DiffSampling: Enhancing Diversity and Accuracy in Neural Text Generation

Anonymous ACL submission

Abstract

002 Despite their growing capabilities, language models still frequently reproduce content from their training data, generate repetitive text, and favor common grammatical patterns and vocabulary. A possible cause is the decoding strategy: the most common strategies either 007 consider only the most probable tokens, which reduces output diversity, or increase the likelihood of unlikely tokens, compromising output accuracy and correctness. In this paper, we propose three new decoding methods that lever-013 age a mathematical analysis of the token probability distribution to ensure the generation of contextually appropriate text. In particular, the difference between consecutive, sorted probabilities can be used to truncate incorrect tokens. 017 Experiments concerning math problem solving, extreme summarization, and the divergent as-019 sociation task demonstrate that our approach consistently performs at least as well as existing methods in terms of quality and diversity.

1 Introduction

037

041

In recent years, large language models (LLMs) have demonstrated remarkable performance (Bubeck et al., 2023), driven by the availability of large-scale datasets, advances in computational power (Bommasani et al., 2021), and the development of innovative learning strategies (e.g., Stiennon et al., 2020; Rafailov et al., 2023). While training provides LLMs with the information and skills required to process natural language, another aspect plays a key role at generation time: the decoding strategy, that is, the method used to extract text sequences from the model. The choice of decoding scheme significantly impacts the generated output, as there is a pronounced trade-off between quality and diversity (Ippolito et al., 2019). The most straightforward strategies, such as greedy decoding (always selecting the highest-probability token) or sampling, tend to



Figure 1: A graphical representation of the effects of our *DiffSampling* methods. In the top-left square, the original distribution. In the top-right square, *DiffSampling-cut* truncates after the minimum discrete derivative. In the bottom-left square, *DiffSampling-lb* also imposes a total probability lower bound p_{lb} (here $p_{lb} = 0.8$). In the bottom-right square, *DiffSampling-minp* applies truncation only among tokens with a probability less than p_{min} times the highest probability (here $p_{min} = 0.3$).

repeat the same tokens multiple times (Su et al., 2022), reproduce training data (Carlini et al., 2021; Franceschelli et al., 2024), or flatten the lexicon in favor of the most common grammatical structures and words (Fleisig et al., 2024; Reviriego et al., 2023). Although the temperature parameter may increase the likelihood of less frequent tokens, it also raises the chance of syntactically incorrect ones by flattening their probabilities, regardless of their actual ranking. An ideal solution should concentrate on where the critical mass of the probability distribution resides. More precisely, with critical mass, we refer here to the portion of the probability distribution that collectively accounts for the majority of the probability mass of the tokens. In this direction, a common approach is nucleus sampling (Holtzman et al., 2020), which removes the tail of the distribution by focusing

042

043

045

047

048

051

on the smallest subset of tokens whose global probability exceeds a given threshold. However, key issues remain: first, nucleus sampling is sensitive to the choice of the threshold; second, it can either preserve incorrect tokens or exclude appropriate ones if the critical mass is smaller or larger than the threshold, respectively. As suggested by Hewitt et al. (2022), the learned probability distribution can be viewed as a mixture of the true distribution, which assigns a non-zero probability only to appropriate tokens (the critical mass), and a smoothing distribution, which assigns a small but non-zero probability to incorrect tokens. This smoothing is necessary for learning purposes.

060

061

062

065

072

074

075

077

078

079

081

084

100

102

103

104 105

107

109

In this paper, we introduce a family of decoding strategies called *DiffSampling*, based on the analysis of the probability distribution of the tokens, and in particular, on the minimum discrete derivative (i.e., the largest difference between consecutive probabilities in a sorted distribution). We propose a method for excluding incorrect tokens introduced by the smoothing distribution, along with two relaxed variants designed to promote output diversity (see Figure 1). We then provide a comprehensive evaluation of them under three different tasks (namely, math problem solving, extreme summarization, and the divergent association task) and discuss their advantages and limitations. We show that DiffSampling consistently performs better in either quality or diversity.

