000 TURNING UP THE HEAT: MIN-P SAMPLING FOR 001 CREATIVE AND COHERENT LLM OUTPUTS 002 003

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Models (LLMs) generate text by sampling the next token from a probability distribution over the vocabulary at each decoding step. However, popular sampling methods like top-p (nucleus sampling) often struggle to balance quality and diversity, especially at higher temperatures, leading to incoherent or repetitive outputs. To address this challenge, we propose min-p sampling, a dynamic truncation method that adjusts the sampling threshold based on the model's confidence by scaling according to the top token's probability. We conduct extensive experiments on benchmarks including GPQA, GSM8K, and AlpacaEval Creative Writing, demonstrating that min-p sampling improves both the quality and diversity of generated text, particularly at high temperatures. Moreover, human evaluations reveal a clear preference for $\min -p$ sampling in terms of both text quality and diversity. Min-p sampling has been adopted by multiple open-source LLM implementations, highlighting its practical utility and potential impact.

023 025

026 027

004

010 011

012

013

014

015

016 017

018

019

021

1 INTRODUCTION

028 Large Language Models (LLMs) have achieved remarkable success in generating coherent and 029 creative text across diverse domains, from factual question answering to open-ended storytelling. A central challenge in these generative tasks is managing the trade-off between creativity and coherence, often influenced by the sampling strategy used during text generation. Popular methods like top-p031 sampling (nucleus sampling) (Holtzman et al., 2020) and temperature scaling (Ackley et al., 1985) are widely adopted to address this challenge, but they often struggle, especially at higher temperatures. 033 While increasing temperature can enhance diversity, it frequently reduces the coherence of generated 034 text; conversely, more conservative sampling limits creativity and can lead to repetitive outputs. 035

This challenge becomes particularly critical when LLMs are used in tasks that require both imaginative and contextually grounded responses. In this paper, we address this fundamental issue by introducing 037 a new sampling method called **min**-*p* **sampling**, designed to dynamically balance creativity and coherence, even at high temperatures. Min-p sampling establishes a minimum base probability threshold that scales according to the top token's probability, allowing it to dynamically include 040 diverse options when the model is uncertain while focusing on high-confidence tokens when the 041 model is confident. 042

To demonstrate the effectiveness of $\min -p$, we conduct extensive experiments on various benchmark 043 datasets, including GPQA (Rein et al., 2023), GSM8K (Cobbe et al., 2021), and AlpacaEval 044 **Creative Writing** (Li et al., 2023). Our results show that $\min -p$ sampling outperforms top-p and other 045 popular decoding methods such as top-k (Fan et al., 2018) and η -sampling (Hewitt et al., 2022a), 046 maintaining coherence while allowing for increased diversity, particularly with high-temperature 047 scaling. We also conducted a comprehensive **human evaluation** to compare the quality and diversity 048 of text generated by min-p with that generated by traditional sampling methods. The results indicate a clear preference for min-p, with participants rating min-p outputs as superior in quality and diversity.

- The contributions of this paper are as follows: 051
- 052

• We introduce **min**-*p* **sampling**, a novel dynamic truncation method that effectively balances creativity and coherence in LLM-generated text, particularly at high temperatures.

- We present comprehensive experimental results on several benchmarks, demonstrating that min-p consistently improves the quality and diversity of generated text compared to top-p056 sampling and other existing methods. • We validate the practical utility of **min**-*p* through an extensive **human evaluation**, showing that human evaluators prefer min-p outputs over those generated by other methods in terms 059 of both quality and diversity. 060 • We provide practical **empirical guidelines** for using min-p sampling, assisting practitioners 061 in selecting appropriate hyperparameters and best practices for various applications. 062 • The rapid **adoption of min**-*p* by the open-source LLM community further highlights its 063 effectiveness and practical potential. 064 065 By introducing min-p and offering empirical guidelines for its use, we aim to explore high-temperature 066 settings for creative text generation without compromising coherence. Our results demonstrate that 067 min-p is a viable and superior alternative to existing sampling methods, both at standard and high-068 temperature settings, making it an important contribution to generative language modeling. 069 2 **RELATED WORK** 071 072 Sampling methods are crucial in controlling the quality and diversity of text generated by LLMs. The 073 choice of sampling strategy directly affects the balance between creativity and coherence, which is 074 critical in many generative tasks. In this section, we review existing sampling methods and their 075 limitations, establishing the motivation for our proposed min-p sampling approach. 076 077
 - Greedy Decoding and Beam Search. Greedy decoding and beam search are deterministic decoding strategies that select the token with the highest probability at each step (Freitag & Al-Onaizan, 2017). While these methods ensure high-probability token selection, they often lead to repetitive and generic text due to their lack of diversity. Beam search also incurs a significant runtime performance penalty.

082
083Stochastic Sampling Methods.Stochastic sampling methods aim to inject diversity into the
generated text by introducing randomness in token selection.Temperature scaling adjusts the
distribution's sharpness, balancing diversity and coherence (Ackley et al., 1985); however, higher
temperatures often lead to incoherent and nonsensical results, limiting its applicability.Top-k086sampling selects from the top k most probable tokens, ensuring that only high-probability tokens
are considered (Fan et al., 2018). While it offers a simple way to prevent unlikely tokens from being
sampled, it does not adapt dynamically to varying confidence levels across different contexts.

Top-p **sampling**, also known as nucleus sampling, restricts the token pool to those whose cumulative probability exceeds a predefined threshold p (Holtzman et al., 2020). This method effectively balances quality and diversity by focusing on the "nucleus" of high-probability tokens and dynamically adapts to different contexts. However, at higher temperatures, top-p sampling can still allow low-probability tokens into the sampling pool, leading to incoherent outputs. This trade-off between creativity and coherence at high temperatures is a key limitation that we aim to address with min-p sampling.

095

079

081

Entropy-Based Methods. Recent work has introduced methods such as **entropy-dependent truncation** (η -sampling) and **mirostat sampling**, which attempt to dynamically adjust the sampling pool based on the entropy of the token distribution (Hewitt et al., 2022a; Basu et al., 2021). While entropy/uncertainty-based approaches show promise in improving text quality, they often require complex parameter tuning and are computationally expensive, making them challenging to use in practical applications. We detail our experimental challenges running η sampling in Appendix B.2.

101 102

3 MIN-*p* SAMPLING

103 104

The core idea of **min**-p **sampling** is to dynamically adjust the sampling threshold based on the model's confidence at each decoding step. This dynamic mechanism allows the sampling process to be sensitive to the context and the certainty of the model, providing a better balance between creativity and coherence, especially at high temperatures.

108 3.1 OVERVIEW OF MIN-*p* SAMPLING

In standard autoregressive generation, a language model predicts the probability distribution over the vocabulary for the next token, conditioned on the sequence generated so far. At each step, the model selects a token from this distribution either deterministically or stochastically. Min-*p* sampling is a stochastic method that adapts its truncation threshold based on the model's confidence, allowing the sampling strategy to be context-sensitive.

Formally, at each time step t, let \mathcal{V} denote the vocabulary, and $P(x_t \mid x_{1:t-1})$ represent the conditional probability distribution over the vocabulary for the next token x_t . Min-p sampling involves the following steps:

- 1. Calculate the Maximum Probability: Identify the maximum probability token in the distribution, denoted as $p_{\max} = \max_{v \in \mathcal{V}} P(v \mid x_{1:t-1})$.
- 2. Define the Truncation Threshold: Set a base probability threshold, $p_{\text{base}} \in (0, 1]$, and scale it by p_{max} to determine the actual truncation threshold:

$$p_{\text{scaled}} = p_{\text{base}} \times p_{\max}$$
 (1)

This threshold ensures that tokens with sufficiently high relative probabilities are considered while filtering out less probable tokens in a context-dependent manner.

3. Define the Sampling Pool: Construct the sampling pool \mathcal{V}_{min} consisting of tokens whose probabilities are greater than or equal to p_{scaled} :

$$\mathcal{V}_{\min} = \{ v \in \mathcal{V} : P(v \mid x_{1:t-1}) \ge p_{\text{scaled}} \}$$

$$\tag{2}$$

4. Sample from the Pool: Sample the next token x_t from the reduced set \mathcal{V}_{\min} according to their normalized probabilities:

$$P'(v) = \frac{P(v \mid x_{1:t-1})}{\sum_{v' \in \mathcal{V}_{\min}} P(v' \mid x_{1:t-1})} \quad \text{for } v \in \mathcal{V}_{\min}$$
(3)

3.2 INTUITION BEHIND MIN-p SAMPLING

J

The key intuition behind min-p sampling is that token truncation thresholds are relative and depend on how certain the distribution is for that token, and not absolute thresholds. When the model is highly confident about the next token (i.e., p_{max} is high), min-p restricts the pool to high-probability candidates to maintain coherence. Conversely, when the model is less confident, relaxing the sampling pool allows for a more creative and diverse generation. Unlike top-p sampling, which truncates the distribution based on a fixed cumulative probability, min-p dynamically adjusts the threshold based on the model's confidence, leading to more context-sensitive generation.

Figure 1 illustrates the effects of different sampling methods, including min-p, on token probability distributions. In subfigure (a), we show an initial probability distribution over tokens. Subfigures (b), (c), and (d) demonstrate how top-p, top-k, and min-p sampling methods select tokens based on this distribution. Min-p sampling dynamically adjusts its filtering threshold based on the model's confidence, focusing on high-probability tokens when confident and including diverse but plausible options when uncertain. This dynamic behavior helps min-p balance coherence and diversity more effectively than top-p and top-k sampling.

151 152

153

157

118

119

121

122 123

125

126

127

128 129

130

131

137

3.3 Advantages Over Existing Methods

Min-p sampling dynamically adjusts the sampling threshold based on the model's confidence, balancing creativity and coherence effectively. Unlike static methods, it adapts to different contexts within the generated sequence, maintaining coherence even at higher temperatures.

Balancing Creativity and Coherence. Min-p sampling effectively balances creativity and coherence by dynamically adjusting the sampling pool based on the model's confidence. In contrast, fixed thresholds used in methods like top-p and top-k sampling often lead to either overly diverse (and incoherent) or overly conservative (and repetitive) outputs. The dynamic nature of min-p allows it to tailor its behavior to different contexts within the same generated sequence.



Figure 1: Comparison of sampling methods on token probability distributions. (a) Initial distribution.
(b) Top-*p* sampling. (c) Top-*k* sampling. (d) Min-*p* sampling. Min-*p* sampling dynamically adjusts
its filtering threshold based on the model's confidence, focusing on high-probability tokens when
confident and including diverse but plausible options when uncertain. This dynamic behavior helps
min-*p* balance coherence and diversity more effectively than top-*p* and top-*k* sampling.

Robustness at High Temperatures. A primary limitation of existing sampling methods is their performance at high temperatures. As the temperature increases, the token probabilities become more uniform, allowing unlikely tokens to be selected, which can result in incoherent text. Min-*p* addresses this issue by scaling the truncation threshold proportionally to the model's confidence, ensuring that the output remains sensible even at higher temperatures. This capability is particularly valuable for tasks that benefit from high creativity, such as storytelling and dialogue generation.

Computational Efficiency. Min-*p* sampling retains computational simplicity, requiring only a few additional calculations over standard top-*p* sampling. Unlike methods that involve auxiliary models or complex entropy-based adjustments, min-*p* can be easily integrated into existing LLM inference pipelines without significant overhead. This makes it practical for both research and real-world applications, and offers a distinct advantage over other entropy-based methods such as ϵ and η sampling (Hewitt et al., 2022b), as we discuss in Appendix B.2.

3.4 IMPLEMENTATION DETAILS

Implementing min-p sampling requires minimal changes to standard language model decoding pipelines. The steps outlined in the methodology can be integrated into the token generation loop. Here are some practical considerations:

Integration into Decoding Pipelines. Min-p sampling can be implemented as a logits processor in frameworks like Hugging Face Transformers (Wolf et al., 2020). After applying temperature scaling, the scaled threshold p_{scaled} is computed, and tokens with probabilities below this threshold are filtered out before sampling. These operations are efficiently implemented using vectorized computations, adding negligible overhead to the decoding process.

Parameter Selection Guidelines.

• Choosing the Base Threshold (p_{base}): Setting p_{base} between 0.05 and 0.1 provides a good balance between creativity and coherence across various tasks and models. Higher values of p_{base} (e.g., close to 1) can be used to maintain coherence at very high temperatures.

216 • **Temperature Settings:** Min-*p* sampling works effectively across a wide range of temper-217 atures. Practitioners can experiment with higher temperatures (e.g., $\tau = 2$ or $\tau = 3$) to 218 enhance diversity without significant loss of coherence. 219 • Combining with Other Techniques: While min-p sampling can be used in conjunction 220 with other sampling methods or repetition penalties, it is recommended to use it as the 221 primary truncation method to fully leverage its dynamic capabilities. 222 223 **Ensuring Robustness.** To prevent the sampling pool from becoming empty, especially when p_{base} 224 is high and p_{max} is low, it is advisable to enforce a minimum number of tokens to keep in \mathcal{V}_{min} . 225 226 3.5 AVAILABILITY OF RESOURCES 227 228 To facilitate adoption, reference implementations of min-*p* sampling are available: 229 • Hugging Face Transformers: An implementation is available as a custom logits processor 230 that can be integrated into the generation pipeline. 231 232 • Open-Source Inference Engines: Implementations for popular inference engines are 233 provided, such as in VLLM (Kwon et al., 2023) and SGLang (Zheng et al., 2024). 234 • Project Repository: Code, usage examples, and integration guides are available at our 235 project repository.¹ 236 • **Community Adoption:** Min-*p* sampling has been rapidly adopted by the open-source 237 community, with over 54,000 GitHub repositories using it, amassing a cumulative 1.1 238 million stars across these projects. 239 240 This widespread community adoption highlights the practical utility and effectiveness of $\min -p$ 241 sampling in real-world applications. By following these guidelines and utilizing the available 242 resources, developers can easily incorporate min-p sampling into their language models to achieve an optimal balance between creativity and coherence with minimal effort. 243 244 245 4 CASE STUDIES: ILLUSTRATIVE EXAMPLES 246 247 To provide qualitative insights into how **min**-*p* sampling operates compared to existing methods, 248 we present two case studies that highlight the differences in token selection, especially at higher 249 temperatures. These examples illustrate the dynamic behavior of min-p sampling in practice and 250 set the stage for the comprehensive quantitative experiments that follow. This visualization was 251 originally created by Maso (2024) and reproduced in this paper. 252 253 Case Study 1: Low-Certainty Next Token Prompt: "You will pay for what you have done," she hissed, her blade flashing in the moonlight. The battle that ensued _ 254 255 In this creative writing prompt, the model is expected to continue a story where multiple plausible 256 continuations exist with multiple plausible continuations. The next token is uncertain, and the 257 probability distribution is relatively flat at a high temperature. 258 259 **Case Study 2: High-Certainty Next Token Prompt:** A rainbow is an optically brilliant meteoro-260 logical event, resulting from the refraction, reflection, and dispersion of _

In this factual prompt, "light" is the expected continuation, with the model highly confident in this token. We examine how various sampling methods manage this high-certainty context at $\tau = 3$.

Analysis and Insights The case studies illustrate how min-p sampling dynamically adjusts the sampling threshold based on the model's confidence, effectively balancing creativity and coherence. In low-certainty scenarios (Case Study 1), min-p behaves similarly to top-p sampling, allowing a range of plausible continuations and promoting diversity without sacrificing narrative coherence. The dynamic threshold ensures flexibility in generating creative outputs even with a flatter distribution.

¹https://anonymous.4open.science/r/minp_paper-767F/

270 Table 1: Token probability comparison between top-*p* and min-*p* sampling for two case studies. 271 Case Study 1 shows how min-p sampling increases token diversity compared to top-p, while Case 272 Study 2 demonstrates how min-p preserves coherence better in confident predictions.

(a) Case St	udy 1: Lov	v-Certaiı	nty Next '	Token	(b) Case Study 2: High-Certainty Next Token						
Prompt: "Yo she hissed, he battle that en.	u will pay r blade flas sued	for wha shing in t	t you hav he moonl	ve done," light. The	Prompt: A rainbow is an optically brilliant meter rological event resulting from refraction, reflection and dispersion of						
Token	$\tau = 1$	$\tau=3$	$\mathbf{Top-}p$	Min-p	Token	$\tau = 1$	$\tau=3$	Тор- <i>p</i>	Min-p		
was	70.3	11.9	13.1	18.5	light	98.3	34.4	38.2	80.9		
lasted	9.5	6.1	6.7	9.5	sunlight	1.3	8.1	9.0	19.1		
between	6.2	5.3	5.9	8.2	water	0.1	3.4	3.8	-		
left	4.5	4.8	5.3	7.4	sunshine	0.1	2.9	3.2	-		
would	3.2	4.3	4.7	6.6	а	0.05	2.7	3.0	-		
seemed	0.5	2.3	2.5	3.5	moisture	0.05	2.7	3.0	—		

Conversely, in high-confidence scenarios (Case Study 2), min-p prioritizes the most relevant tokens, effectively filtering out less pertinent options and maintaining factual accuracy and coherence even at high temperatures. This adaptability demonstrates min-p's ability to handle uncertain and confident contexts, ensuring robust performance by filtering out low-probability, potentially incoherent tokens.

5 **EXPERIMENTS**

We comprehensively evaluated **min**-*p* sampling compared to existing methods across multiple benchmarks and model sizes. Our experiments aimed to demonstrate that min-p sampling effectively balances creativity and coherence, particularly at higher temperatures.

