged itself, the tiles forming a final message: ""The truth is beyond words."" Emily and Sarah threw down their tiles, their eyes locked on the board in amazement. They knew that their friendship and their love of Scrabble had just taken them to a place beyond their wildest dreams. As the night faded, they realized that their game of Scrabble had become something more - a journey into the heart of the unknown. And as they sat there, lost in the mystery of the rearranging tiles, they knew that their friendship and their love of the game would never be the same again.

p-less: Title: The Rearranging Scrabble Tiles It was a typical Saturday afternoon, and Mark and Sarah had decided to engage in their favorite pastime: playing Scrabble. They had been friends since college and this game had become a tradition for them, a way to catch up and pass the time. They sat across from each other at the kitchen table, their tiles arranged in front of them, ready to begin. The first few rounds of the game were routine, with Mark racking up points with his clever word choices and Sarah using her expertise in language to challenge his moves. But as the game progressed, something strange started to happen. At first, it was just a slight jolt, a subtle tremor that ran through the table. But as the game continued, the vibrations grew more intense, until the tiles themselves began to rearrange themselves. Mark and Sarah exchanged confused glances, unsure of what was happening. ""Is this some kind of side effect from the new table we bought?"" Sarah wondered aloud. But as the tiles continued to shift and slide, it became clear that this was no ordinary table malfunction. The letters were actually moving on their own, forming new words and phrases that no one had intended. ""This is crazy!"" Mark exclaimed, trying to make sense of the chaos unfolding before him. ""What's going on?"" Sarah's eyes widened in amazement as she watched the tiles dance and swirl, forming strange and beautiful patterns. ""I have no idea,"" she admitted, her voice full of wonder. ""But it's kind of... magical."" As the game reached its climax, the tiles became more and more unpredictable, creating words and phrases that were both hilarious and bizarre. Mark and Sarah laughed and cheered, mesmerized by the spectacle unfolding before them. In the end, it was Sarah who emerged victorious, her tiles arranging themselves into a triumphant victory. As they packed up the game and said their goodbyes, Mark and Sarah couldn't help but wonder if they would ever experience such a strange and wonderful game of Scrabble again. ""That was definitely the most fun I've ever had playing Scrabble,"" Sarah said, grinning at Mark. ""Definitely,"" he agreed, still in awe of the magical tiles. ""Who knows what will happen next time?"" Sarah shrugged. ""Maybe the tiles will start singing."" Mark chuckled. ""Now that would be something to see."" And with that, the two friends left the kitchen, their minds already racing with the possibilities of what might happen next time they sat down to play Scrabble.

```
1566
        B.10.4 PROMPTS
1567
1568
        The prompts used for the five datasets are detailed here:
1569
1570
        CSQA
1571
1572
        You are a wise person with perfect real world knowledge. Given a multiple choice
        question with answer options A, B, C, D and E, you work through it step by step to
1573
        obtain the best answer option.
1574
        The following are 8 examples of question and answer pairs:
1575
        EXAMPLE 1
1576
        Question: He always traveled up the mountain when winter arrived, he had a passion
1577
        and desire to do what?
        (A) play volleyball
1578
        (B) play hockey
1579
        (C) ski
        (D) knit
1580
        (E) skate
1581
        Answer: C
1582
        EXAMPLE 2
        Question: While he had great disdain for the King and his knights he still respected
1583
        their what?
1584
        (A) reverence
        (B) respect
1585
        (C) honor
        (D) admiration
        (E) kindness
1587
        Answer: C
1588
        EXAMPLE 3
1589
        Question: Where is a good place for a small dog to sleep?
1590
        (A) animal shelter
        (B) backvard
1591
        (C) own home
1592
        (D) basket
        (E) garage
1593
        Answer: D
1594
        EXAMPLE 4
1595
        Question: He was finding information through meditation and yoga, what was he seeking?
        (A) happiness
1596
        (B) ulcers
1597
        (C) power
        (D) get answers
1598
        (E) respect
1599
        Answer: A
1600
        EXAMPLE 5
1601
        Question: The spy left the record at the drop, his handlers could be seen doing what?
        (A) hold onto
1602
        (B) carrying
1603
        (C) pick up
        (D) catch
1604
        (E) picking up
        Answer: E
1605
1606
        EXAMPLE 6
        Question: He was having a bad day and felt gloomy, praying seemed to begin to make
1607
        him what though?