The remainder of this paper is structured as follows. First, we introduce the decoding problem from a neural language model perspective and discuss existing approaches (Section 2). Then, we present our discrete derivative-based sampling strategy and three different methods to employ it (Section 3). Finally, in Section 5 we evaluate our methods on mathematical problem-solving tasks, extreme summarization, and the divergent association task against the most common baselines, finding that *DiffSampling* is a simple yet effective way to generate appropriate and diverse text.

2 Background

Language Modeling 2.1

An autoregressive language model (LM) is a probability distribution $p_{\theta}(\mathbf{x})$ parameterized by θ over a variable-length text sequence $\mathbf{x} = (x_1 \dots x_T)$, 106 where T is the sequence length and each token x_t is in a finite vocabulary \mathcal{V} . The probability distribution is factorized as $p_{\theta}(\mathbf{x})$ =

 $\prod_{t=1}^{T} p_{\theta}(x_t | x_1 \dots x_{t-1})$, and the LM is usually trained to maximize the likelihood of the true distribution $p^*(\mathbf{x})$ for any \mathbf{x} from a reference dataset (the training set). In other words, given as input $x_1 \dots x_t$, the model learns to approximate the probability of each token from \mathcal{V} being x_{t+1} . While this makes the model immediately capable of scoring the probability of a given text, it also allows for the generation of new sentences. Given a commonly human-written prefix (also known as a prompt) $\mathbf{x} = (x_1 \dots x_P)$ of length P, we can decode a continuation $\hat{\mathbf{x}} = (x_{P+1} \dots x_{T+P})$ from the LM through its factorized representation introduced before. However, we must remember that the model is trained to score and not to generate sentences. A given text might have zero probability for generation purposes (e.g., the text is syntactically incorrect), but non-zero probability for ranking purposes (Hewitt et al., 2022).

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

2.2 Decoding Strategies

The decoding of tokens from the probability distribution learned by a neural language model can occur in several ways. The greedy strategy involves selecting the most probable token each time. However, this can lead to a consistent lack of diversity and several repetitions. The standard approach involves sampling from the probability distribution, which can be transformed through a temperature parameter τ . The temperature scales the differences among the various probabilities: a temperature lower than 1 will increase the probability of the most-probable tokens (a zero temperature degenerates to greedy strategy), while a temperature higher than 1 will increase the probability of the least-probable tokens, allowing for more diversity in generation (Peeperkorn et al., 2024). However, this might lead to the selection of tokens that are not syntactically appropriate for the current input. Top-k sampling (Fan et al., 2018) reduces the token space to the k most probable ones.

To generate more natural and coherent solutions, contrastive search (Su et al., 2022) employs a top-ksampling method combined with a degeneration penalty. This promotes the selection of tokens that differ from those already generated, enhancing the diversity and quality of the output. Nevertheless, limiting the number of selected tokens a priori can lead to the exclusion of meaningful tokens or the inclusion of inappropriate ones. A possible solution is to set k dynamically, as in Mirostat (Basu et al., 2021): to maintain the perplexity of generated text

256

257

259

211

212

213

214

at a desired value, the k parameter is actively tuned based on the current cross-entropy.

161

162

163

164

166

169

170

171

174

175

176

177

178

179

180

181

182

183

187

192

193

194

197

198

199

203

206

207

210

Alternatively, nucleus (or top-p) sampling (Holtzman et al., 2020) reduces the token space to the smallest subset of tokens with a total probability no less than p. To restrict the nucleus to tokens whose information content is close to the expected one given prior context, locally typical sampling (Meister et al., 2023) focuses on the tokens with negative log-probability within a certain absolute range from the conditional entropy (and a total probability no less than p). Finally, Hewitt et al. (2022) assert that a language model learns a mixture of the true token distribution and a smoothing distribution to avoid infinite perplexity. For desmoothing the distribution, they propose ϵ - and η sampling, which truncate tokens with a probability smaller than a threshold set a priori or dynamically through the entropy of the distribution, respectively. This threshold can also be set according to the magnitude of the highest probability as in min-p (Minh et al., 2025). However, such strategies do not guarantee the exclusion of the smoothing-induced tail. Contrastive decoding (Li et al., 2023) leverages the difference in likelihood between a large language model and a smaller, less capable one to prioritize tokens with sufficiently high probability under the expert model. However, it requires access to a smaller model with an identical vocabulary, which is not always available. While conceptually aligned, our method simplifies the threshold computation and provides more intuitive guarantees on the suitability of selected tokens.