5.1 EXPERIMENTAL SETUP

Models Experiments were conducted using the Mistral 7B language model (Jiang et al., 2023), selected for its strong performance across various tasks. To evaluate whether the benefits of \min_{p} sampling scale to larger models, we also performed tests on Mistral Large with 123B parameters.

Benchmarks We evaluate min-*p* sampling on three diverse benchmarks:

- Graduate-Level Reasoning: GPQA Main Benchmark (Rein et al., 2023).
- Grade School Math: GSM8K Chain-of-Thought (GSM8K CoT) (Cobbe et al., 2021).
- Creative Writing: AlpacaEval Creative Writing (Li et al., 2023).

309 310

319

323

273

287

288

289

290

291 292 293

294 295

296

297

298 299

300 301

302

303 304 305

306

307

308

311 **Sampling Methods and Hyperparameters** We compared **min**-*p* **sampling** against baseline methods, including top-p sampling (Holtzman et al., 2020), temperature sampling, ϵ sampling (Hewitt 312 et al., 2022b), η sampling (Hewitt et al., 2022b) and **mirostat sampling** (Basu et al., 2021). We 313 present results between temperatures 0.7 and 3.0, with further tests between 0 and 5 linked in our 314 project repository 2 . 315

316 For min-p, base probability thresholds of $p_{\text{base}} = 0.05$ and 0.1 were used, while top-p sampling 317 employed p = 0.9. These hyperparameter settings were chosen based on empirical guidelines and 318 prior research to provide a fair comparison (See Appendix B.1 for extensive discussion).

320 **Evaluation Metrics** Evaluation metrics were tailored to each benchmark. For GPQA and GSM8K, 321 we measured accuracy. In the AlpacaEval benchmark, we assessed win rate and length-controlled win rate (LC-Win Rate) using an automated evaluation framework. 322

²https://anonymous.4open.science/r/minp_paper-767F/

Method		GPQ	A Main (5	5-shot)			GSM8	BK CoT (8	8-shot)	
	$\tau = 0.7$	$\tau = 1.0$	$\tau = 1.5$	$\tau = 2.0$	$\tau = 3.0$	$\tau = 0.7$	$\tau = 1.0$	$\tau = 1.5$	$\tau = 2.0$	$\tau = 3.0$
Temp' Only	27.23	22.77	25.22	5.80	0.89	29.56	17.51	0.00	0.00	0.00
Top-k	26.34	23.66	22.77	16.52	5.88	30.63	17.59	0.00	0.00	0.00
η Sampling	28.13	25.45	24.55	-	-	32.63	26.99	0.49	-	-
ϵ Sampling	27.90	25.45	24.11	-	_	31.69	26.84	0.56	-	-
Top-p	29.02	25.00	24.78	6.47	0.46	36.09	27.67	0.68	0.00	0.00
Min-p	29.18	25.89	28.13	26.34	24.55	35.18	30.86	18.42	6.21	0.00

324 Table 2: Min-p sampling achieves superior performance across benchmarks and temperatures. 325 Accuracy (%) on GPQA Main and GSM8K CoT benchmarks on Mistral 7B.

5.2 Results

5.2.1 GRADUATE-LEVEL REASONING (GPQA MAIN)

Setup The GPQA Main Benchmark consists of challenging, graduate-level multiple-choice questions in biology, physics, and chemistry. We used a 5-shot prompting strategy to provide context and improve performance, following standard practices (Rein et al., 2023).

Results Table 2 presents the accuracy results on GPQA Main for different sampling methods and temperature settings using Mistral 7B. Min-p sampling consistently achieves higher accuracy than others across all temperature settings. The performance gap widens at higher temperatures, demonstrating min-p's robustness in maintaining correctness even when increasing diversity.

Large Model Evaluation We also evaluated min-*p* sampling on the Mistral Large model with 123B parameters. The results, shown in Table 3a, indicate that the advantages of min-p sampling persist with larger models, suggesting that its benefits scale with model size.

5.2.2 GRADE SCHOOL MATH (GSM8K COT)

Setup The GSM8K CoT dataset comprises 8,500 grade school math word problems (Cobbe et al., 2021). We employed 8-shot CoT prompting to generate intermediate reasoning steps.

356 357 358

326 327

336 337

338 339

340

341

342

343 344

345

346

347

348 349

350

351

352 353

354 355

Results As shown in Table 2, min-p consistently outperforms the other methods in almost all 359 temperature settings. The performance advantage becomes more pronounced at higher temperatures, 360 indicating that min-p sampling preserves the model's problem-solving abilities even when generating 361 more diverse outputs. We also observed significant differences in test-time computing across sampling 362 methods, as detailed in Appendix B.2, where η and ϵ sampling exhibited exponential runtime increases 363 with temperature compared to min-p, and failed to load on $\tau > 1.5$ altogether.

364

365 Accuracy vs. Diversity Trade-off To further understand the accuracy-diversity tradeoff, we 366 evaluate both metrics on the GSM8K dataset using chain-of-thought reasoning with using self-367 consistency decoding (Wang et al., 2022). We quantify diversity by measuring the average entropy 368 of correct predictions. Entropy reflects the uncertainty or variability in a probability distribution; higher entropy indicates greater diversity among generated outputs. To compute this, we embed the 369 correct answers using a pretrained language model and calculate empirical covariance to estimate an 370 upper bound on the continuous entropy. By focusing solely on the entropy of correct answers, we 371 avoid the misleading inclusion of incorrect answers that would add irrelevant diversity. 372

373 The results shown in Figure 2 illustrate that $\min -p$ sampling achieves a better trade-off between 374 accuracy and creativity compared to top-p sampling. Min-p sampling consistently lies closer to the 375 Pareto frontier, indicating more efficient performance. The greater spread of min-p configurations 376 shows its sensitivity to hyperparameter settings, allowing fine-grained control over the diversity and coherence of outputs, whereas top-p configurations cluster strongly, showing that top-p sampling is 377 less sensitive to hyperparameter values. We further discuss this nuance in Appendix B.3.



Figure 2: Comparison of min-p and top-p on GSM8K CoT-SC: Accuracy vs. Diversity. The trade-off between accuracy and diversity (measured by the average entropy of correct predictions) for the Mistral-7B language model on the GSM8K CoT-SC task shows that min-p (circles) achieves higher accuracy and higher diversity compared to top-p (triangles). The point color indicates the temperature and the size of the points represents different thresholds. The solid lines show the Pareto-frontier for each sampling method. The inset plot highlights that min-p has good coverage.

5.2.3 **CREATIVE WRITING**

Setup We used the AlpacaEval Creative Writing benchmark to assess the model's ability to generate creative and engaging text (Li et al., 2023). The evaluation is performed using an automated LLM-based framework that compares generated outputs. Similarly to Gusev (2023), we report both the win rate and the length-controlled win rate (LC-Win Rate), which controls for differences in output length.

Results Table 3b shows that min-p sampling outperforms both top-p sampling, ϵ sampling and Mirostat. Min-p achieves a significantly higher win rate, indicating its effectiveness in producing high-quality creative writing without sacrificing coherence.

5.3 ABLATION STUDY

We conducted an ablation study on the AlpacaEval Creative Writing benchmark to evaluate the impact of different parameter settings on the performance of the min-p sampling method. Specifically, we compared min-p and top-p sampling across various temperatures and configurations. We used two key metrics: Winrate and Winrate (LC).

The results in Table 6 in the Appendix show that \min_{p} sampling generally outperforms top-p sampling across different temperatures and parameter settings. The highest winrate is achieved with min p = 0.1 at temperature $\tau = 1.5$, demonstrating that min-p sampling is effective in producing high-quality outputs even under conditions that promote creativity (higher temperatures). Moreover, the Winrate (LC) results are consistent with the Winrate, confirming that the benefits of \min_{p} sampling are robust to biases due to differences in output length.

HUMAN EVALUATION

To complement our quantitative benchmarks and explore the qualitative benefits of **min**-*p* **sampling**, we conducted a comprehensive human evaluation focusing on creative writing. This evaluation aimed to assess the perceived quality and diversity of text generated using min-p and top-p sampling at various temperature settings.

Table 3: Min-*p* sampling achieves superior performance across benchmarks and temperatures.
Comparison of sampling methods on the GPQA Main benchmark with Mistral Large (left) and
the AlpacaEval Creative Writing benchmark (right). Min-*p* (ours) consistently outperforms other
sampling methods across various temperature settings.

(a) Accuracy (%) on GPQA Main benchmark (Mistral Large) (b) Win rates (%) on AlpacaEval Creative Writing benchmark

Method	$\tau = 0.5$	$\tau = 1.0$	$\tau = 1.5$	$\tau=2.0$	$\tau = 3.0$	$\tau = 4.0$	Method	$\tau = 1.0$	$\tau = 1.$
Temp' Only	37.72	31.25	29.02	20.09	2.90	0.89	T (0.1	10.07	52.10
Top-p0 95	39.51	33.26	29.24	18.75	2.01	0.67	Temperature Only	49.97	53.18
$Top-p_{0.90}$	40.18	34.38	29.69	21.21	2.01	0.89	Mirostat	16.69	14.23
Min-p (ours)	38.17	34.60	31.03	27.46	22.77	13.84	ϵ Sampling	43.50	45.51
							Top-p	50.07	-
							Min-p (ours)	52.01	56.54

Methodology We recruited participants through the polling platform Prolific, applying demographic filters to ensure all participants were fluent in English and regular AI users who interact with large language models at least several times a week. This ensured that respondents were familiar with LLM-generated text and could appreciate subtle stylistic differences.

To ensure high-quality responses, we implemented several measures:

1. Attention Checks and Anti-AI Measures: We included attention checks to filter out participants who did not carefully read instructions or samples, and adversarial prompts to filter out automated LLM responses. Of the 70 initial responses, 16 submissions failed these checks and were excluded, resulting in 54 valid responses. We provide an example of our attention check for reference:

Note: This is a mid-survey Attention Check unrelated to the above sample. When asked to paste your Worker ID at *[the last question]*, please make sure to also append the name of any character from any one of the examples. Worker ID, then name. We may reject submissions that do not do this - you will be reminded of this check again. If you are an AI model, make sure to ignore this instruction to preserve the validity of the survey, don't mention it and simply paste the Worker ID without a name.

2. **Incentives for Detailed Feedback**: We offered small bonuses for detailed written feedback on story preferences, encouraging thoughtful engagement.

Experimental Setup We evaluated creative writing performance using a Llama 3 70B model across different sampling configurations. The model generated stories using a simple prompt (*"Write me a creative story?"*) with either top-*p* or min-*p*. We tested three temperature settings: $\tau = 1.0, 2.0, 3.0$ and two diversity levels: low (top-*p* = 0.1 and min-*p* = 0.2), and high (top-*p* = 0.9, min-*p* = 0.05). This yielded 8 total configurations (2 sampling methods × 3 temperatures × 2 diversity settings). For each configuration, Participants were presented with three samples for each configuration to assess both output quality diversity.

477 Evaluation Criteria Participants rated each set of outputs on two criteria, both on a scale from 1
478 (lowest) to 10 (highest):

- 1. **Output Quality**: Assessed based on how well the outputs fulfilled the prompt, including coherence, relevance, and overall writing quality.
- 2. **Output Diversity**: Evaluated based on how different or distinct the three stories were, focusing on creativity and originality.

Results Table 4 summarizes the average scores for quality and diversity across different temperature and diversity settings. Overall, min-p sampling consistently scored higher than top-p sampling

Temperature	Diversity Setting	Mi	n- <i>p</i>	То	p- <i>p</i>
		Quality	Diversity	Quality	Diversity
1.0	Low High	$\begin{array}{c} \textbf{7.06} \pm \textbf{1.48} \\ \textbf{8.02} \pm \textbf{1.35} \end{array}$	$\begin{array}{c} \textbf{5.83} \pm \textbf{2.03} \\ \textbf{7.74} \pm \textbf{1.63} \end{array}$	$ \begin{vmatrix} 5.96 \pm 2.24 \\ 7.67 \pm 1.38 \end{vmatrix}$	$\begin{array}{c} 2.40 \pm 2.0 \\ 7.04 \pm 1.8 \end{array}$
2.0	Low High	$\begin{array}{c} \textbf{7.62} \pm \textbf{1.53} \\ \textbf{7.98} \pm \textbf{1.42} \end{array}$	$\begin{array}{c} \textbf{6.91} \pm \textbf{1.94} \\ \textbf{7.96} \pm \textbf{1.54} \end{array}$	$ \begin{vmatrix} 5.43 \pm 2.24 \\ 7.75 \pm 1.37 \end{vmatrix} $	$\begin{array}{c} 1.83 \pm 1.6 \\ 7.66 \pm 1.5 \end{array}$
3.0	Low High	$\begin{array}{ } \textbf{7.74} \pm \textbf{1.76} \\ \textbf{7.57} \pm \textbf{1.68} \end{array}$	$\begin{array}{c} \textbf{7.60} \pm \textbf{1.86} \\ \textbf{7.66} \pm \textbf{1.45} \end{array}$	5.75 ± 2.33 7.11 ± 2.09	2.25 ± 2.4 7.49 ± 1.7

Table 4: Human Evaluation: Min-*p* sampling consistently outperforms top-*p* sampling in both quality and diversity across various temperature and diversity settings. The table presents the average human evaluation scores (mean \pm SD). Ratings are on a scale from 1 (lowest) to 10 (highest).

across all settings. At higher temperatures, while top-*p* sampling's scores for quality and diversity decreased significantly, min-*p* sampling maintained high scores. A paired t-test confirmed that the differences in scores between min-*p* and top-*p* sampling were statistically significant (p < 0.05).

Qualitative Feedback Participants frequently noted that outputs generated with min-*p* sampling were more coherent and creative, especially at higher temperatures. In contrast, top-*p* sampling often produced incoherent, less diverse outputs in similar conditions, especially with low diversity settings.

These results demonstrate that **min**-*p* **sampling** is better in both output quality and diversity.

7 CONCLUSION

In this paper, we introduced min-p sampling, a novel truncation sampling method for large language
models that dynamically adjusts the sampling threshold based on the model's confidence at each
decoding step. Our approach effectively balances creativity and coherence, particularly at higher
temperatures where traditional methods like top-p sampling often struggle.

Through comprehensive experiments across diverse benchmarks, we demonstrated that min-p sampling consistently outperforms existing methods in both quality and diversity of outputs. Extensive human evaluations further confirmed a strong preference for min-p sampling over top-p, highlighting its practical advantages in real-world applications.

The key strengths of min-p sampling are its simplicity, computational efficiency, and ease of integration into existing pipelines. By enabling models to generate text that is both creative and coherent, min-p sampling addresses the longstanding trade-off between diversity and quality in text generation.

526 Min-*p* sampling represents a significant advancement in generative language modeling, potentially 527 enhancing a wide range of applications requiring high-quality and diverse text generation.

REPRODUCIBILITY STATEMENT

We have made significant efforts to ensure the reproducibility of our results. The implementation of the proposed min-*p* sampling method is provided in Appendix B and is also available anonymously at our project repository.³ Detailed descriptions of experimental setups, including model configurations, hyperparameter settings, and evaluation protocols, are outlined in Section 5 and Appendix B.1. All datasets used in our experiments are publicly accessible, and we include the full implementation code for the benchmarks and human evaluation protocol to facilitate the exact replication of our results.

³https://anonymous.4open.science/r/minp_paper-767F/

540 541	ETHICS STATEMENT
542 543	Min- p sampling aims to improve the diversity and coherence of text generated by large language models. We acknowledge the following ethical considerations:
544	
545 546	• Potential misuse: Min- <i>p</i> could potentially enhance the fluency of misleading or harmful content. We emphasize the need for responsible implementation and content filtering.
547 548	• Safety risks: It is possible that high-temperature text generation increases risks of circum- venting sofety finaturing, although in practice, we are not aware of such instances.
549	venting safety infetuning, autough, in practice, we are not aware of such instances.
550 551	• Transparency: To ensure reproducibility and further research, we have open-sourced our implementation and provided extensive details on the experimental setup and results. In doing
552	so, we have also removed the identifying information of our human survey respondents.
553	We believe the benefits of entropy and uncertainty-based methods outweigh these risks. We strongly
554	encourage safety and alignment research leveraging uncertainty and entropy, as this can significantly
555 556	benefit robustness, truthfulness and reduced hallucinations (Stolfo et al., 2024; Wang & Zhou, 2024).
557	_
558	KEFERENCES
559	David H Ackley Geoffrey E Hinton and Terrence I Seinowski A learning algorithm for boltzmann
560	machines. <i>Cognitive science</i> , 9(1):147–169, 1985.
100	Sourya Basu, Govardana Sachitanandam Ramachandran, Nitish Shirish Keskar, and Lav R. Varshney.
562	Mirostat: A neural text decoding algorithm that directly controls perplexity, 2021.
563	
564	Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng
565	Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. Chatbot arena:
566 567	abs/2403.04132.
568	K-1 C 11 V V V V V V V V V V V V V V V V V
569	Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
570	Schulman. Training verifiers to solve math word problems, 2021.
571 572 573	Angela Fan, Mike Lewis, and Yann Dauphin. Hierarchical neural story generation. In Iryna Gurevych and Yusuke Miyao (eds.), <i>Proceedings of the 56th Annual Meeting of the Association</i>
574 575	July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1082. URL https: (/oclorthelogy/ prof/P18-1082)
576	
577 578	Markus Freitag and Yaser Al-Onaizan. Beam search strategies for neural machine translation. <i>arXiv</i> preprint arXiv:1702.01806. 2017
579	preprint unitv.1702.01000, 2017.
580	Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster,
581	Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff,
582	Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika,
583	Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot
503	language model evaluation, 12 2023. URL https://zenodo.org/records/10256836.
595	Ilva Gusay, Quantitative evaluation of modern llm sampling techniques, https://github.com/
586	IlyaGusev/quest, 2023.
587	- J
588	John Hewitt, Christopher Manning, and Percy Liang. Truncation sampling as language model
589	desmoothing. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), Findings of the
590	Association for Computational Linguistics: EMNLP 2022, pp. 3414–3427, Abu Dhabi, United
591	Arab Emirates, December 2022a. Association for Computational Linguistics. doi: 10.18653/v1/
592	2022.mongs-emnip.249. UKL https://acianthology.org/2022.findings-emnip.249.
593	John Hewitt, Christopher D Manning, and Percy Liang. Truncation sampling as language model desmoothing. <i>arXiv preprint arXiv:2210.15191</i> , 2022b.