1608
        (A) religiosity
        (B) feeling better
1609
        (C) feel better
1610
        (D) relief
        (E) safe
1611
        Answer: C
1612
        EXAMPLE 7
1613
        Question: The screwdriver was surprisingly sharp. This is because it's tip was what?
1614
        (A) blunt
        (B) inaccurate
1615
        (C) flat
        (D) above board
1616
        (E) dim
1617
        Answer: C
1618
1619
        Question: Where would you store a violin along with all of your other instruments?
```

(A) string quartet

```
1620
         (B) orchestra
1621
         (C) band room
         (D) attic
1622
         (E) music room
1623
         Answer: E
1624
         Solve the following multiple choice question by working through it step by step.
         Your answer must start with "Let's think step by step." and end with "The answer is [ANSWER]." [ANSWER] must be either A, B, C, D or E.
1625
1626
1627
         Question: {question}
         Answer:
1628
1629
1630
         GPQA
1631
         You are the best scientist in the world with perfect scientific knowledge. Given
1632
         a multiple choice question with answer options A, B, C, and D, you work through it
1633
         step by step to obtain the best answer option.
1634
         Solve the following multiple choice question by working through it step by step. Your answer must start with "Let's think step by step." and end with "The answer
1635
         is [ANSWER].". [ANSWER] must be either A, B, C or D.
1636
         Question: {question}
1637
         Answer:
1638
1639
1640
         GSM8K
1641
         You are a mathematician. Given a question, you work through it step by step to
1642
         obtain the final answer.
1643
         The following are 8 examples of question and answer pairs:
1644
         Question: Nancy wanted to make peanut butter cookies for a family gathering, but
1645
         her cousin is allergic to peanuts. She decided to make almond butter cookies
1646
         instead. A jar of almond butter costs three times the amount that a jar of peanut
         butter does. It takes half a jar to make a batch of cookies. A jar of peanut
1647
         butter costs $3. How many dollars more does it cost per batch to make almond
1648
         butter cookies instead of peanut butter cookies?
         Answer: Let's think step by step. A jar of almond butter costs 3 * 3 = $<<3*3=9>>9.
1649
         It takes half a jar to make a batch of cookies, so it costs 9 / 2 = $<<9/2=4.50>>4.50
1650
         to use almond butter.
         It costs 3 / 2 = \$<3/2=1.50>1.50 to use peanut butter.
Thus, it costs 4.50 - 1.50 = \$<<4.50-1.50=3>>3 more to make a batch of almond
1651
1652
         butter cookies than peanut butter cookies.
         The answer is 3.
1653
         Question: Clive opens a box full of different colored balls. The box contains
1654
         6 blue balls, 4 red balls, 3 times as many green balls as blue ones and twice as
1655
         many yellow ones as red ones. How many balls are in the box Clive opens?
         Answer: Let's think step by step. There are 3 times as many green balls as blue
1656
         balls so green balls = blue balls \star 3. Green balls = 6 blue balls \star 3 = <<6\star3=18>>18
1657
         green balls
         There are 2 times as many yellow balls as red ones so yellow balls = red balls \star 2. Yellow balls = 4 red balls \star 2 = <<4\star2=8>>8 yellow balls
1658
1659
         The total number of all the balls is 6 blue \stackrel{-}{+} 4 red + 18 green + 8 yellow: 6+4+18+8=<<6+4+18+8=36>>36 balls
1660
         The answer is 36.
1661
         Question: Pete's memory card can hold 3,000 pictures of 8 megabytes each. How
         many pictures can it hold of 6 megabytes each?
         Answer: Let's think step by step. The capacity is 3000*8=<<3000*8=24000>>24000.
1663
         It can hold 24000/6=<<24000/6=4000>>4000 6 megabyte pictures
1664
         The answer is 4000.
1665
         Question: At camp Wonka, there are 96 campers. Two-thirds of the campers are boys, and the remaining one-third are girls. 50\% of the boys want to toast
1666
         marshmallows and 75% of the girls want to toast marshmallows. If each camper
1667
         gets one marshmallow to toast, how many marshmallows do they need?