3 DiffSampling

Given the probability distribution of the next token, let us imagine sorting it to have tokens in descending order based on their probability. Following Hewitt et al. (2022), only the first D tokens have a positive probability under the true token distribution, while the remaining $|\mathcal{V}| - D$ tokens receive a non-zero final probability solely due to the smoothing distribution, which prevents infinite perplexity. To generate correct text, we need to limit our sampling among the first D tokens, thus, we need to identify a cutting point that is as close as possible to the D-th token. We propose to achieve this by truncating after the largest difference between probabilities: the token to its left should be the least probable token that our model considers correct.

From a mathematical analysis perspective, this

point is characterized simply and elegantly as the location where the derivative reaches its minimum. Let us consider a probability distribution p(x) defined for a limited number of $x_1 \dots x_N$, with p() monotonically decreasing. According to the forward difference approximation, the discrete derivative of a function $f(x_n)$ is defined as $\Delta f(x_n) = f(x_{n+1}) - f(x_n)$, thus we have:

$$\Delta p(x_n) = \begin{cases} p(x_{n+1}) - p(x_n) & \text{if } n < N\\ -p(x_n) & \text{if } n = N \end{cases}$$
(1)

which is always non-positive. $\operatorname{argmin}(\Delta p(x_n))$ represents the index of the last token before the point characterized by the largest difference between consecutive probabilities.

In particular, it seems plausible that $\operatorname{argmin}(\Delta p(x_n)) \leq D$, i.e., it either marks the point where the true distribution ends and smoothing begins to take effect, or a point within the true distribution that separates tokens with significantly higher probabilities from the rest. Indeed, due to the inner nature of smoothing, it seems unreasonable that the maximum difference is between tokens with zero probability under the true distribution, and thus only because of the smoothing distribution.

Building on this intuition, we propose *DiffSampling*, a family of three decoding strategies. The first one, which we call *DiffSampling-cut*, leverages the aforementioned property and consists of cutting the distribution tail at the right side of the minimum discrete derivative, i.e., sampling among the tokens $x_i, i \leq \operatorname{argmin}(\Delta p(x_n))$. Due to the guarantee of selecting a correct token, this approach can be seen as an improved greedy strategy: when the model has high confidence in a single token, it degenerates into the greedy strategy; otherwise, it preserves other appropriate tokens, increasing diversity.

Since the minimum discrete derivative should guarantee the correctness of the truncation, *any* preserved token should come from the true distribution: we can sample at a higher temperature to foster diversity without the usual trade-off with quality. Note that although temperature scaling is typically applied before truncation, doing so alters the probability distribution, potentially shifting the minimum of the discrete derivative forward possibly into the region of tokens that have zero probability under the true distribution. To preserve the mathematical properties discussed above, we instead apply temperature scaling *after* truncation.

However, as previously discussed, this cutoff 260 point can fall within the true distribution, thereby 261 excluding tokens that are still correct; a frequent scenario consists of the first token minimizing $\Delta p(x_n)$, but still having a quite low probability. To address this issue, we propose two relaxations 265 to *right-move* the truncation. The first one builds upon top-p sampling and introduces a lower bound on the saved mass probability. DiffSampling-lb considers truncating based on $\Delta p(x_n)$ in such a 269 way that the resulting tokens have a total probability at least equal to the lower bound p_{lb} . In other 271 words, given k cardinality of the smallest subset of tokens whose total probability is not lower than p_{lb} , 273 it computes the $\operatorname{argmin}(\Delta p(x_n))$ for $n \geq k$ (i.e., 274 the cutting point is between tokens not included in the top-p nucleus). This approach can be seen as an improved top-p sampling: it *corrects* the p parameter via our derivative-based approach to include 278 279 appropriate tokens after the selected nucleus.

> Alternatively, we can build upon min-p sampling by introducing a dynamic upper bound on the probability of truncated tokens. DiffSampling-minp considers truncating based on $\Delta p(x_n)$ in such a way that no discarded tokens have a probability greater than $p_{min} \cdot p(x_0)$. In other words, given j index of the lowest-probable token with a probability greater than $p_{min} \cdot p(x_0)$, it computes the $\operatorname{argmin}(\Delta p(x_n))$ for $n \geq j$. This approach can be seen as an improved min-p sampling: if there are tokens after index j with a probability very close to the threshold, it still preserves them.