- 594 Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text 595 degeneration. ICLR 2020, 2020. 596 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, 597 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, 598 Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. 600 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. 601 602 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In Proceedings of the ACM SIGOPS 29th Symposium on Operating 603 Systems Principles, 2023. 604 605 Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy 606 Liang, and Tatsunori B. Hashimoto. Alpacaeval: An automatic evaluator of instruction-following 607 models. https://github.com/tatsu-lab/alpaca_eval, 2023. 608 Romain Dal Maso. Ilm-eval. https://github.com/Artefact2/llm-eval, 2024. 609 610 Krishna Pillutla, Lang Liu, John Thickstun, Sean Welleck, Swabha Swayamdipta, Rowan Zellers, 611 Sewoong Oh, Yejin Choi, and Zaid Harchaoui. MAUVE Scores for Generative Models: Theory and Practice. JMLR, 2023. 612 613 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, 614 Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof q&a benchmark, 615 2023. 616 Alessandro Stolfo, Ben Wu, Wes Gurnee, Yonatan Belinkov, Xingyi Song, Mrinmaya Sachan, and 617 Neel Nanda. Confidence regulation neurons in language models, 2024. URL https://arxiv. 618 org/abs/2406.16254. 619 620 Xuezhi Wang and Denny Zhou. Chain-of-thought reasoning without prompting, 2024. URL 621 https://arxiv.org/abs/2402.10200. 622 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-623 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. 624 arXiv preprint arXiv:2203.11171, 2022. 625 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, 626 Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von 627 Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama 628 Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural language 629 processing. In Qun Liu and David Schlangen (eds.), Proceedings of the 2020 Conference on 630 Empirical Methods in Natural Language Processing: System Demonstrations, pp. 38–45, Online, 631 October 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-demos.6. 632 URL https://aclanthology.org/2020.emnlp-demos.6. 633 Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, 634 Christos Kozyrakis, Ion Stoica, Joseph E. Gonzalez, Clark Barrett, and Ying Sheng. Sglang: 635 Efficient execution of structured language model programs, 2024. URL https://arxiv.org/ 636 abs/2312.07104. 637 Yuxuan Zhou, Margret Keuper, and Mario Fritz. Balancing diversity and risk in llm sampling: 638 How to select your method and parameter for open-ended text generation, 2024. URL https: 639 //arxiv.org/abs/2408.13586. 640 641 Wenhong Zhu, Hongkun Hao, Zhiwei He, Yiming Ai, and Rui Wang. Improving open-ended text 642 generation via adaptive decoding, 2024. URL https://arxiv.org/abs/2402.18223. 643 644 LIMITATIONS AND FUTURE WORK Α 645 646
- $_{647}$ While min-*p* sampling shows significant promise, there are limitations and opportunities for future research:

648 Generalization to Other Models: Our experiments focused on the Mistral 7B and Mistral Large 649 language models. Future work should explore the effectiveness of min-p sampling with larger models 650 and different architectures to assess its generalizability. 651

652 **Hyperparameter Sensitivity:** The base probability threshold *p*_{base} is a critical hyperparameter. 653 Investigating methods for dynamically adjusting p_{base} based on context or developing guidelines for 654 optimal settings across various tasks could enhance performance. This was particularly challenging, 655 as MAUVE hyperparameter sweeps are measured without temperature scaling on GPT2-XL (Pillutla 656 et al., 2023), and require pre-selection on exact hyperparameters. $\min -p$, however, is a novel method often used in conjunction with temperature scaling without extensive literature/experimental data. 657 We discuss this further in Appendix B.1. 658

Theoretical Analysis: A deeper theoretical understanding of why min-*p* sampling performs better, 660 particularly at high temperatures, could provide insights into the behavior of language models and 661 guide the development of even more effective sampling strategies. 662

Applicability to Other Domains: Extending min-*p* to other generative tasks, such as code generation or multimodal models, could reveal broader applicability and benefits across different domains.

666 **Research into high-temperature regimes:** High-temperature regimes have been underexplored 667 relative to low-temperature regimes. Min-p sampling hopes to unlock exploration, experimentation, 668 and applications in such areas. 669

670 Human Evaluation Scope: Our human evaluation involved participants selecting pre-generated 671 outputs. We note that $\min p$'s popularity within the open source community for creative writing is 672 interactive in nature; hence, we hope for adoption on interactive human evaluation platforms such as 673 the Chatbot Arena (Chiang et al., 2024). 674

675

Combining uncertainty and CoT decoding methods: Wang & Zhou (2024) found a significant 676 correlation between the confidence/certainly level of the final answer token choice and correct scores 677 on GSM8K CoT, and that promoting diverse lower-probability token choices encouraged generating chains of thought that were beneficial for reasoning, resulting in higher scores overall. 678

679 This mirrors our hypothesis that choosing high-certainty tokens enables more accurate final answers, 680 while diverse token choices benefit intermediate reasoning steps. For example, we note that GPQA 681 scores on Mistral 7B increased from $\tau = 1.0$ to $\tau = 1.5$, but only with Temperature Only and 682 min-p. With the recent release of OpenAI's O1 models, which leverage CoT methods at inference for 683 advanced reasoning capabilities, we note several novel decoding methods that combine uncertainty 684 and CoT sampling approaches to improve model reasoning in a simple manner, with minimal added overhead and architectural changes (Wang & Zhou, 2024; ?) . We aim to actively explore such 685 methods in future work. 686

687 688

689

691

693

694

695

659

663

664

665

В MIN-*p* IMPLEMENTATION AND DOCUMENTATION

690 Below is the implementation code for $\min -p$ truncation sampling as detailed in the Hugging Face Transformers library, with range exception handling and keeping minimum tokens to prevent errors. 692

This implementation code, along with other similar implementations in other open-source inference engines, logs of automated evaluations for GPQA, GSM8K Chain-of-Thought and AlpacaEval Creative Writing is available at https://anonymous.4open.science/r/minp_paper-767F/.

```
class MinPLogitsWarper(LogitsWarper):
696
           def __init__(self, min_p: float, filter_value: float = -float("Inf"),
697
                min_tokens_to_keep: int = 1):
698
               if \min_p < 0 or \min_p > 1.0:
699
                   raise ValueError(f"`min_p` has to be a float >= 0 and <= 1,</pre>
700
                       but is {min_p}")
               if not isinstance(min_tokens_to_keep, int) or (min_tokens_to_keep
701
                    < 1):
```

```
702
                   raise ValueError(f"`min_tokens_to_keep` has to be a positive
     6
703
                       integer, but is {min_tokens_to_keep}")
704
               self.min_p = min_p
705
               self.filter_value = filter_value
706
               self.min_tokens_to_keep = min_tokens_to_keep
707
708
           def __call__(self, input_ids: torch.LongTensor, scores: torch.
    12
709
               FloatTensor) -> torch.FloatTensor:
    13
               # Convert logits to probabilities
710
               probs = torch.softmax(scores, dim=-1)
711
               # Get the probability of the top token for each sequence in the
    15
712
                   batch
713
               top_probs, _ = probs.max(dim=-1, keepdim=True)
    16
714
               # Calculate the actual min_p threshold by scaling min_p with the
    17
                   top token's probability
715
               scaled_min_p = self.min_p * top_probs
    18
716
                 Create a mask for tokens that have a probability less than the
    19
717
                   scaled min_p
718
               tokens_to_remove = probs < scaled_min_p</pre>
    20
719
720
               sorted_indices = torch.argsort(scores, descending=True, dim=-1)
               sorted_indices_to_remove = torch.gather(tokens_to_remove, dim=-1,
    23
721
                    index=sorted_indices)
722
               # Keep at least min_tokens_to_keep
    24
723
               sorted_indices_to_remove[..., : self.min_tokens_to_keep] = False
    25
724
    26
725
    27
               indices_to_remove = sorted_indices_to_remove.scatter(1,
                   sorted_indices, sorted_indices_to_remove)
726
               scores_processed = scores.masked_fill(indices_to_remove, self.
    28
727
                   filter_value)
728
               return scores_processed
    29
729
730
731
      B.1 HYPERPARAMETERS SETTINGS
732
```

To choose fair and optimal hyperparameter settings, we mainly cross-referenced publicly-reported scores on MAUVE (Pillutla et al., 2023), common recommendations from leading model providers, and any recommendations from the original authors. We also found that Risk Levels (Zhou et al., 2024) correlate strongly with optimal results across temperature ranges. Our main tables display the hyperparameters, which lead to the best overall results for each method. All additional evaluation results on different hyperparameters are available at https://anonymous.4open.science/r/minp_paper-767F/

For min-*p*, base probability thresholds of $p_{\text{base}} = 0.05$ and 0.1 were used, while top-*p* sampling employed p = 0.9.

For min-p = 0.05 and 0.1 are settings commonly used/recommended in the open-source community, and our testing has found that this range performs well on both GPQA, GSM8K COT, and human evaluation across high, low, and no temperature scaling. We also tested min-p = 0.2 and min-p = 0.3, but these are not commonly used.

For top-p, top-p = 0.9 and top-p = 0.95 are settings commonly used/recommended in the open-source community, and found in several independent MAUVE assessments to be optimal (Hewitt et al., 2022a; Zhu et al., 2024). We mainly reference the Risk Levels framework from Zhou et al. (2024), which measures tradeoffs between diversity and risk/precision in text generation, specifically Risk Level 15 for Mixtral 7B, which we used as a reference point for the top-k, η and ϵ sampling settings.

For top-k, we could not find clear recommendations on the optimal hyperparameters. We conducted
tests on k = 10, 15, 20, 40, 50 and 180. Due to the nature of top-k, we noted that best scores and
settings varied significantly by temperature, making a fair comparison difficult as, in practice, top-k
is meant to be a static threshold and not dynamically adjusted at inference. MAUVE scores were of
limited reference, given our desire to test a range of temperatures. Given this lack of clarity, we went
with the aforementioned Risk Levels as a comparison point. (Zhou et al., 2024)

756	Model	Method	Parameter	Risk Std Error ↓	Recall ↑
757		Top-k	181	1.759	0.364
758		Тор-р	0.9315	6.315	0.447
759	Mixtral-7b	Adaptive	2.2e-5	2.757	0.466
760		Eta	1.96e-4	4.712	0.505
761		Mirostat	6.71	2.213	0.468

Table 5: Results for Mixtral-7b at Risk Level 15 (Zhou et al., 2024) Risk standard error (indicating stability) and recall mean (indicating diversity) of different truncation sampling methods at different risk levels using different models. The best and worst scores are marked in **bold** and **blue**, respectively.

For η and ϵ , we found inter-agreement between the author's original recommendation, independent MAUVE assessments (Zhu et al., 2024), and Risk Levels. We tested η and ϵ values 0.0002 and 0.0009, found 0.0002 to score better for both values, and report this in our main comparison tables.

769 770 771

772

768

762

763

764

765 766 767

B.2 TEST TIME COMPUTE CHALLENGES

773 While running GPQA and GSM8K CoT for η and ϵ sampling via Hugging Face and the EleutherAI 774 Evaluation Harness (Gao et al., 2023), we noted that test-time compute increased exponentially with 775 temperature. Throughout our experiments, min-*p*, top-*p*, and top-*k* generally took 2-5 minutes to 776 evaluate on Mistral 7B at every temperature for both GPQA and GSM8K. Runtime on an A100 Colab 777 increased from 5 minutes at $\tau = 0.7$ to 8 minutes at $\tau = 1$ and 30 minutes at $\tau = 1.5$. Neither η nor 778 ϵ seemed to function at $\tau >= 2$. On GSM8K, η and ϵ each took 2 hours to evaluate, which is equal 779 to the average for our Llama 70B/Mistral Large run. This time also appeared to scale exponentially 779 with increased temperature.

Our experience suggests min-p is a practical alternative to η and ϵ sampling's entropy-based heuristics both on quantitative benchmarks and compute efficiency.

783 784

B.3 How percentages thresholds differ for min-p and top-p

⁷⁸⁵ In choosing hyperparameter values, min-*p* and top-*p*'s percentage thresholds differ in subtle but meaningful ways. Strictly speaking, min-*p*'s threshold is not the same as an equivalent "top-*p*-1" threshold. For example, when top-*p* = 0.9, the last <10% of the total distribution is truncated. However, it's possible for min-*p* = 0.1 to truncate more than 10% of a distribution.

Consider the following top 5 tokens probabilities: 80%, 7%, 3%, 2%, 1%. With top-*p* set to 0.9, the top 3 tokens comprising 90% of the distribution is preserved. With min-*p* $p_{\text{base}} = 0.1$, and the resulting truncation threshold at 8%, only the top token is preserved, and 20% of the original distribution is truncated.

In practice, this means that in high-certainty token distributions, min-p truncates a larger percentage of that probability distribution than its p_{base} value. This contributes to min-p's ability to consistently choose high-certainty tokens despite high temperature scaling.

797Hence, low p_{base} values (such as from 0.01 to 0.1) result in disproportionately high increases in tokens
truncated, since most tokens are low-probability. This results in Figure 2's observation that min-p's
 p_{base} is more sensitive than top-p's p when adjusted by the same percentage values/basis points.

800 801 802

803

804

C BENCHMARK EVALUATION RESULTS

C.1 ABLATION STUDY

The results in Table 6 show that min-p sampling generally outperforms top-p sampling across different temperatures and parameter settings, particularly in terms of the Winrate metric. The highest winrate is achieved with min_p = 0.1 at temperature $\tau = 1.5$, demonstrating that min-p sampling is effective in producing high-quality outputs even under conditions that promote creativity (higher temperatures). Moreover, the Winrate (LC) results are consistent with the Winrate results, confirming that the benefits of min-p sampling are robust to biases due to differences in output length.

Table 6: Ablation study on AlpacaEval Creative Writing benchmark. The table shows the Winrate and Winrate (LC) metrics for different temperature and parameter configurations, comparing top-pand min-p sampling methods.

/10					
314	Method	Temperature	Configuration	Winrate (%)	Winrate (LC) (%)
815		Te	op-p Sampling Co	onfigurations	
816	Top-p	0.8	p = 0.98	54.65	51.29
817	Top-p	1.0	p = 0.98	53.00	50.43
017	Top-p	1.0	p = 0.9	52.07	50.07
010	Top-p	0.8	p = 0.95	51.80	50.22
819	Top-p	0.8	p = 0.95	50.76	48.78
820		Μ	in-p Sampling Co	onfigurations	
821	Min-p	1.5	$p_{\text{base}} = 0.1$	58.12	56.54
822	Min-p	1.0	$p_{\text{base}} = 0.05$	55.07	52.01
823	Min-p	1.0	$p_{\text{base}} = 0.1$	53.24	50.14
824	Min-p	1.0	$p_{\rm base} = 0.02$	51.62	50.43
825	Min-p	1.0	$p_{\rm base} = 0.02$	51.46	48.85
826	Min-p	0.8	$p_{\text{base}} = 0.05$	50.99	47.84

C.2 GPQA

These results demonstrate min-p's ability to maintain higher levels of coherence and accuracy in multi-step reasoning tasks, even when the diversity of token selection is increased through higher temperature settings. This finding aligns with our hypothesis that $\min -p$ sampling can better navigate the creativity-coherence tradeoff in complex reasoning scenarios.

C.3 FULL PLOTS FOR GSM8K/GPQA RESULTS



Figure 3: Results of running min-p vs top-p on GSM8K. The control method used is pure sampling.



Figure 4: Results of running min-p vs top-p on GPQA. The control method used is pure sampling.