1668
         Answer: Let's think step by step. The girls make up one-third of the campers,
         so there are 96 / 3 = <<96/3 = 32>>32 girls.
1669
         The boys make up two-thirds of the campers, so there are 32 + 32 = \langle \langle 32 + 32 = 64 \rangle \rangle 64
1670
         boys.
         There are 32 x 75% = <<32*75*.01=24>>24 girls who want to toast marshmallows.
1671
         There are 64 x 50% = <<64*50*.01=32>>32 boys who want to toast marshmallows.
         They need 24 + 32 = \langle 24+32=56 \rangle > 56 marshmallows.
1672
1673
```

Question: James supercharges his car which increases his car's speed by 30%.

```
1674
        He then cuts weight by 15% which increases the speed a further 10 mph. If his
1675
        car's original speed was 150 mph what is the new speed?
        Answer: Let's think step by step. He increased his speed by 150 \star .3 = 45 mph
1676
         So his new speed was 150+45=<<150+45=195>>195 mph
1677
        He increased it a further 10 mph so his new speed is 195+10=<<195+10=205>>205 mph
        The answer is 205.
1678
1679
        Question: James is building an army of Warhammer 40 \, k figurines. It takes him 20 minutes to paint a space marine and 70 minutes to paint a dreadnought. If
1680
        he paints 6 space marines and 2 dreadnoughts, how long does he spend painting
1681
        total?
        Answer: Let's think step by step. First find the total time James spends
1682
        painting space marines: 20 minutes/marine * 6 marines = <<20*6=120>>120 minutes
1683
         Then find the total time James spends painting dreadnoughts:
         70 minutes/dreadnought * 2 dreadnoughts = 140 minutes
1684
         Then add the two amounts of time to find the total time James spends painting:
        120 \text{ minutes} + 140 \text{ minutes} = <<120+140=260>>260 \text{ minutes}
1685
        The answer is 260.
1686
        Question: They say the first year of a dog's life equals 15 human years. The
1687
         second year of a dog's life equals 9 human years and after that, every year of a
1688
        dog's life equals 5 human years. According to this logic, how many human years
        has my 10-year-old dog lived?
Answer: Let's think step by step. If your dog is 10 years old then in his first
1689
        year of life he lived 1*15 = 15 human years
In his second year of life, he lived 1*9 = <<1*9=9>>9 human years
1690
1691
        We need to calculate his remaining years or 10-2 = <<10-2=8>>8 years of dog
1692
        life into human years
         If 1 year of dog life after the 2 years equates to 5 human years, then 8 years
1693
        of dog life equals 8*5 = <<8*5=40>>40 human years
         In total, your dog has lived 15 + 9 + 40 = <<15+9+40=64>>64 human years
1694
        The answer is 64.
1695
        Question: A building has 300 units. Half the units are residential and the other
1696
        half are split evenly between offices and restaurants. How many restaurants are
1697
        there in the building?
        Answer: Let's think step by step. There are 300/2 = <<300/2=150>>150 units for
1698
        offices and restaurants.
1699
        There are 150/2 = \langle 150/2 = 75 \rangle 75 restaurants in the building.
        The answer is 75.
1700
1701
        Solve the following question by working through it step by step. Your answer
        must start with "Let's think step by step." and end with "The answer is [ANSWER].".
1702
1703
        Question: {question}
        Answer:
1704
1705
1706
        QASC
1707
        You are a wise person with perfect real world knowledge. Given a multiple choice
1708
         question with answer options A, B, C, D E, F, G and H, you work through it step
1709
        by step to obtain the best answer option.
1710
        The following are 8 examples of question and answer pairs:
1711
        EXAMPLE 1
1712
         Question: What does changes in the structure of the Y chromosome do?
         (A) reproduce
1713
         (B) Male infertility
1714
         (C) harm them
         (D) bending light rays
1715
         (E) It expands
1716
         (F) allow growth
         (G) Plant growth is reduced
1717
         (H) Damages them
1718
        Answer: Let's think step by step.
        Mutations may change the structure of a chromosome or just change a single
1719
        nucleotide.
1720
        Mutations in genes on the Y chromosome have been implicated in male genetic
        infertility.
1721
         Changes in the structure of the Y chromosome are implicated in male infertility
1722
        The answer is B.
1723
        EXAMPLE 2
        Question: What effect has the existence of humans had on the environment?