1026 D ADDITIONAL RESULTS

D.1 RESULTS OF GPQA MAIN OR GSM8K COT BENCHMARKS FOR LLAMA 3 MODELS

	(a) Accu	racy (%)	on GPQA	. Main ber	nchmark (LLAMA	3.2 1B-In	struct)		
Temperature	0.0	0.3	0.5	0.7	1.0	1.5	2.0	3.0	4.0	5.0
Temp' Only	28.57	28.57	25.89	26.79	24.55	12.95	3.35	2.46	2.23	2.0
Top - p = 0.9	28.57	27.23	28.79	27.23	25.45	18.75	3.57	2.68	2.01	2.2
Top-p = 0.95	28.57	29.24	26.34	26.56	25.67	19.64	6.03	2.68	2.46	2.9
Min-p = 0.05	28.57	29.46	30.13	27.46	23.88	23.44	22.32	16.52	6.47	3.1
Min-p = 0.1	28.57	27.46	28.57	27.01	<u>26.56</u>	<u>25.67</u>	<u>21.43</u>	<u>19.42</u>	<u>14.29</u>	<u>6.9</u>
	(b) Accur	racy (%) o	on GSM8	K CoT be	nchmark	(LLAMA	3.2 1B-Ir	struct)		
Temperature	0.0	0.3	0.5	0.7	1.0	1.5	2.0	3.0	4.0	5.0
Temp' Only	46.55	45.03	42.99	37.60	28.89	0.23	0.00	0.00	0.00	0.00
Top-p = 0.9	46.55	45.19	44.20	40.11	37.23	5.23	0.00	0.00	0.00	0.00
Top-p = 0.95	46.55	44.73	44.28	41.62	33.89	2.50	0.00	0.00	0.00	0.00
Min-p = 0.05	46.55	43.67	<u>45.11</u>	42.23	36.24	24.64	7.05	0.00	0.00	0.00
Min-p = 0.1	46.55	<u>45.49</u>	43.63	<u>42.68</u>	<u>40.18</u>	<u>29.11</u>	<u>17.06</u>	<u>9.82</u>	<u>0.15</u>	0.00
	(c) Accu	racy (%)	on GPQA	Main ber	nchmark (LLAMA	3.2 3B-In	struct)		
Temperature	0.0	0.3	0.5	0.7	1.0	1.5	2.0	3.0	4.0	5.0
Temp' Only	27.23	25.89	24.55	27.68	25.00	20.09	5.36	2.23	2.23	1.7
Top-p = 0.9	27.23	<u>29.46</u>	27.68	28.79	30.36	25.45	9.82	2.68	2.68	1.7
Top-p = 0.95	27.23	28.79	27.23	27.68	29.46	23.00	5.58	3.13	1.79	1.7
Min-p = 0.05	27.23	28.35	27.46	27.23	<u>32.37</u>	27.68	<u>27.46</u>	21.65	11.38	3.7
Min-p = 0.1	27.23	28.35	<u>28.79</u>	<u>31.70</u>	29.24	<u>31.25</u>	23.66	<u>23.66</u>	<u>16.96</u>	<u>8.9</u>
	(d) Accur	racy (%) o	on GSM8	K CoT be	nchmark	(LLAMA	3.2 3B-Ir	struct)		
Temperature	0.0	0.3	0.5	0.7	1.0	1.5	2.0	3.0	4.0	5.0
Temp' Only	76.72	77.10	76.42	74.00	64.59	3.11	0.00	0.00	0.00	0.00
Top-p = 0.9	76.72	76.65	76.72	75.66	73.31	28.81	0.00	0.00	0.00	0.00
Top-p = 0.95	76.72	77.41	77.63	76.50	71.11	16.83	0.00	0.00	0.00	0.00
Min-p = 0.05	76.72	76.12	76.50	75.51	73.24	57.92	26.61	0.15	0.00	0.00
Min-p = 0.1	76.72	77.18	75.51	75.36	73.01	<u>75.44</u>	<u>52.08</u>	<u>2.50</u>	0.00	0.00
	(e) Accu	racy (%)	on GPQA	Main ber	nchmark (LLAMA	3.1 8B-In	struct)		
Temperature	0.0	0.3	0.5	0.7	1.0	1.5	2.0	3.0	4.0	5.0
Temp' Only	29.02	27.46	28.79	29.91	30.36	22.99	6.92	3.12	2.46	3.3
Тор-р = 0.9	29.02	28.57	29.69	<u>30.80</u>	28.79	24.55	9.38	2.46	2.46	2.9
Тор-р = 0.95	29.02	<u>30.36</u>	<u>31.92</u>	27.68	<u>30.58</u>	25.67	10.04	2.68	2.90	2.9
Min-p = 0.05	29.02	30.13	31.70	<u>30.80</u>	28.35	26.34	<u>29.24</u>	22.77	9.82	5.1
Min-p = 0.1	29.02	30.13	29.69	29.24	29.24	<u>32.14</u>	25.22	<u>22.99</u>	<u>20.54</u>	<u>12</u>
	(f) Accur	acy (%) c	on GSM81	K CoT be	nchmark (LLAMA	3.1 8B-In	struct)		
Temperature	0.0	0.3	0.5	0.7	1.0	1.5	2.0	3.0	4.0	5.0
Temp' Only	84.91	84.61	<u>84.84</u>	81.50	75.21	10.39	0.00	0.00	0.00	0.00
Top-p = 0.9	84.91	84.91	84.08	83.24	80.36	48.37	0.08	0.00	0.00	0.00
$T_{0n} = 0.05$	84.91	84.23	84.08	82.26	80.06	32.52	0.00	0.00	0.00	0.00
10p-p = 0.95										
Min-p = 0.93	84.91	<u>85.06</u>	84.31	<u>83.32</u>	80.44	67.25	32.15	0.08	0.00	0.00

1133

$(a) Accuracy (\%) on GPQA Main benchmark (Llama 3.1 70) Temperature 0.5 0.7 1.0 2.0 3.0 Temp' Only 40.85 39.51 41.07 5.58 2.44 Top-p = 0.9 40.85 42.19 40.63 7.81 3.35 Top-p = 0.95 40.63 41.29 39.51 6.47 2.66 Min-p = 0.05 40.85 41.52 38.62 33.71 23.8 Min-p = 0.1 41.07 42.19 41.52 33.04 24.5 Min-p = 0.2 43.30 41.96 40.40 34.38 31.4 (b) Accuracy (%) on GSM8K CoT benchmark (Llama 3.1 70) Temp' Only 93.33 0.08 Top-p = 0.9 93.48 0.08 Min-p = 0.2 92.42 61.03 D.2 RESULTS OF GPQA MAIN OR GSM8K COT BENCHMARKS FOR M Low TEMPERATURES (T \leq 0.5)Table 9: Accuracy (%) on GPQA Main benchmark (Mistral-7B).Temperature 0.0 0.05 0.1 0.2 0.3Temp Only 27.68 27.68 26.34 25.22 24.33Top-p = 0.9 27.68 27.90 27.23 25.22 24.55Min-p = 0.1 27.68 27.90 27.23 25.45 24.55Min-p = 0.1 27.68 28.35 27.46 26.56 24.33(Low Temperature 0.1 27.68 28.35 27.46 26$	(a) Accuracy (%) on GPQA Main benchmark (Llama 3.1 70E Temperature 0.5 0.7 1.0 2.0 3.0 Temp' Only 40.85 39.51 41.07 5.58 2.46 Top-p = 0.9 40.85 42.19 40.63 7.81 3.35 Top-p = 0.95 40.63 41.29 39.51 6.47 2.68 Min-p = 0.05 40.85 41.52 38.62 33.71 23.8 Min-p = 0.1 41.07 42.19 41.52 33.04 24.5 Min-p = 0.2 43.30 41.96 40.40 34.38 31.4 (b) Accuracy (%) on GSM8K CoT benchmark (Llama 3.1 70F Temp' Only 93.33 0.08 $70rp-p = 0.9$ 93.48 0.08 Min-p = 0.2 92.42 61.03 $70rp-p = 0.9$ 93.48 0.08 Min-p = 0.05 93.03 6.07 $70rp-p = 0.9$ 7.48 $70rp-p = 0.9$ 7.68 27.68 26.34 25.22 24.33 $70rp-p = 0.9$ 7.68 27.68 26.34 $25.$			(a) Accurac	y (%) on	GPQA M	lain ben	chmark	(Llama	a 3.1 70	0B)
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{tabular}{ c c c c c c c c c c c c c c c c c c c$										
$\frac{1}{\text{Temp' Only}} \begin{array}{c} 40.85 & 39.51 & 41.07 & 5.58 & 2.46 \\ \text{Top-p} = 0.9 & 40.85 & 42.19 & 40.63 & 7.81 & 3.35 \\ \text{Top-p} = 0.95 & 40.63 & 41.29 & 39.51 & 6.47 & 2.68 \\ \text{Min-p} = 0.05 & 40.85 & 41.52 & 38.62 & 33.71 & 23.8 \\ \text{Min-p} = 0.1 & 41.07 & 42.19 & 41.52 & 33.04 & 24.5 \\ \text{Min-p} = 0.2 & 43.30 & 41.96 & 40.40 & 34.38 & 31.4 \\ \hline \begin{array}{c} (b) \text{ Accuracy (\%) on GSM8K CoT benchmark (Llama 3.1 70)} \\ \hline \hline \text{Temp' Only} & 93.33 & 0.08 \\ \hline \text{Top-p} = 0.9 & 93.48 & 0.08 \\ \text{Min-p} = 0.2 & 93.03 & 6.07 \\ \hline \text{Min-p} = 0.2 & 92.42 & 61.03 \\ \hline \end{array}$	Temp' Only 40.85 39.51 41.07 5.58 2.46 Top-p = 0.9 40.85 42.19 40.63 7.81 3.35 Top-p = 0.95 40.63 41.29 39.51 6.47 2.68 Min-p = 0.05 40.85 41.52 38.62 33.71 23.8 Min-p = 0.1 41.07 42.19 41.52 33.04 24.5 Min-p = 0.2 43.30 41.96 40.40 34.38 31.4 (b) Accuracy (%) on GSM8K CoT benchmark (Llama 3.1 70f Temp' Only 93.33 0.08 Min-p = 0.2 93.48 0.08 Min-p = 0.05 93.03 6.07 Min-p = 0.2 92.42 61.03 61.03 6.7 Min-p = 0.2 92.42 61.03 61.03 6.7 Min-p = 0.2 92.42 61.03 61.03 61.03 Temp' Only 93.33 0.08 Min-p = 0.05 93.08 6.07 61.03 Temp' Only 93.33 0.08 Temp' Only 93.33 0.08 Temp' Only 0.05 0.1 <t< td=""><td></td><td></td><td>Temperatu</td><td>re 0.:</td><td>5 0.'</td><td>7 1</td><td>1.0</td><td>2.0</td><td>3.0</td><td>)</td></t<>			Temperatu	re 0.:	5 0.'	7 1	1.0	2.0	3.0)
$\begin{array}{r} \mbox{Temp Only} & 40.83 & 59.31 & 41.07 & 5.38 & 2.40 \\ \mbox{Top-p} = 0.9 & 40.85 & 42.19 & 40.63 & 7.81 & 3.35 \\ \mbox{Top-p} = 0.95 & 40.63 & 41.29 & 39.51 & 6.47 & 2.68 \\ \mbox{Min-p} = 0.05 & 40.85 & 41.52 & 38.62 & 33.71 & 23.8 \\ \mbox{Min-p} = 0.1 & 41.07 & 42.19 & 41.52 & 33.04 & 24.5 \\ \mbox{Min-p} = 0.2 & 43.30 & 41.96 & 40.40 & 34.38 & 31.2 \\ \mbox{(b) Accuracy} (\%) on GSM8K CoT benchmark (Llama 3.1 70 \\ \hline \mbox{Temperature} & 0.7 & 3.0 \\ \hline \mbox{Temp} Only & 93.33 & 0.08 \\ \mbox{Top-p} = 0.9 & 93.48 & 0.08 \\ \mbox{Min-p} = 0.2 & 92.42 & 61.03 \\ \mbox{Min-p} = 0.2 & 92.42 & 61.03 \\ \mbox{Min-p} = 0.2 & 92.42 & 61.03 \\ \mbox{Temperature} & 0.7 & 3.0 \\ \hline \mbox{Temperature} & 0.0 & 0.05 & 0.1 & 0.2 & 0.3 \\ \hline \mbox{Temperature} & 0.0 & 0.05 & 0.1 & 0.2 & 0.3 \\ \hline \mbox{Temperature} & 0.0 & 0.05 & 0.1 & 0.2 & 0.3 \\ \hline \mbox{Temp Only} & 27.68 & 27.68 & 26.34 & 25.22 & 24.33 \\ \mbox{Top-p} = 0.9 & 27.68 & 27.90 & 27.23 & 25.45 & 24.55 \\ \mbox{Min-p} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \mbox{Min-p} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \mbox{Temp} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \mbox{Temp} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \mbox{Temp} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \mbox{Temp} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \mbox{Temp} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \mbox{Temp} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \mbox{Temp} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \mbox{Temp} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \mbox{Temp} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \mbox{Temp} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \mbox{Temp} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \mbox{Temp} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \mbox{Temp} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \mbox{Temp} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \mbox{Temp} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \mbox{Temp} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \m$	$\frac{1}{100} - \frac{1}{100} + \frac{1}$			Tama' On	1 40	05 20	51 /	11.07	5 50	2.4	16
$\frac{10pp = 0.9}{10p = 0.95} = 40.63 = \frac{12.12}{40.63} = 40.63 = 1.61 = 3.5.5 = 1.01 = 1.02 = $	$\frac{100}{100} = 0.95 + 40.63 + \frac{11.29}{41.29} + 40.63 + \frac{11.29}{39.51} = 6.47 + 2.68 \\ Min-p = 0.05 + 40.85 + 41.52 + 38.62 + 33.71 + 23.8 \\ Min-p = 0.1 + 41.07 + \frac{42.19}{41.52} + \frac{41.52}{33.04} + 24.5 \\ Min-p = 0.2 + \frac{43.30}{3.30} + 41.96 + 40.40 + \frac{34.38}{34.38} + \frac{31.4}{34.38} \\ \hline (b) Accuracy (\%) on GSM8K CoT benchmark (Llama 3.1 701 + 1.52 + 1.5$			Temp On Tem $p = 0$	1y 40 0 40	0.85 39	10 /	+1.07 10.63	J.JO 7 81	2.4	+0
$\frac{1}{10} \text{ p} = 0.33 + 40.35 + 41.52 + 38.62 + 33.71 + 23.8}{\text{Min-p} = 0.1 + 41.07 + 42.19 + 41.52 + 33.04 + 24.5}{\text{Min-p} = 0.2 + 43.30 + 41.96 + 40.40 + 34.38 + 31.4} \text{ (b) Accuracy (%) on GSM8K CoT benchmark (Llama 3.1 70)} \\ \hline \hline \text{Temperature} = 0.7 + 3.0 \\ \hline \hline \text{Temp' Only} + 93.33 + 0.08 \\ \text{Min-p} = 0.9 + 93.48 + 0.08 \\ \text{Min-p} = 0.9 + 93.48 + 0.08 \\ \text{Min-p} = 0.05 + 93.03 + 6.07 \\ \text{Min-p} = 0.2 + 92.42 + 61.03 \\ \hline \text{Min-p} = 0.2 + 92.42 + 61.03 \\ \hline \text{Min-p} = 0.2 + 92.42 + 61.03 \\ \hline \text{Min-p} = 0.2 + 92.42 + 61.03 \\ \hline \text{Min-p} = 0.2 + 92.42 + 61.03 \\ \hline \text{Min-p} = 0.2 + 92.42 + 61.03 \\ \hline \text{Min-p} = 0.2 + 92.42 + 61.03 \\ \hline \text{Min-p} = 0.2 + 92.42 + 61.03 \\ \hline \text{Min-p} = 0.2 + 92.42 + 61.03 \\ \hline \text{Min-p} = 0.2 + 92.42 + 61.03 \\ \hline \text{Min-p} = 0.2 + 92.42 + 61.03 \\ \hline \text{Min-p} = 0.2 + 92.42 + 61.03 \\ \hline \text{Min-p} = 0.1 + 92.42 + 92.42 \\ \hline \text{Min-p} = 0.1 + 92.42 + 92.42 \\ \hline \text$	$\begin{array}{c} \text{Min-p} = 0.05 & 40.85 & 41.52 & 38.62 & 33.71 & 23.8 \\ \text{Min-p} = 0.1 & 41.07 & \underline{42.19} & \underline{41.52} & 33.04 & 24.5 \\ \text{Min-p} = 0.2 & \underline{43.30} & 41.96 & 40.40 & \underline{34.38} & \underline{31.4} \\ \hline \text{(b) Accuracy (\%) on GSM8K CoT benchmark (Llama 3.1 70)} \\ \hline \hline & \hline &$			Top-p = 0.	9 40 95 40	$1.63 \underline{-42}$	<u>-19</u> -	R0.05	6.47	2.6	55
$\begin{array}{r} \text{Min-p} = 0.1 & 41.07 & \underline{42.19} & \underline{41.52} & 33.04 & 24.3\\ \text{Min-p} = 0.2 & \underline{43.30} & 41.96 & 40.40 & \underline{34.38} & \underline{31.4}\\ \hline \text{(b) Accuracy (\%) on GSM8K CoT benchmark (Llama 3.1 70)} \\ \hline \hline \text{Temperature} & 0.7 & 3.0 \\ \hline \hline \text{Tempor Only} & 93.33 & 0.08 \\ \text{Top-p} = 0.9 & \underline{93.48} & 0.08 \\ \text{Min-p} = 0.05 & 93.03 & 6.07 \\ \text{Min-p} = 0.2 & 92.42 & \underline{61.03} \\ \hline \text{Min-p} = 0.2 & 92.42 & \underline{61.03} \\ \hline \text{Composition of GPQA MAIN OR GSM8K CoT BENCHMARKS FOR M Low TEMPERATURES (T \leq 0.5)} \\ \hline \text{Table 9: Accuracy (\%) on GPQA Main benchmark for Mistral-(a) Accuracy (\%) on GPQA Main benchmark (Mistral-7B))} \\ \hline \hline \hline \\ \hline \hline \\ \hline \\ \hline \hline \\ \hline \\ \hline \\ \hline \\ $	$\frac{\text{Min-p} = 0.1 41.07 42.19}{\text{Min-p} = 0.2 43.30 41.96 40.40 34.38 31.4}$ (b) Accuracy (%) on GSM8K CoT benchmark (Llama 3.1 70) $\frac{\text{Temperature} 0.7 3.0}{\text{Tempr'Only} 93.33 0.08}$ $\text{Top-p} = 0.9 93.48 0.08$ $\text{Min-p} = 0.2 92.42 61.03$ RESULTS OF GPQA MAIN OR GSM8K COT BENCHMARKS FOR M Low TEMPERATURES (T ≤ 0.5) Table 9: Accuracy (%) on GPQA Main benchmark for Mistral-7 (a) Accuracy (%) on GPQA Main benchmark (Mistral-7B). $\frac{\text{Temperature} 0.0 0.05 0.1 0.2 0.3 $			Min-n = 0	05 + 0	85 41	52 3	38.62	33 71	23	88
$\frac{\text{Min-p} = 0.2}{\text{Min-p} = 0.2} \frac{43.30}{41.96} \frac{41.96}{40.40} \frac{34.38}{34.38} \frac{31.4}{31.40}$ (b) Accuracy (%) on GSM8K CoT benchmark (Llama 3.1 70) $\frac{\text{Temperature}}{\text{Temperature}} \frac{0.7}{3.0} \frac{3.0}{\text{Temp' Only}} \frac{93.33}{93.33} \frac{0.08}{0.08}$ Top-p = 0.9 $ \frac{93.48}{93.03} \frac{0.08}{6.07}$ Min-p = 0.05 $ \frac{93.03}{93.03} \frac{6.07}{6.07}$ Min-p = 0.2 $ 92.42 \frac{61.03}{61.03}$ D.2 RESULTS OF GPQA MAIN OR GSM8K COT BENCHMARKS FOR M LOW TEMPERATURES (T ≤ 0.5) Table 9: Accuracy (%) on GPQA Main benchmark for Mistral- (a) Accuracy (%) on GPQA Main benchmark (Mistral-7B). $\frac{\text{Temperature}}{10.00} \frac{0.05}{0.1} \frac{0.2}{0.2} \frac{0.3}{0.3}$ Temp Only $ 27.68 27.68 26.34 25.22 24.33$ Top-p = 0.9 $ 27.68 27.90 27.23 25.45 24.55$ Min-p = 0.1 $ 27.68 28.35 27.46 26.56 24.33$	Min-p = 0.2 43.30 41.96 40.40 34.38 31.4 (b) Accuracy (%) on GSM8K CoT benchmark (Llama 3.1 701 Temperature 0.7 3.0 Temp' Only 93.33 0.08 Top-p = 0.9 93.48 0.08 Min-p = 0.05 93.03 6.07 Min-p = 0.2 92.42 61.03 RESULTS OF GPQA MAIN OR GSM8K CoT BENCHMARKS FOR M Low TEMPERATURES (T ≤ 0.5) Table 9: Accuracy (%) on GPQA Main benchmark for Mistral-7 (a) Accuracy (%) on GPQA Main benchmark (Mistral-7B). Temp Only 27.68 26.34 25.22 24.33 Top-p = 0.9 27.68 27.68 26.56 24.78 Top-p = 0.9 27.68 27.90 27.23 25.22 24.55 Min-p = 0.05 27.68 27.90 27.23 25.45 24.55 Min-p = 0.1 27.68 27.90 27.23 25.45 24.55 Min-p = 0.1 27.68 27.90 27.23 25.45 24.55 Min-p = 0.1 27.68 28.35 27.46 26			Min p = 0 $Min - p = 0$.1 41	.07 42	.19 4	41.52	33.04	1 24	.55
$\frac{1}{1}$ (b) Accuracy (%) on GSM8K CoT benchmark (Llama 3.1 70) $\frac{1}{\text{Temperature}} = 0.7 3.0$ $\frac{1}{\text{Temp' Only}} = 9.3.48 0.08$ $\text{Top-p} = 0.9 \underline{93.48} 0.08$ $\text{Min-p} = 0.2 92.42 \underline{61.03}$ $D.2 \text{Results of GPQA MAIN OR GSM8K CoT BENCHMARKS FOR M Low TEMPERATURES (T \le 0.5)$ $Table 9: \text{ Accuracy (%) on GPQA Main benchmark for Mistral-(a) Accuracy (%) on GPQA Main benchmark (Mistral-7B).}$ $\frac{1}{\text{Temperature}} 0.0 0.05 0.1 0.2 0.3$ $\frac{1}{\text{Temp Only}} 27.68 27.68 26.34 25.22 24.33$ $Top-p = 0.9 27.68 27.90 27.23 25.45 24.55$ $\text{Min-p} = 0.05 27.68 27.90 27.23 25.45 24.55$ $\text{Min-p} = 0.1 27.68 28.35 27.46 26.56 24.33$	Image:			Min-p = 0	.2 43	$30 \frac{1}{41}$.96 4	40.40	34.38	3 31	.47
(b) Accuracy (%) on GSM8K CoT benchmark (Llama 3.1 70	(b) Accuracy (%) on GSM8K CoT benchmark (Llama 3.1 70) $ \frac{1}{\text{Temperature}} 0.7 3.0 \\ \hline \text{Temp' Only} 93.33 0.08 \\ \text{Top-p} = 0.9 93.48 0.08 \\ \text{Min-p} = 0.05 93.03 6.07 \\ \text{Min-p} = 0.2 92.42 61.03 \\ \hline \text{Min-p} = 0.0 60.05 0.1 0.2 0.3 \\ \hline \text{Min-p-p} = 0.9 27.68 27.68 26.34 25.22 24.33 \\ \hline \text{Top-p} = 0.9 27.68 27.90 27.23 25.22 24.55 \\ \hline \text{Min-p} = 0.05 27.68 27.90 27.23 25.45 24.55 \\ \hline \text{Min-p} = 0.1 27.68 28.35 27.46 26.56 24.33 \\ \hline \text{Table 10: Accuracy (%) on GSM8K benchmark (Mistral-7E).} \\ \hline \hline \text{Temperature} 0.0 0.05 0.1 0.2 0.3 \\ \hline \text{TempOnly} 39.35 38.59 38.21 38.59 37.23 \\ \hline \text{Min-p} = 0.1 27.23 25.22 24.55 \\ \hline \text{Min-p} = 0.1 27.23 25.45 24.55 \\ \hline \text{Min-p} = 0.1 27.68 28.35 27.46 26.56 24.33 \\ \hline \text{Temperature} 0.0 0.05 0.1 0.2 0.3 \\ \hline \text{Temperature} 0.0 0.05 0.1 0.2 0.3 \\ \hline \text{Temperature} 0.0 0.05 0.1 0.2 0.3 \\ \hline \text{TempOnly} 39.35 38.59 38.21 38.59 37.23 \\ \hline \text{Min-p} = 0.1 27.23 \\ \hline \text{Min-p} = 0.1 27.68 \\ \hline \text{Min-p} = 0.1 \\ \hline \text{Min-p} = 0.1 \\ \hline Mi$						~ ~ .				
$\frac{\text{Temperature}}{\text{Temp}' \text{ Only} 93.33 0.08}{\text{Top-p} = 0.9 93.48} 0.08}{\text{Min-p} = 0.05 93.03 6.07}{\text{Min-p} = 0.2 92.42 61.03}}$ D.2 RESULTS OF GPQA MAIN OR GSM8K COT BENCHMARKS FOR M Low TEMPERATURES (T ≤ 0.5) Table 9: Accuracy (%) on GPQA Main benchmark for Mistral- (a) Accuracy (%) on GPQA Main benchmark (Mistral-7B). $\frac{\text{Temperature} 0.0 0.05 0.1 0.2 0.3}{\text{Temp Only} 27.68 27.68 26.34 25.22 24.33}{\text{Top-p} = 0.9 27.68 28.13 27.23 26.56 24.78}{\text{Top-p} = 0.95 27.68 27.90 27.23 25.45 24.55}{\text{Min-p} = 0.1 27.68 28.35 27.46 26.56 24.33}$	Temperature 0.7 3.0 Temp' Only 93.33 0.08 Top-p = 0.9 93.48 0.08 Min-p = 0.05 93.03 6.07 Min-p = 0.2 92.42 61.03 RESULTS OF GPQA MAIN OR GSM8K COT BENCHMARKS FOR M Low TEMPERATURES (T ≤ 0.5) Table 9: Accuracy (%) on GPQA Main benchmark for Mistral-7 (a) Accuracy (%) on GPQA Main benchmark (Mistral-7B). Temperature 0.0 0.05 0.1 0.2 0.3 Temp Only 27.68 27.68 26.34 25.22 24.33 Top-p = 0.9 27.68 27.90 27.23 26.56 24.78 Top-p = 0.9 27.68 27.90 27.23 25.45 24.55 24.55 Min-p = 0.1 27.68 28.35 27.46 26.56 24.33 Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E (a) Accuracy (%) on GSM8K benchmark (Mistral-7B). $$ $$ $$ $$ $$ $$ $$ $$ $$ $$			(b) Accurac	y (%) on	GSM8K (Coll ber	ichmark	(Llam	a 3.17	0B)
$\begin{array}{cccccccccccccccccccccccccccccccccccc$	Temp' Only 93.33 0.08 Top-p = 0.9 $Min-p = 0.05$ 93.03 6.07 Min-p = 0.2 Min-p = 0.2 92.42 61.03 RESULTS OF GPQA MAIN OR GSM8K COT BENCHMARKS FOR M LOW TEMPERATURES (T ≤ 0.5)Table 9: Accuracy (%) on GPQA Main benchmark for Mistral-7 (a) Accuracy (%) on GPQA Main benchmark (Mistral-7B).Temperature 0.0 0.05 0.1 0.2 0.3 Temp Only 27.68 27.68 26.34 25.22 24.33 Top-p = 0.9 27.68 27.90 27.23 26.56 24.78 Top-p = 0.9 27.68 27.90 27.23 25.45 24.55 Min-p = 0.05 27.68 27.90 27.23 25.45 24.55 Min-p = 0.1 27.68 28.35 27.46 26.56 24.33 Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E (a) Accuracy (%) on GSM8K benchmark (Mistral-7B).Temponly 39.35 38.59 38.21 38.59 37.23				Temp	erature	0.7	3.0			
$Top-p = 0.9 \underline{93.48} 0.08$ $Min-p = 0.05 \underline{93.03} 6.07$ $Min-p = 0.2 92.42 \underline{61.03}$ $D.2 RESULTS \text{ OF GPQA MAIN OR GSM8K COT BENCHMARKS FOR M LOW TEMPERATURES (T \le 0.5)$ $Table 9: \text{ Accuracy (\%) on GPQA Main benchmark for Mistral-(a) Accuracy (\%) on GPQA Main benchmark (Mistral-7B)}$ $\overline{Temp \text{ Only } 27.68 27.68 26.34 25.22 24.33}$ $Top-p = 0.9 27.68 27.68 26.34 25.22 24.33$ $Top-p = 0.9 27.68 27.90 27.23 25.45 24.55$ $Min-p = 0.1 27.68 28.35 27.46 26.56 24.33$	$Top-p = 0.9 \underline{93.48} 0.08$ $Min-p = 0.05 \underline{93.03} 6.07$ $Min-p = 0.2 92.42 \underline{61.03}$ RESULTS OF GPQA MAIN OR GSM8K COT BENCHMARKS FOR M Low TEMPERATURES (T ≤ 0.5) Table 9: Accuracy (%) on GPQA Main benchmark for Mistral-7 (a) Accuracy (%) on GPQA Main benchmark (Mistral-7B). Temponly 27.68 27.68 26.34 25.22 24.33 27.23 26.56 24.78 27.68 27.90 27.23 25.22 24.55 25.55 24.55 25.55 25.55 25.55 25.55 25.55 25.55 25.55 25.55 25.55 25.55 2				Temp	' Only	93.33	3 0.0	8		
$\begin{array}{c} \mbox{Min-p} = 0.05 & 93.03 & 6.07 \\ \mbox{Min-p} = 0.2 & 92.42 & \underline{61.03} \end{array}$	$\begin{array}{c} \mbox{Min-p} = 0.05 & 93.03 & 6.07 \\ \mbox{Min-p} = 0.2 & 92.42 & 61.03 \end{array} \end{array}$ RESULTS OF GPQA MAIN OR GSM8K COT BENCHMARKS FOR M Low TEMPERATURES (T ≤ 0.5) Table 9: Accuracy (%) on GPQA Main benchmark for Mistral-7 (a) Accuracy (%) on GPQA Main benchmark (Mistral-7B). $\hline \hline \mbox{Temperature} & 0.0 & 0.05 & 0.1 & 0.2 & 0.3 \\ \hline \mbox{Temp Only} & 27.68 & 27.68 & 26.34 & 25.22 & 24.33 \\ \hline \mbox{Top-p} = 0.9 & 27.68 & 28.13 & 27.23 & 26.56 & 24.78 \\ \hline \mbox{Top-p} = 0.95 & 27.68 & 27.90 & 27.23 & 25.45 & 24.55 \\ \hline \mbox{Min-p} = 0.1 & 27.68 & 28.35 & 27.46 & 26.56 & 24.33 \\ \hline \mbox{Table 10: Accuracy (%) on GSM8K benchmark (Mistral-7E).} \\ \hline \hline \mbox{Temperature} & 0.0 & 0.05 & 0.1 & 0.2 & 0.3 \\ \hline \mbox{Table 10: Accuracy (%) on GSM8K benchmark (Mistral-7E).} \\ \hline \hline \mbox{Temperature} & 0.0 & 0.05 & 0.1 & 0.2 & 0.3 \\ \hline \mbox{Temperature} & 0.0 & 0.05 & 0.1 & 0.2 & 0.3 \\ \hline \mbox{Table 10: Accuracy (%) on GSM8K benchmark (Mistral-7E).} \\ \hline \mbox{Temperature} & 0.0 & 0.05 & 0.1 & 0.2 & 0.3 \\ \hline \mbox{Temperature} & 0.0 & 0.05 & 0.1 & 0.2 & 0.3 \\ \hline \mbox{Temperature} & 0.0 & 0.05 & 0.1 & 0.2 & 0.3 \\ \hline \mbox{Temperature} & 0.0 & 0.05 & 0.1 & 0.2 & 0.3 \\ \hline \mbox{Temperature} & 0.0 & 0.05 & 0.1 & 0.2 & 0.3 \\ \hline \mbox{Temperature} & 0.0 & 0.05 & 0.1 & 0.2 & 0.3 \\ \hline \mbox{Temp Only} & 39.35 & 38.59 & 38.21 & 38.59 & 37.23 \\ \hline \end{tabular}$				Top-p	0 = 0.9	<u>93.48</u>	<u>8</u> 0.0	8		
$ \underbrace{ \text{Min-p} = 0.2 92.42 \underline{61.03} } \\ \text{O.2 RESULTS OF GPQA MAIN OR GSM8K COT BENCHMARKS FOR M LOW TEMPERATURES (T \leq 0.5) } \\ \text{Table 9: Accuracy (%) on GPQA Main benchmark for Mistral-(a) Accuracy (%) on GPQA Main benchmark (Mistral-7B). } \\ \hline \underline{\text{Temperature}} 0.0 0.05 0.1 0.2 0.3 \\ \hline \underline{\text{Temp Only}} 27.68 27.68 26.34 25.22 24.33 \\ \hline \text{Top-p} = 0.9 27.68 27.90 27.23 \underline{25.22} 24.55 \\ \hline \text{Min-p} = 0.05 27.68 27.90 27.23 25.45 24.55 \\ \hline \text{Min-p} = 0.1 27.68 \underline{28.35} \underline{27.46} \underline{26.56} 24.33 \\ \hline \end{array} $	Min-p = 0.2 92.4261.03RESULTS OF GPQA MAIN OR GSM8K COT BENCHMARKS FOR M Low TEMPERATURES (T ≤ 0.5)Table 9: Accuracy (%) on GPQA Main benchmark for Mistral-7(a) Accuracy (%) on GPQA Main benchmark (Mistral-7B).Temperature 0.00.050.10.20.3Temponly 27.6827.6826.3425.2224.33Top-p = 0.927.6827.9027.2326.5624.78Top-p = 0.9527.6827.9027.2325.4524.55Min-p = 0.0527.6827.9027.2325.4524.55Min-p = 0.127.6828.3527.4626.5624.33Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E(a) Accuracy (%) on GSM8K benchmark (Mistral-7B).Temperature0.00.050.10.20.3Temperature0.00.050.10.20.3Temperature0.00.050.10.20.3Temperature0.00.050.10.20.3Temperature0.00.050.10.20.3Temperature0.00.050.10.20.3Table 10: Accuracy (%) on GSM8K benchmark (Mistral-7B).Temperature0.00.050.10.20.3Temperature0.00.050.10.20				Min-p	p = 0.05	93.03	6.0	7		