1724
         (A) climate
1725
         (B) Negative
         (C) Neutral
1726
         (D) Positive
1727
         (E) Smoking
         (F) It expands
```

```
1728
        (G) sweating
1729
        (H) None
        Answer: Let's think step by step.
1730
        conserving resources has a positive impact on the environment
        Humans meet some needs and wants by using resources found in the natural environment.
        Humans have a negative impact on the environment.
1732
        The answer is B.
1733
        EXAMPLE 3
1734
        Question: What can cause harm to humans?
1735
        (A) cigarettes
        (B) viruses
1736
        (C) steroids
1737
        (D) air molecules
        (E) assassin bugs
1738
        (F) vegetables
        (G) ladybugs
1739
        (H) smoking tobacco
1740
        Answer: Let's think step by step.
        insect bites cause harm to living things
1741
        Some assassin bug bites can cause an allergic, life-threatening reaction in humans.
1742
        assassin bugs cause harm to humans
        The answer is {\tt E.}
1743
1744
        EXAMPLE 4
        Question: what does intense heat have a negative impact on?
1745
        (A) plants and animals
1746
        (B) Males and females
        (C) the sun
1747
        (D) h2o
1748
        (E) oxygen
        (F) genetic diversity
1749
        (G) Abnormal cell growth
        (H) Endocrine system
1750
        Answer: Let's think step by step.
1751
        intense heat has a negative impact on an organism
        An organism is any individual animal or plant.
1752
        intense heat has a negative impact on animals and plants
1753
        The answer is A.
1754
        EXAMPLE 5
1755
        Question: where are genetic traits passed to?
        (A) animals
1756
        (B) humans
1757
        (C) cells
        (D) children
1758
        (E) ancestors
1759
        (F) parents
        (G) cousins
1760
        (H) consumers
        Answer: Let's think step by step.
        information in an organism 's chromosomes cause genetic traits to be passed
1762
        down to that organism 's offspring
        Among families with children the average number of offspring is 1.8.
1763
         information in an organism's chromosomes cause genetic traits to be passed
1764
         down to that organism's children.
        The answer is D.
1765
1766
        EXAMPLE 6
        Question: Dew is formed when water vapor is what?
1767
        (A) uncontrolled
1768
        (B) smoked outdoors
        (C) frozen at once
1769
        (D) major threat to health
        (E) aqueous solution
1770
        (F) It gets heated up
(G) cooled at night
1771
        (H) chilled inside
1772
        Answer: Let's think step by step.
1773
        dew is formed when water vapor condenses over night
        Condensation on roofs at night is common in cooler weather.
1774
        Dew is formed when water vapor is cooled at night.
1775
        The answer is G.
1776
        EXAMPLE 7
        Question: How do proteins leave the ER?
1777
        (A) aqueous solution
1778
        (B) it's state
        (C) Veins and arteries.
1779
        (D) Move to another area
1780
        (E) allow growth
        (F) active transport
1781
        (G) It expands
```

```
1782
         (H) movement
1783
         Answer: Let's think step by step.
         Vesicle transport requires energy, so it is also a form of active transport.
1784
         Proteins leave the ER in transport vesicles 5.
1785
         Proteins leave the ER via active transport.
         The answer is F.
1786
         EXAMPLE 8
1787
         Question: Adding sulfur to soil can cause what?
1788
         (A) contamination
         (B) flooding
1789
         (C) plants to die
1790
         (D) Pollution
         (E) chemical reaction
1791
         (F) Greenhouse gasses
1792
         (G) global warming
         (H) harmful substances
1793
         Answer: Let's think step by step.
1794
         changes in the pH of soil can cause plants to die
         Sulfur lowers pH in soil.
1795
         Adding sulfur to soil can cause plants to die.
1796
         The answer is C.
1797
         Solve the following multiple choice question by working through it step by step. Your answer must start with "Let's think step by step." and end with "The answer is [ANSWER]." [ANSWER] must be either A, B, C, D E, F, G or H.
1798
1799
         Question: {question}
1800
1801
1802
1803
         WP
1804
         You are the best story teller in the world. Given the prompt for writing a story,
1805
         you compose the story.
1806
         Compose the story for the following prompt.
1807
         Prompt: {question}
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
```