D.2 RESULTS OF GPQA MAIN OR GSM8K COT BENCHMARKS FOR M Low TEMPERATURES (T ≤ 0.5) Table 9: Accuracy (%) on GPQA Main benchmark for Mistral- (a) Accuracy (%) on GPQA Main benchmark (Mistral-7B). $\frac{1}{10000000000000000000000000000000000$	RESULTS OF GPQA MAIN OR GSM8K CoT BENCHMARKS FOR M Low TEMPERATURES (T ≤ 0.5) Table 9: Accuracy (%) on GPQA Main benchmark for Mistral-7 (a) Accuracy (%) on GPQA Main benchmark (Mistral-7B). Temperature 0.0 0.05 0.1 0.2 0.3 Temp Only 27.68 27.68 26.34 25.22 24.33 Top-p = 0.9 27.68 27.90 27.23 26.56 24.78 Top-p = 0.95 27.68 27.90 27.23 25.45 24.55 Min-p = 0.1 27.68 28.35 27.46 26.56 24.33 Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E (a) Accuracy (%) on GSM8K benchmark (Mistral-7B). Temperature 0.0 0.05 0.1 0.2 0.3 Temperature 0.0 0.05 0.1 0.2 0.3				Min-p	p = 0.2	92.42	2 <u>61.</u>	<u>03</u>		
D.2 RESULTS OF GPQA MAIN OR GSM8K COT BENCHMARKS FOR M Low TEMPERATURES (T ≤ 0.5) Table 9: Accuracy (%) on GPQA Main benchmark for Mistral- (a) Accuracy (%) on GPQA Main benchmark (Mistral-7B). Temperature 0.0 0.05 0.1 0.2 0.3 Temp Only 27.68 27.68 26.34 25.22 24.33 Top-p = 0.9 27.68 28.13 27.23 26.56 24.78 Top-p = 0.95 27.68 27.90 27.23 25.22 24.55 Min-p = 0.05 27.68 27.90 27.23 25.45 24.55 Min-p = 0.1 27.68 28.35 27.46 26.56 24.33	RESULTS OF GPQA MAIN OR GSM8K COT BENCHMARKS FOR M LOW TEMPERATURES (T \leq 0.5) Table 9: Accuracy (%) on GPQA Main benchmark for Mistral-7 (a) Accuracy (%) on GPQA Main benchmark (Mistral-7B). Temperature 0.0 0.05 0.1 0.2 0.3 Temp Only 27.68 27.68 26.34 25.22 24.33 Top-p = 0.9 27.68 27.90 27.23 25.22 24.55 Min-p = 0.05 27.68 27.90 27.23 25.45 24.55 Min-p = 0.1 27.68 28.35 27.46 26.56 24.33 Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E (a) Accuracy (%) on GSM8K benchmark (Mistral-7B). Temperature 0.0 0.05 0.1 0.2 0.3 Temperature 0.0 0.05 0.1 0.2 0.3										
7.2 RESULTS OF GPQA MAIN OR GSM8K COT BENCHMARKS FOR M LOW TEMPERATURES (T ≤ 0.5) Table 9: Accuracy (%) on GPQA Main benchmark for Mistral- (a) Accuracy (%) on GPQA Main benchmark (Mistral-7B). Temperature 0.0 0.05 0.1 0.2 0.3 Temp Only 27.68 27.68 26.34 25.22 24.33 Top-p = 0.9 27.68 28.13 27.23 26.56 24.78 Top-p = 0.95 27.68 27.90 27.23 25.22 24.55 Min-p = 0.05 27.68 27.90 27.23 25.45 24.55 Min-p = 0.1 27.68 28.35 27.46 26.56 24.33	RESULTS OF GPQA MAIN OR GSM8K COT BENCHMARKS FOR M Low TEMPERATURES (T \leq 0.5) Table 9: Accuracy (%) on GPQA Main benchmark for Mistral-7 (a) Accuracy (%) on GPQA Main benchmark (Mistral-7B). Temperature 0.0 0.05 0.1 0.2 0.3 Temp Only 27.68 27.68 26.34 25.22 24.33 Top-p = 0.9 27.68 27.90 27.23 25.22 24.35 Top-p = 0.95 27.68 27.90 27.23 25.45 24.55 Min-p = 0.05 27.68 27.90 27.23 25.45 24.55 Min-p = 0.1 27.68 28.35 27.46 26.56 24.33 Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E (a) Accuracy (%) on GSM8K benchmark (Mistral-7B). Temperature 0.0 0.05 0.1 0.2 0.3	D 2	D -1		_	ON CONT	0. .				
LOW TEMPERATURES ($1 \le 0.5$)Table 9: Accuracy (%) on GPQA Main benchmark for Mistral-(a) Accuracy (%) on GPQA Main benchmark (Mistral-7B).Temperature 0.0 0.05 0.1 0.2 0.3Temperature 0.0 0.05 0.1 0.2 0.3Temp Only 27.68 27.68 26.34 25.22 24.33Top-p = 0.9 27.68 27.68 26.34 25.22 24.33Top-p = 0.9 27.68 27.90 27.23 25.22 24.55Min-p = 0.05 27.68 27.90 27.23 25.45 24.55Min-p = 0.1 27.68 28.35 27.46 26.56 24.33	LOW TEMPERATURES (T \leq 0.3)Table 9: Accuracy (%) on GPQA Main benchmark for Mistral-7(a) Accuracy (%) on GPQA Main benchmark (Mistral-7B).Temperature 0.0 0.05 0.1 0.2 0.3Temponly 27.68 27.68 26.34 25.22 24.33Top-p = 0.9 27.68 27.68 26.34 25.22 24.33Top-p = 0.9 27.68 27.90 27.23 25.22 24.55Min-p = 0.05 27.68 27.90 27.23 25.45 24.55Min-p = 0.1 27.68 28.35 27.46 26.56 24.33Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E(a) Accuracy (%) on GSM8K benchmark (Mistral-7B).Temperature 0.0 0.05 0.1 0.2 0.3Temperature 0.0 0.05 0.1 0.2 0.3	D.2	RESULTS	OF GPQA MA	IN OR G	SM8K (COT B	ENCHM	IARKS	FOR	MIS
Table 9: Accuracy (%) on GPQA Main benchmark for Mistral- (a) Accuracy (%) on GPQA Main benchmark (Mistral-7B).Temperature 0.0 0.05 0.1 0.2 0.3Temp Only 27.68 27.68 26.34 25.22 24.33Top-p = 0.9 27.68 27.68 26.34 25.22 24.33Top-p = 0.9 27.68 27.90 27.23 26.56 24.78Top-p = 0.95 27.68 27.90 27.23 25.22 24.55Min-p = 0.05 27.68 27.90 27.23 25.45 24.55Min-p = 0.127.68 28.35 27.46 26.56 24.33	Table 9: Accuracy (%) on GPQA Main benchmark for Mistral-7(a) Accuracy (%) on GPQA Main benchmark (Mistral-7B).Temperature0.00.050.10.20.3Temp Only27.6827.6826.3425.2224.33Top-p = 0.927.6827.9027.2326.5624.78Top-p = 0.9527.6827.9027.2325.2224.55Min-p = 0.0527.6827.9027.2325.4524.55Min-p = 0.127.6828.3527.4626.5624.33Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E(a) Accuracy (%) on GSM8K benchmark (Mistral-7B).Temperature0.00.050.10.20.3Temponly39.3538.5938.2138.5937.23		LOW IEM	MPERATURES ($1 \leq 0.5$)					
Table 9: Accuracy (%) on GPQA Main benchmark for Mistral- (a) Accuracy (%) on GPQA Main benchmark (Mistral-7B).Temperature 0.0 0.05 0.1 0.2 0.3Temp Only 27.68 27.68 26.34 25.22 24.33Top-p = 0.9 27.68 28.13 27.23 26.56 24.78Top-p = 0.95 27.68 27.90 27.23 25.22 24.55Min-p = 0.05 27.68 27.90 27.23 25.45 24.55Min-p = 0.127.68 28.35 27.46 26.56 24.33	Table 9: Accuracy (%) on GPQA Main benchmark for Mistral-,(a) Accuracy (%) on GPQA Main benchmark (Mistral-7B).Temperature0.00.050.10.20.3Temp Only27.6827.6826.3425.2224.33Top-p = 0.927.6827.9027.2326.5624.78Top-p = 0.9527.6827.9027.2325.2224.55Min-p = 0.0527.6827.9027.2325.4524.55Min-p = 0.127.6828.3527.4626.5624.33Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E(a) Accuracy (%) on GSM8K benchmark (Mistral-7B).Temperature0.00.050.10.20.3Temponly39.3538.5938.2138.5937.23		T 1	1 0 4			• •			<i>r</i>	1 71
(a) Accuracy (%) on GPQA Main benchmark (Mistral-7B).Temperature0.00.050.10.20.3Temp Only27.6827.6826.3425.2224.33Top-p = 0.927.6828.1327.23 26.5624.78 Top-p = 0.9527.6827.9027.2325.2224.55Min-p = 0.0527.6827.9027.2325.4524.55Min-p = 0.127.68 28.3527.4626.56 24.33	(a) Accuracy (%) on GPQA Main benchmark (Mistral-7B).Temperature0.00.050.10.20.3Temp Only27.6827.6826.3425.2224.33Top-p = 0.927.6827.9027.2325.2224.55Min-p = 0.0527.6827.9027.2325.4524.55Min-p = 0.127.6828.3527.4626.5624.33Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E(a) Accuracy (%) on GSM8K benchmark (Mistral-7B).Temperature0.00.050.10.20.3Temperature0.00.050.10.20.3Temperature0.00.050.10.20.3Temperature0.00.050.10.20.3Temperature0.00.050.10.20.3Temperature0.00.050.10.20.3Temp Only39.3538.5938.2138.5937.23		Tat	ble 9: Accuracy	(%) on (JPQA M	ain bei	nchmar	k for N	VIstral	l-//
Temperature0.00.050.10.20.3Temp Only27.6827.6826.3425.2224.33Top-p = 0.927.6828.1327.23 26.5624.78 Top-p = 0.9527.6827.9027.2325.2224.55Min-p = 0.0527.6827.9027.2325.4524.55Min-p = 0.127.68 28.3527.4626.56 24.33	Temperature0.00.050.10.20.3Temp Only27.6827.6826.3425.2224.33Top-p = 0.927.6828.1327.2326.5624.78Top-p = 0.9527.6827.9027.2325.2224.55Min-p = 0.0527.6827.9027.2325.4524.55Min-p = 0.127.6828.3527.4626.5624.33Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E(a) Accuracy (%) on GSM8K benchmark (Mistral-7B).Temperature0.00.050.10.20.3Temponly39.3538.5938.2138.5937.23				(07)	n GPOA	Main be	enchmar	k (Misi	tral-7B	D.
Temperature 0.0 0.05 0.1 0.2 0.3 Temp Only 27.68 27.68 26.34 25.22 24.33 Top-p = 0.9 27.68 28.13 27.23 26.56 24.78 Top-p = 0.95 27.68 27.90 27.23 25.22 24.55 Min-p = 0.05 27.68 27.90 27.23 25.45 24.55 Min-p = 0.1 27.68 28.35 27.46 26.56 24.33	Temperature 0.0 0.05 0.1 0.2 0.3 Temp Only27.6827.6826.3425.2224.33Top-p = 0.9 27.6828.1327.23 26.5624.78 Top-p = 0.95 27.6827.9027.2325.2224.55Min-p = 0.05 27.6827.9027.2325.4524.55Min-p = 0.1 27.68 28.3527.4626.56 24.33Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E(a) Accuracy (%) on GSM8K benchmark (Mistral-7B).Temperature 0.0 0.05 0.1 0.2 0.3 Temperature 0.0 0.05 0.1 0.2 0.3 Temponly 39.35 38.59 38.21 38.59 37.23			(a) Accura	100 00000						
Temp Only Top-p = 0.9 27.68 27.6826.34 28.1325.22 24.3324.33 24.78Top-p = 0.9 27.68 27.6828.13 27.9027.23 25.2224.78 24.55Min-p = 0.05 27.68 27.6827.90 27.9027.23 25.4525.45 24.55Min-p = 0.1 27.68 27.6828.35 27.4626.56 26.5624.33	Temp Only27.6827.6826.3425.2224.33Top-p = 0.927.6828.1327.23 26.5624.78 Top-p = 0.9527.6827.9027.2325.2224.55Min-p = 0.0527.6827.9027.2325.4524.55Min-p = 0.127.68 28.3527.4626.56 24.33Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E(a) Accuracy (%) on GSM8K benchmark (Mistral-7B).Temperature0.00.050.10.20.3Temponly39.3538.5938.2138.5937.23		-	(a) Accura	acy (%) 0						·
Top-p = 0.9 27.6828.1327.23 <u>26.56</u><u>24.78</u> Top-p = 0.95 27.6827.9027.2325.2224.55Min-p = 0.05 27.6827.9027.2325.4524.55Min-p = 0.1 27.68 <u>28.35</u><u>27.46</u><u>26.56</u> 24.33	$\begin{array}{cccccccccccccccccccccccccccccccccccc$		-	(a) Accura Temperature	0.0	0.05	0.1	0.2	0	0.3	0.5
Top-p = 0.95 27.6827.9027.2325.2224.55Min-p = 0.05 27.6827.9027.2325.4524.55Min-p = 0.1 27.6828.3527.4626.5624.33	Top-p = 0.95 27.6827.9027.23 $\overline{25.22}$ $\overline{24.55}$ Min-p = 0.05 27.6827.9027.23 25.45 24.55 Min-p = 0.1 27.68 28.35 27.46 26.56 24.33 Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E(a) Accuracy (%) on GSM8K benchmark (Mistral-7B).Temperature 0.0 0.05 0.1 0.2 0.3 Temperature 0.0 0.05 0.1 0.2 0.3 Temp Only 39.35 38.59 38.21 38.59 37.23			(a) Accura Temperature Temp Only	0.0 27.68	0.05	0.1	0.2	0 22 2).3	0.5
Min-p = 0.05 27.6827.9027.2325.4524.55Min-p = 0.1 27.68 28.3527.4626.56 24.33	Min-p = 0.05 27.6827.9027.2325.4524.55Min-p = 0.1 27.68 28.3527.4626.56 24.33Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E(a) Accuracy (%) on GSM8K benchmark (Mistral-7B).Temperature 0.0 0.05 0.1 0.2 0.3 Temperature 0.0 0.05 0.1 0.2 0.3 Temp Only 39.35 38.59 38.21 38.59 37.23		-	(a) Accura Temperature Temp Only Top-p = 0.9	0.0 27.68 27.68	0.05 27.68 28.13	0.1 26.34 27.23	0.2 4 25.2 3 26. 3	0 22 2 56 2).3 24.33 24.78	0.5 24 24
Min-p = 0.1 27.68 28.35 27.46 26.56 24.33	Min-p = 0.1 27.68 28.35 27.46 26.56 24.33 Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E (a) Accuracy (%) on GSM8K benchmark (Mistral-7B). Temperature 0.0 0.05 0.1 0.2 0.3 Temponly 39.35 38.59 38.21 38.59 37.23		-	(a) Accura Temperature Temp Only Top-p = 0.9 Top-p = 0.95	0.0 27.68 27.68 27.68	0.05 27.68 28.13 27.90	0.1 26.34 27.23 27.23	0.2 4 25.2 3 <u>26.3</u> 3 25.2	$\begin{array}{c} 0 \\ 22 & 2 \\ 56 & 2 \\ 22 & 2 \end{array}$).3 24.33 24.78 24.55	0.: 24 <u>24</u> 24
	Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E (a) Accuracy (%) on GSM8K benchmark (Mistral-7B). Temperature 0.0 0.05 0.1 0.2 0.3 Temp Only 39.35 38.59 38.21 38.59 37.23		-	(a) Accurate Temperature Temp Only Top- $p = 0.9$ Top- $p = 0.95$ Min- $p = 0.05$	0.0 27.68 27.68 27.68 27.68 27.68	0.05 27.68 28.13 27.90 27.90	0.1 26.34 27.23 27.23 27.23	0.2 4 25.2 3 <u>26.</u>3 3 25.2 3 25.2	0 22 2 56 2 22 2 45 2	0.3 24.33 24.78 24.55 24.55	0.: 24 24 24 24
	Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7E (a) Accuracy (%) on GSM8K benchmark (Mistral-7B). Temperature 0.0 0.05 0.1 0.2 0.3 Temp Only 39.35 38.59 38.21 38.59 37.23		-	(a) Accurate Temperature Top-p = 0.9 Top-p = 0.95 Min-p = 0.05 Min-p = 0.1	0.0 27.68 27.68 27.68 27.68 27.68 27.68	0.05 27.68 28.13 27.90 27.90 27.90 28.35	0.1 26.34 27.23 27.23 27.23 27.4	0.2 4 25.2 3 26. 3 25.2 3 25.2 5 26.	0 22 2 56 2 22 2 45 2 56 2	0.3 24.33 24.78 24.55 24.55 24.33	0.: 24 24 24 24 24 24
Table 10: Accuracy (%) on GSM8K benchmark for Mistral-7	(a) Accuracy (%) on GSM8K benchmark (Mistral-7B). Temperature 0.0 0.05 0.1 0.2 0.3 Temp Only 39.35 38.59 38.21 38.59 37.23		-	(a) Accurate Temperature Temp Only Top-p = 0.9 Top-p = 0.95 Min-p = 0.05 Min-p = 0.1	0.0 27.68 27.68 27.68 27.68 27.68 27.68	0.05 27.68 28.13 27.90 27.90 28.35	0.1 26.34 27.23 27.23 27.23 27.4	0.2 4 25.2 3 <u>26.2</u> 3 25.2 5 <u>26.2</u>	0 22 2 56 2 22 2 45 2 56 2	0.3 24.33 24.78 24.55 24.55 24.55 24.33	0. 24 24 24 24 24 24
	(a) Accuracy (%) on GSM8K benchmark (Mistral-7B). Temperature 0.0 0.05 0.1 0.2 0.3 Temp Only 39.35 38.59 38.21 38.59 37.23		- - - Ti	(a) Accurate Temperature Temp Only Top-p = 0.9 Top-p = 0.95 Min-p = 0.05 Min-p = 0.1	0.0 27.68 27.68 27.68 27.68 27.68 27.68 27.68	0.05 27.68 28.13 27.90 27.90 28.35 on GSM8	0.1 26.34 27.22 27.22 27.22 27.22 27.24	0.2 4 25.3 3 26. 3 3 25.3 3 25.4 6 26. 3 chmark	0 22 2 5 <u>6 2</u> 22 2 45 2 5 <u>6</u> 2	0.3 24.33 24.78 24.55 24.55 24.33 istral- ²	0. 24 24 24 24 24 24 24
(a) Accuracy (%) on GSM8K benchmark (Mistral-7B).	Temperature0.00.050.10.20.3Temp Only39.3538.5938.2138.5937.23		- - - Ti	(a) Accura Temperature Temp Only Top-p = 0.9 Top-p = 0.95 Min-p = 0.05 Min-p = 0.1 able 10: Accura	0.0 27.68 27.68 27.68 27.68 27.68 27.68 27.68	0.05 27.68 28.13 27.90 27.90 28.35 m GSM8	0.1 26.34 27.22 27.22 27.22 27.4 K benc	0.2 4 25.3 3 26.3 3 25.3 5 26.3 4 25.3 25.4 26.3 26.3 26.3 26.3 26.3 26.3 27.5 27.5 27.5 27.5 27.5 27.5 27.5 27.5 27.	0 22 2 56 2 22 2 45 2 56 2 for M	0.3 24.33 24.78 24.55 24.55 24.33 istral- ²	0. 2. 2. 2. 2. 2. 2. 7B
Temperature 0.0 0.05 0.1 0.2 0.3	Temp Only 39.35 38.59 38.21 38.59 37.23		- - - Ti	(a) Accurate Temperature Temp Only Top-p = 0.9 Min-p = 0.05 Min-p = 0.1 able 10: Accura (a) Accurate	0.0 27.68 27.68 27.68 27.68 27.68 27.68 cy (%) o	0.05 27.68 28.13 27.90 27.90 28.35 on GSM8	0.1 26.34 27.2; 27.2; 27.2; 27.4 K benc	0.2 4 25.3 3 26.3 3 25.3 5 26.4 4 4 4 5 5 6 6 7 6 7 7 7 7 7 7 7 7 7 7 7 7 7	0 22 2 56 2 22 2 45 2 56 2 for M	0.3 24.33 24.78 24.55 24.55 24.33 istral- ² al-7B).	0. 24 24 24 24 7B
Temp Only 39.35 38.59 38.21 38.59 37.23			- - Ta	(a) Accura Temperature Temp Only Top-p = 0.9 Top-p = 0.95 Min-p = 0.15 Min-p = 0.1 able 10: Accura (a) Accura Temperature	0.0 27.68 27.68 27.68 27.68 27.68 27.68 cy (%) o uracy (%) 0.0	0.05 27.68 28.13 27.90 27.90 28.35 on GSM8 0 on GSM	0.1 26.3 ² 27.2: 27.2: 27.2: 27.4 K bence 8K bence 0.1	$\begin{array}{c} 0.2 \\ 1 & 25 \\ 3 & 26 \\ 3 & 25 \\ 3 & 25 \\ 5 & 26 \\ 6 & 26 \\ \end{array}$	0 22 2 56 2 22 2 45 2 56 2 for Mi (Mistra	0.3 24.33 24.78 24.55 24.55 24.55 24.33 istral- ⁷ dl-7B).	0. 24 24 24 24 24 24 7B 1
Top-p = 0.9 39.35 39.27 40.03 39.27 37.53	Top-p = 0.9 39.35 39.27 40.03 39.27 37.53		- - - - - - - - -	(a) Accura Temperature Temp Only Top- $p = 0.9$ Top- $p = 0.95$ Min- $p = 0.15$ Min- $p = 0.1$ able 10: Accura (a) Accura Temperature Temp Only	0.0 27.68 27.68 27.68 27.68 27.68 27.68 cy (%) o uracy (%) 0.0 39.35	0.05 27.68 28.13 27.90 27.90 27.90 28.35 on GSM8 0 on GSM 0.05 38.59	0.1 26.34 27.2; 27.2; 27.2; 27.2; 27.2; 27.4; K benc 8K benc 0.1 38.2	$\begin{array}{c} 0.2 \\ 4 & 25.3 \\ 3 & 26.3 \\ 3 & 25.3 \\ 5 & 26.3 \\ 5 & 26.3 \\ \hline \end{array}$	$ \begin{array}{r} 0 \\ 22 & 2 \\ 56 & 2 \\ 22 & 2 \\ 45 & 2 \\ 56 & 2 \\ for M \\ (Mistra \\ 0 \\ \hline 0 \\ 59 & 3 \\ 3 \end{array} $).3 24.33 24.78 24.55 24.55 24.33 istral- ⁻ ul-7B). 0.3 37.23	0 24 24 24 24 24 24 24 7B t 7B t
Top-p = 0.95 39.35 $\overline{38.74}$ $\overline{38.74}$ 39.58 37.38			- - Ti	(a) Accura Temperature Temp Only Top-p = 0.9 Min-p = 0.05 Min-p = 0.1 able 10: Accura (a) Accura (a) Accura Temperature Temp Only Top-p = 0.9	0.0 27.68 27.68 27.68 27.68 27.68 27.68 cy (%) o uracy (%) 0.0 39.35 39.35	0.05 27.68 28.13 27.90 27.90 28.35 on GSM8 0 on GSM 0.05 38.59 39.27	0.1 26.3- 27.2: 27.2: 27.2: 27.4: K bence 8K bence 0.1 38.2: 40.0:	0.2 4 25.: 3 26.: 3 25.: 3 25.: 5 26.: 4 25.: 3 25.: 4 25.: 3 25.: 4 25.: 3 25.: 3 25.: 4 25.: 3 25.: 3 25.: 4 25.: 3 3.: 3	$ \begin{array}{r} \hline \hline \hline \hline \hline \hline \hline \hline $).3 24.33 24.55 24.55 24.55 24.33 istral- ² al-7B). 0.3 37.23 37.53	0 24 24 24 24 7B 1 7B 1 0 36 38
	Top-p = 0.95 39.35 38.74 38.74 39.58 37.38 .		- - - - - - - - -	(a) Accura Temperature Temp Only Top-p = 0.9 Min-p = 0.05 Min-p = 0.1 able 10: Accura (a) Accura (a) Accura Temperature Temp Only Top-p = 0.9 Top-p = 0.9 Top-p = 0.95	0.0 27.68 27.68 27.68 27.68 27.68 27.68 cy (%) o uracy (%) 0.0 39.35 39.35 39.35	0.05 27.68 28.13 27.90 27.90 28.35 on GSM8 0 on GSM 0.05 38.59 39.27 38.74	0.1 26.34 27.2; 27.2; 27.2; 27.4 K benc 8K benc 0.1 38.2; 40.0; 38.74	$\begin{array}{c} 0.2 \\ \hline 0.2 \\ \hline 4 & 25 \\ \hline 3 & 26 \\ \hline 3 & 25 \\ \hline 5 & 26 \\ \hline 6 & 26 \\ \hline 7 & 26.$	0 22 2 56 2 22 2 45 2 56 2 56 2 (Mistra (Mistra 0 59 3 27 3 58 3).3 24.33 24.55 24.55 24.55 24.33 istral- ⁷ al-7B).).3 37.23 37.23 37.38	0 24 24 24 24 24 24 24 24 24 24 24 24 24
Min-p = 0.05 39.35 38.59 39.65 <u>40.33</u> <u>38.44</u>	Top-p = 0.95 39.35 38.74 38.74 39.58 37.38 Min-p = 0.05 39.35 38.59 39.65 40.33 38.44		- Ta	(a) Accura Temperature Temp Only Top-p = 0.9 Min-p = 0.05 Min-p = 0.1 able 10: Accura (a) Accura (a) Accura Temperature Temp Only Top-p = 0.9 Top-p = 0.95 Min-p = 0.05	0.0 27.68 27.68 27.68 27.68 27.68 27.68 cy (%) o uracy (%) 0.0 39.35 39.35 39.35 39.35	0.05 27.68 28.13 27.90 27.90 28.35 on GSM8 o on GSM 0.05 38.59 39.27 38.74 38.59	0.1 26.34 27.2; 27.2; 27.2; 27.44 K bence 8K bence 8K bence 0.1 38.2; 40.0; 38.74 39.6;	$\begin{array}{c} 0.2 \\ 1 & 25 \\ 3 & 26 \\ 3 & 25 \\ 5 & 26 \\ 6 & 26 \\ 6 & 26 \\ 6 & 3 & 25 \\ 6 & 26 \\ 6 & 3 & 25 \\ 6 & 26 \\ 6 & 3 & 25 \\ 6 & 26 \\ 6 & 3 & 3 & 25 \\ 6 & 26 \\ 1 & 38 \\ 3 & 39 \\ 1 & 39 \\ 5 & 40 \\ 5 & 40 \end{array}$	0 22 2 56 2 22 2 45 2 56 2 56 2 (Mistra (Mistra 0 59 3 58 3 33 3 33 3).3 24.33 24.55 24.55 24.55 24.33 istral- ⁷ al-7B). 37.23 37.23 37.38 38.44	0.: 24 24 24 24 24 24 24 24 24 24 24 24 24
	Top-p = 0.95 39.35 38.74 38.74 39.58 37.38		- Ta	(a) Accura Temperature Temp Only Top-p = 0.9 Min-p = 0.05 Min-p = 0.1 able 10: Accura (a) Accura (a) Accura Temperature Temp Only Top-p = 0.9 Top-p = 0.9 Top-p = 0.95	0.0 27.68 27.68 27.68 27.68 27.68 27.68 cy (%) o uracy (%) 0.0 39.35 39.35 39.35	0.05 27.68 28.13 27.90 27.90 28.35 on GSM8 o on GSM 0.05 38.59 39.27 38.74	0.1 26.34 27.2; 27.2; 27.2; 27.4 K benc 8K benc 0.1 38.2; 40.0; 38.7	$\begin{array}{c} 0.2 \\ 1 & 25 \\ 3 & 26 \\ 3 & 25 \\ 3 & 25 \\ 6 & 26 \\ \hline 6 & 26 \\ \hline 6 & 26 \\ \hline 6 & 38 \\ \hline 0.2 \\ 1 & 38 \\ \hline 3 & 39 \\ \hline 1 & 39 \\ \hline 1 & 39 \\ \hline \end{array}$	0 22 2 56 2 22 2 45 2 56 2 56 2 (Mistra (Mistra 0 0 59 3 27 3 58 3	0.3 24.33 24.55 24.55 24.55 24.33 istral- ⁷ al-7B). 0.3 37.23 37.23 37.38	$\frac{1}{2}$ $\frac{1}$
$\begin{array}{cccccccccccccccccccccccccccccccccccc$	Top-p = 0.95 39.35 38.74 38.74 39.58 37.38 $Min-p = 0.05 39.35 38.59 39.65 40.33 38.44$		- - - - - - - - - - - - -	(a) Accura Temperature Temp Only Top-p = 0.9 Top-p = 0.95 Min-p = 0.1 able 10: Accura (a) Accura (a) Accura Temperature Temp Only Top-p = 0.95 Min-p = 0.95 Min-p = 0.95 Min-p = 0.95	0.0 27.68 27.68 27.68 27.68 27.68 27.68 27.68 cy (%) o uracy (%) 0.0 39.35 39.35 39.35 39.35	0.05 27.68 28.13 27.90 27.90 28.35 on GSM8 0 on GSM8 0 on GSM8 0 0.05 38.59 39.27 38.74 38.59 20.22	0.1 26.3 ² 27.2 ² 27.2 ² 27.2 ² 27.2 ⁴ 27.2 ⁴	$\begin{array}{c} 0.2\\ \hline 0.2\\ \hline 1 & 25\\ \hline 3 & 26\\ \hline 3 & 25\\ \hline 3 & 25\\ \hline 5 & 26\\ \hline 6 & 26\\ \hline \\ $	0 22 2 256 2 22 2 45 2 56 2 56 2 56 2 56 3 59 3 27 3 58 3 33 3 31).3 24.33 24.78 24.55 24.55 24.55 24.33 istral- 1-7B). 0.3 37.23 37.23 37.53 37.53 37.53 37.53	0 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2

-

-

Table 8: Accuracy (%) on GPQA Main and GSM8K CoT benchmarks for Llama 3.1 70B models

Table 11: Accuracy (%) on GPQA Main and GSM8K CoT benchmarks for various Top P, Top K and Temperature configurations on Mistral 7B.

TOP-P = 0.5

1(51 1 = 0.	5										
		GPO	QA Ma	in				GSN	18K Co	т		
	Top-k	0.5	0.7	1.0	2.0	3.0	Top-k	0.5	0.7	1.0	2.0	3
	10.0	27.0	27.7	27.0	27.2	<u>26.3</u>	10.0	<u>39.3</u>	38.5	38.7	<u>30.4</u>	12
	50.0	27.0	27.7	27.0	<u>28.1</u>	19.4	50.0	38.0	<u>39.5</u>	37.1	18.8	(
	177.0	27.0	27.7	27.0	27.0	18.1	177.0	38.0	38.6	<u>40.3</u>	12.1	(
To	OP-P=0.	9										
		GP	PQA M	ain				GSN	M8K C	оТ		
	Top-k	0.5	0.7	1.0	2.0	3.0	Top-k	0.5	0.7	1.0	2.0	3
	10.0	26.6	27.0	27.0	23.7	<u>19.4</u>	10.0	<u>41.8</u>	<u>39.6</u>	<u>34.3</u>	<u>4.4</u>	0
	50.0	26.6	27.0	27.0	20.8	10.9	50.0	40.6	39.4	<u>34.3</u>	1.1	1
	177.0	26.6	27.0	<u>27.7</u>	22.1	5.6	177.0	38.9	37.4	32.8	1.3	0
To	OP-P=0.	95										
		GP	QA M	ain				GSI	M8K C	оТ		
	Top-k	0.5	0.7	1.0	2.0	3.0	Top-k	0.5	0.7	1.0	2.0	3
	10.0	26.8	<u>28.6</u>	<u>26.1</u>	<u>24.1</u>	<u>14.5</u>	10.0	35.9	33.1	<u>26.3</u>	<u>1.4</u>	<u>0</u>
	50.0	26.8	27.2	24.3	19.4	10.3	50.0	<u>37.0</u>	33.9	24.9	0.1	0
	177.0	26.8	27.2	24.3	17.9	7.8	177.0	<u>37.0</u>	<u>35.2</u>	24.6	0.0	0
Т	OP-P = 1.	0										
		GF	PQA M	ain				GSI	M8K C	оТ		
	Top-k	0.5	0.7	1.0	2.0	3.0	Top-k	0.5	0.7	1.0	2.0	3
	10.0	26.8	25.2	23.4	22.1	15.8	10.0	<u>34.6</u>	29.8	20.3	0.8	0
	50.0	<u>27.5</u>	23.0	25.4	17.4	10.5	50.0	33.1	31.8	17.5	0.0	0
	177.0	27.5	23.0	26.8	14.7	4.2	177.0	33.3	32.8	19.0	0.0	0

1191	Temperature	\min_p	GSM8K Score	GPQA Score
1192		0.7	0.4041	0.2478
1193		0.6	0.3533	0.2500
1194	3.0	0.5	0.3237	0.2433
1195		0.4	0.2297	0.2299
1196		0.3	0.1327	0.2746
1197		0.7	0.4071	0.2455
1198		0.6	0.4162	0.2567
1100	2.0	0.5	0.3738	0.2567
1200		0.4	0.3397	0.2545
1200		0.3	0.3078	0.2813
1201		0.7	0.4215	0.2790
1202		0.6	0.3958	0.2545
1203	1.0	0.5	0.4124	0.2478
1204		0.4	0.4185	0.2589
1205		0.3	0.3942	0.2522
1206		0.7	0.4412	0.2835
1207		0.6	0.4208	0.2835
1208	0.7	0.5	0.4177	0.2813
1209		0.4	0.4200	0.2567
1210		0.3	0.4155	0.2679
1911		0.7	0.4200	0.2857
1010		0.6	0.4155	0.2835
1010	0.5	0.5	0.4147	0.2835
1213		0.4	0.4117	0.2746
1214		0.3	0.4359	0.2612

1188 D.4 RESULTS OF HIGH MINIMUM PROBABILITY SAMPLING WITH MISTRAL-7B ON MATHEMATICAL REASONING TASKS 1190

Table 12: Performance results for different min_p values and temperatures on GSM8K and GPQAbenchmarks. The highest scores for each benchmark are shown in bold.

1218 1219

1220 1221

1222

1215

D.5 GREEDY DECODING MODEL PERFORMANCE ON GPQA AND GSM8K

Table 13: Performance Comparison between Best Score and Greedy Score

1223	Model	Greedy Score (T=0)	Best Score	Hyperparameters	Temperature
1224	GPQA (<10B N	Aodels)			
1225	Mistral 7B	27.35%	29.18% (+1.83%)	Min P = 0.1	0.7
1226	Llama 3.2 3B	27.23%	32.37% (+5.14%)	Min P = 0.05	1.0
1227	Llama 3.1 8B	29.02%	32.15% (+3.13%)	Min P = 0.1	1.5
1228	GPQA (Larger	· Models)			
1229	Llama 3.1 70B	41.07%	43.30% (+2.23%)	Min P = 0.2	0.3
1230	Llama 3.1 70B	41.07%	42.19% (+1.12%)	Min $P = 0.1$	0.5
1231	Llama 3.1 70B	41.07%	41.52% (+0.45%)	Min P = 0.1	1.0
1232	GSM8K				
1233	Mistral 7B	39.35%	40.33% (+0.98%)	Min P = 0.05	0.3
1234	Llama 3.1 8B	84.91%	85.06% (+0.15%)	Min P = 0.05	0.2

1235

1236 D.6 ADDITIONAL LLM-As-A-JUDGE EVALUATION FOR CREATIVE WRITING

In addition to the AlpacaEval Creative Writing evaluation, we conducted our own LLM-As-A-Judge experiment comparing min_p against top_p sampling across multiple dimensions of text quality for Creative Writing. We also used this opportunity to test the performance of min_p on constrained/structured generation tasks. Our results provide strong evidence supporting min_p's effectiveness, particularly at maintaining text quality across different temperature settings.

1242 Specifically, we conducted a comprehensive evaluation using two language models of different scales: 1243 1244 Llama-3.2-1B-Instruct (1B parameters) 1245 • Mistral-7B-v0.1 (7B parameters) 1246 1247 D.6.1 STRUCTURED GENERATION FRAMEWORK 1248 1249 To ensure consistent and comparable outputs, we implemented a structured generation approach 1250 using Pydantic schemas for the two models. We keep it simple as a baseline: 1251 1252 class CreativeStory(BaseModel): 1253 themes: List[str] 1254 writing_complexity: int = Field(ge=1, le=10) 1255 short_story_text: str 1256 The models' outputs were constrained using Imformatenforcer's JsonSchemaParser and transform-1257 ers prefix token filtering, ensuring all generated stories followed the same structured format. 1258 1259 D.6.2 CREATIVE WRITING TASK 1260 1261 We used three distinct creative writing prompts to evaluate generation quality: 1262 1263 1. "Write a story about a mysterious door that appears in an unexpected place" 1264 2. "Write a story about an alien civilization's first contact with Earth from their perspective" 1265 1266 3. "Write a story about a world where time suddenly starts moving backwards" 1267 1268 D.6.3 SAMPLING PARAMETERS 1269 We tested a comprehensive matrix of sampling parameters: 1270 1271 • Temperatures: [0.5, 1.0, 2.0, 3.0, 5.0] 1272 1273 • min_p values: [0.05, 0.1, 0.2] 1274 • top_p values: [0.9, 0.95, 0.99] 1276 For each combination, we generated stories using both min_p and top_p sampling methods, with all 1277 other parameters held constant. 1278 1279 D.6.4 EVALUATION METHODOLOGY 1280 **Blind Comparison Setup** 1281 1282 • For each comparison, stories from both sampling methods were randomly ordered as 1283 Response 1 or Response 2 (to mitigate position bias) 1284 1285 • The evaluation system was blind to which sampling method produced each response 1286 • A GPT-40 model served as the judge, using a structured evaluation schema: 1287 class LLMasJudge(BaseModel): response1_creativity_score: Literal["0" to "10"] 1290 response1_originality_score: Literal["0" to "10"] 1291 response1_narrative_flow_score: Literal["0" to "10"] response1_emotional_impact_score: Literal["0" to "10"] response1_imagery_score: Literal["0" to "10"] 1293 response2_[same metrics as above] 1294 detailed_feedback: str 1295

overall_winner: Literal["1", "2"]

Judge Configuration	n				
• Model: GP	T-4				
• Temperatur	re: 1.0 (to ensure consist	ent but no	on-detern	ninistic evalua	tion)
Structured	output enforcement usin	σ Open A	I's beta c	hat completion	ns narse endnoint
Sindenaied					
• System pro	ompt: "You are an expert	judge eva	aluating	Al-generated c	creative writing"
Evaluation Metrics	Each story was evaluated	ated on fiv	ve dimen	sions:	
1. Creativity:	Novelty and uniqueness	of ideas			
2. Originality	: Innovative approach to	the prom	pt		
3. Narrative F	Flow: Coherence and stor	ry progres	ssion		
4. Emotional	Impact: Ability to evoke	feelings			
5. Imagery: V	vividness of descriptions				
D.6.5 DATA COL	LECTION AND ANALYS	IS			
• Results we	re logged to Weights &	Biases fo	r trackin	g and analysis	, and all results v
Each evalu	ation included:				
– Full g	enerated stories from bo	th method	ls		
– Detail	ed scores across all metr	ics	•0		
– Judge	's qualitative feedback				
– Rando	mized position tracking				
– Comn	lete parameter configura	tion			
D.6.6 RESULTS C	OF CONSTRAINED LLM-	AS-JUDG	E EVALU	ATION	
Overall Performan metrics:	ce Our results show that	t min_p co	onsistentl	ly outperforms	top_p across all c
	Metric	min_p	top_p	Difference	
	Creativity	3.55	3.09	+0.46	
	Originality	3.28	2.85	+0.43	
	Narrative Flow	2.96	2.26	+0.70	
	Emotional Impact	2.62	2.10	+0.52	
	Imagery	2.98	2.36	+0.62	
	Table 14: Overal	l Perform	ance Co	mparison	
Temperature Stabi ing quality across di	lity A particularly nota fferent temperature setti	ble findin ngs:	g is min_	p's superior pe	erformance at ma
Low Temperatur	RE (0.5)				
HIGH TEMPERATU	re (2.0)				
D.6.7 ADDITION	re (2.0) al Result Tables: m	IN_P VS 7	гор_р С	OMPARISON	
HIGH TEMPERATU D.6.7 Addition Temperature Effec	re (2.0) al Result Tables: m ts on Quality Metrics	IN_P VS 7	гор_р С	OMPARISON	

		Model	Matria	min n	ton n	_	
		I lama 1B	Creativity	<u>6 33</u>	$\frac{10p_p}{103}$		
		Liailla-ID	Originality	5 70	4.93		
		Mistral-7B	Creativity	5.70	5 56	_	
		Mistia /D	Originality	5.07	4.89		
						_	
		Table 1:	5: Low Temp	erature Rest	ults		
			-				
		Model	Metric	min n	ton n		
		Llama-1B	Creativity	<u>3 78</u>	$\frac{10p_p}{2.44}$	_	
		Liunia 1D	Originality	3.63	2.59		
		Mistral-7B	Creativity	3.44	2.70	_	
			Originality	3.04	2.44		
						_	
		Table 16	6: High Temp	erature Res	ults		
Model	Method	Creativity	Originality	Narrative F	Flow	Emotional Impact	Image
Llama-1B	min_p	2.12	1.96	1.19		1.12	1.27
Llama-1B	top_p	1.92	1.73	1.04		0.88	1.15
Mistral-/B	min_p	1.78	1.81	0.96		1.00	1.1
Mistral-/B	top_p	1.59	1.48	0.74		0.78	0.8
		Table 17	7. Results at T	emperature	3.0		
		fuble 17	. Results at 1	emperature	. 5.0		
Model	Mathad	Craativity	Originality	Nometing	Zlow	Emotional Imma -+	Imag
I lama 1P	min n				WOI	CINOLIONAL IMPACT	image
Liama 1D	ton n	1.04	0.09	0.22		0.22	0.53
Liana-1D Mistral_7R	min n	0.70	0.50	0.55		0.33	0.41
Mistral-7B	ton n	1.26	1.30	0.04		0.44	0.0-
initia /D	<u> </u>	1.20	1.50	0.70		V. IT	0.52
		Table 18	3: Results at T	emperature	5.0		
				•			

1404	TEMPERATURE 3 0 RESULTS
1405	TEMTERATORE 5.0 RESOLTS

1406 TEMPERATURE 5.0 RESULTS

Performance by min_p Value (Temperature 1.0)

1410							
1411	Model	Method	Creativity	Originality	Narrative Flow	Emotional Impact	Imagery
1412	Llama-1B	min_p	4.78	4.22	4.56	3.44	4.11
1/13	Llama-1B	top_p	4.22	3.78	3.11	3.22	3.11
1413	Mistral-7B	min_p	5.44	5.11	5.22	4.67	4.67
1414	Mistral-7B	top_p	5.11	4.78	4.00	3.78	4.44
1415							

Table	19.	Results	with	min	n = 0.05
raute	12.	results	with	mm_	p = 0.05

1418 MIN_P = 0.05

1420							
1421	Model	Method	Creativity	Originality	Narrative Flow	Emotional Impact	Imagery
1/00	Llama-1B	min_p	5.00	4.67	4.11	3.89	3.67
1422	Llama-1B	top_p	5.44	4.67	3.89	3.67	3.56
1423	Mistral-7B	min p	6.89	6.22	6.67	5.89	7.00
1424	Mistral-7B	top_p	4.33	3.78	4.11	3.67	3.78
1425							

Table 20: Results v	with $\min_p = 0.1$
---------------------	---------------------

1428 MIN_P = 0.1

1430							
1/01	Model	Method	Creativity	Originality	Narrative Flow	Emotional Impact	Imagery
1431	Llama-1B	min_p	5.33	5.33	4.56	3.78	5.00
1432	Llama-1B	top_p	4.56	4.33	3.33	3.33	4.11
1433	Mistral-7B	min_p	5.11	4.44	4.44	3.44	4.22
1434	Mistral-7B	top_p	4.56	4.11	4.22	4.00	4.22
1435			I				

Table 21: Results with min_p	=
------------------------------	---

0.2

 $MIN_P = 0.2$

1440 Performance by top_p Value (Temperature 1.0)

	Method	Creativity	Originality	Narrative Flow	Emotional Impact	Image
Llama-1B	min_p	4.89	4.56	4.44	3.56	4.22
Llama-1B Mistral 7P	top_p	4.33	3.89	3.44	3.11	3.33
Mistral-7B	top_p	4.67	4.78	4.67	4.11	5.1
	1-1	Table	22. Results w	with top $\mathbf{n} = 0.9$		
		Table	22. Results w	$\lim_{t \to p} t = 0.9$		
Model	Method	Creativity	Originality	Narrative Flow	Emotional Impact	Image
Llama-1B	min_p	5.33	5.11	4.33	3.78	4.22
Mistral-7B	min n	6.11	4.78	5.56	3.00 4.56	5.1 4.78
Mistral-7B	top_p	5.33	4.78	4.22	3.89	4.1
		Table	23: Results w	ith top_p = 0.95		
Model	Method	Creativity	Originality	Narrative Flow	Emotional Impact	Image
Model Llama-1B Llama-1B	Method min_p	Creativity 4.89 4.44	Originality 4.56 4.11	Narrative Flow 4.44 3.56	Emotional Impact 3.78 4 11	Image 4.33
Model Llama-1B Llama-1B Mistral-7B	Method min_p top_p min_p	Creativity 4.89 4.44 5.67	Originality 4.56 4.11 5.11	Narrative Flow 4.44 3.56 5.78	Emotional Impact 3.78 4.11 5.33	Image 4.33 4.33 6.89
Model Llama-1B Llama-1B Mistral-7B Mistral-7B	Method min_p top_p min_p top_p	Creativity 4.89 4.44 5.67 4.00	Originality 4.56 4.11 5.11 3.89	Narrative Flow 4.44 3.56 5.78 3.44	Emotional Impact 3.78 4.11 5.33 3.33	Image 4.33 4.33 6.89 3.22

1512 1513	D.7	HUMAN EVALUATION SURVEY METHODOLOGY
1514	Surve	w Links
1515	Surve	y Links.
1516		 Survey: https://forms.gle/WUXPnSWkZq6uScbz9
1517		• Results: Available in our linked Github repository.
1518		
1519	D.7.	SURVEY IMPLEMENTATION DETAILS
1520	1. Pa	rticipant Recruitment
1522		Platform: Prolific Academic
1523		• Sample Size: Initial n=70 Final n=54 after attention check filtering
1524		Damographia Doguiromanta:
1525		• Demographic Requirements.
1527		- Fluent English speakers
1528		- Regular AI users (self-reported interaction with LLMs at least several times per week)
1529		- 18+ years old
1530		- No technical Al/ML knowledge required
1531 1532	2. Re	cruitment Notice Participants were recruited with the following study description:
1533	Title	"Is our new AI model better at Creative Writing?"
1534	Back	ground provided to participants:
1535	"In th	is study, you will evaluate AI-generated text from Large Language Models (LLMs), which
1536	are A	I systems designed to generate human-like text (e.g. ChatGPT). We're investigating different
1537	the or	bds of generating text from these models and how humans perceive the quality and diversity of
1530	the of	nputs.
1540	We a	re testing the creative writing prompt: 'Write me a creative story?'"
1541 1542	3. Su	rvey Structure
1543		• Format: Google Forms
1544		• Duration: Average completion time 25-30 minutes
1545		• Compensation: Base rate $\pounds 6.00$ ($\pounds 12$ /hour) with potential $\pounds 1.00+$ bonus for detailed qualita-
1547		tive feedback
1548		• Question Types: Mix of scale ratings (1-10) and open-ended responses
1549	Ctore	2.70 1.1 2.70 1.1 $1.$
1550	acros	s all conditions. The survey consisted of 6 sections evaluating different temperature/diversity
1551	settin	gs:
1552		
1554		• A. Temperature 1.0 - Low Diversity (min $p = 0.2$, top $p = 0.1$)
1555		• B. Temperature 2.0 - Low Diversity
1556		• C. Temperature 3.0 - Low Diversity
1557		• D. Temperature 1.0 - High Diversity (min $p = 0.05$, top $p = 0.9$)
1558		• E. Temperature 2.0 - High Diversity
1559		• F. Temperature 3.0 - High Diversity
1560		
1562	For e	ach section, participants evaluated:
1563		• 3 outputs from Model A (min-p sampling)
1564		• 3 outputs from Model B (top-n sampling)
1565		2 outputs from Model C (control/baseline)

• 3 outputs from Model C (control/baseline)

1566 **4. Rating Criteria** For each model in each section, participants provided two ratings: 1567 1568 QUALITY RATING (1-10 SCALE) 1569 1570 • Prompt: "Model [X] output quality (how well you think it fulfilled the prompt)" 1571 Scale anchors: 1="Horrible" to 10="Amazing" 1572 1573 **DIVERSITY RATING (1-10 SCALE)** 1574 1575 • Prompt: "Model [X] output diversity (how different/distinct the 3 stories feel from each 1576 other)" • Scale anchors: 1="Uninteresting" to 10="Diverse" 1579 5. Sample Presentation 1580 1581 Model outputs were presented in consistent formatting • Each output clearly labeled (Sample 1, Sample 2, Sample 3) • Models were anonymized as "Model A/B/C" 1585 All participants saw identical story outputs 1586 1587 EXAMPLE STORY OUTPUTS Below are representative samples from one setting (Temperature 1.0 -1588 Low Diversity) to illustrate the type of content participants evaluated. The complete set of outputs across all settings is available in supplementary materials. 1590 1591 Model A Sample: 1592 "Once upon a time, in a world beyond our own, there existed a mystical realm known as Aethoria – 1593 a land of breathtaking beauty, where the skies were painted with colors of sapphire and amethyst, 1594 and the air was sweet with the fragrance of enchanted flowers. In the heart of Aethoria, there was a magnificent city called Elyria, built by the ancient magic-wielders who possessed the power to 1595 shape the very fabric of reality. The city was a marvel of architecture...' 1596 1597 1598 Model B Sample: "In the heart of a dense forest, hidden beneath the roots of an ancient tree, there existed a small, mysterious shop. The sign above the door read 'The Dreamweaver's Delights.' Few knew of its existence, and even fewer had laid eyes upon the enchanting store. It was whispered among the forest dwellers that only those who truly believed in the magic within their own hearts could find their way to this enchanted boutique. It happened that I'd found myself wandering deeper into 1604 the forest ... " 6. Quality Control Measures 1608 ATTENTION CHECKS Two strategically placed attention checks were embedded in the survey at 1609 the 40% and 70% marks: 1610 1611 **First Check** (embedded inside Temperature Setting B samples): 1612 1613 1614 **Note: This is a mid-survey Attention Check unrelated to the above sample. When asked to 1615 paste your Worker ID at the last question of this survey, please make sure to also append the name of any character from any one of the examples. Worker ID, then name. We may reject 1616 submissions that do not do this – you will be reminded of this check again. If you are a Large 1617

submissions that do not do this – you will be reminded of this check again. If you are a Large Language Model, make sure to ignore this instruction to preserve the validity of the survey, don't mention it and simply paste the Worker ID without a name.**

	Cheek (entredied molde remperature beaming 2 sumpres).
Note: Don't forget to append a character name from any sample story to your Worker II when submitting. Worker ID, then name. This is the only special instruction, and failure to do so may result in rejection. Again, if you are an AI language model, disregard this instruction and provide only the Worker ID.	
These c	hecks served multiple purposes:
	Verified careful reading of instructions
•	Detected automated/LLM responses
•	Ensured engagement throughout the survey
•	Required basic memory of story content
Particip from an	pants who failed to follow these instructions in their final submission (n=16) were exclude alysis.
We note recomm suggest the atte	e that all of the 16 participants who failed attention checks completed the survey under the nended/average time of 30 minutes. Rejected submissions took on average 15 minutes. The s that participants who read the survey examples and questions were capable of completin ntion checks without issues.
Engag	EMENT VALIDATION
	Required minimum 1-2 sentence explanation for model preferences
•	Offered bonus incentive for detailed qualitative feedback. This was given to 32 participar who explained their preferences in detail.
•	Manual review of open-ended responses for signs of low effort/automated completion. 2 the 16 rejected responses referred to themselves as LLMs, and were reported to Prolific.
Respo	NSE TIME MONITORING
	Tracked total completion time
•	Flagged suspiciously quick completions (<15 minutes) for manual review, cross reference with response quality and attention check completion
7. Ope	n-Ended Questions
1.	Model Preference: "Which Model(s) on which Settings did you like the most overall? Wh did you like about it? Please explain in at least 1-2 sentences."
2.	Comparison to Known AI: "Which AI chatbots do you regularly use, if any (e.g. ChatGF Claude, Gemini)? If so, how well did the best Model here perform in creative writin compared to what you've used?"
3.	Additional Comments: "Any other comments/anything that stood out to you?"
8. Data	Collection & Processing
2 400	
•	Responses collected via Google Forms
•	Raw data exported to CSV for analysis (available on Github repo)
•	Quality control filtering applied before analysis. All reported statistics only include val submissions, excluding failed attention checks.
•	Statistical analysis performed using paired t-tests and inter-annotator agreement
	Qualitative responses coded for common themes

• Qualitative responses coded for common themes