> Overall, *DiffSampling* can be seen as a sampling scheme governed by two parameters, i.e., the probability-mass lower bound p_{lb} and the truncated probability upper bound p_{min} (where *DiffSampling-cut* just assumes a value of 0.0 for the first and of 1.0 for the second), plus the additional temperature τ . The full algorithm is reported in Algorithm 1.

4 Illustrative Example

281

290

291

296

297

301

309

To make it easier to understand the advantages of our methods, Table 1 reports an illustrative example of them compared with their most similar methods. For the sake of simplicity, top-*p* and *DiffSampling-lb* consider the same $p = p_{lb} = 0.9$, while min-*p* and *DiffSampling-minp* consider the same $p = p_{min} = 0.1$. As apparent, *DiffSamplingcut* improves upon the greedy strategy by also considering the second-most probable token, while both *DiffSampling-lb* and *DiffSampling-minp* im-

Algorithm 1 DiffSampling

Input: probabilities probs = $[p_1 \dots p_N]$, lower bound p_{lb} , upper bound p_{min} , temperature τ . sorted_probs, indices = sort(probs) $fwd_probs = sorted_probs[1:] + [0.]$ delta_probs = fwd_probs - sorted_probs if $p_{min} > 0.0$ then $th = p_{min} \cdot \text{sorted}_{probs}[0]$ $lim = \operatorname{argmin}(\operatorname{sorted_probs} > th) - 1$ $delta_probs[:lim] = 0.$ else $nucleus = cumsum(sorted_probs) < p_{lb}$ delta probs[nucleus] = 0. end if $cut_i dx = \operatorname{argmin}(\mathsf{delta_probs})$ sorted_probs[$cut_idx + 1$:] = 0. probs = sort_by_idx(sorted_probs, indices) $\log its = \log(probs/sum(probs))/\tau$ probs = softmax(probs)Output: probs

prove upon top-p and min-p by not discarding tokens with very similar probability compared to preserved ones (for example, top-p would discard the 'read' token while having only a 0.014% probability less than ',').

Prompt: Natural language generation (NLG) is the subfield of artificial intelligence and computational linguistics that is concerned with the construction of computer systems that can

101115 1110						
Token	Prob	$\mathbf{Top-}p$	Min-p	D-cut	D-lb	D-minp
generate	37.326	41.366	50.872	59.886	40.929	47.537
produce	25.002	27.709	34.076	40.114	27.416	31.842
understand	7.295	8.084	9.942	-	7.999	9.290
create	3.749	4.154	5.109	-	4.110	4.774
naturally	2.797	3.100	-	-	3.067	3.562
perform	2.352	2.606	-	-	2.579	2.995
reason	1.067	1.182	-	-	1.170	-
be	0.956	1.060	-	-	1.048	-
			-	-		-
recognize	0.350	0.388	-	-	0.384	-
,	0.339	0.375	-	-	0.371	-
read	0.325	-	-	-	0.357	-
respond	0.321	-	-	-	0.352	-
interpret	0.318	-	-	-	0.348	-
interact	0.259	-	-	-	-	-

Table 1: Token probability comparison between topp, min-p, and our methods, showing how they avoid treating differently tokens with very similar probabilities (reported in **bold**). The probabilities are taken from SmolLM-135M-Instruct (Ben Allal et al., 2024).

5 Experiments

To evaluate whether *DiffSampling* helps diversify outputs while maintaining a high accuracy, we test

315

316

407

408

409

410

411

412

413

365

366

it on three case studies: math problem solving, text 318 summarization, and the divergent association task¹. 319 While slightly unconventional, these tasks are very different from each other, and provide meaningful ways to evaluate diversity and quality together, as they have quantifiable goals which can be reached 323 in syntactic and semantic different ways.

Models and Baselines 5.1

324

325

327

328

329

333

335

341

345

347

351

361

362

364

In all our experiments, we start from a state-ofthe-art LLM and test various decoding strategies. For the math problem-solving tasks, we use the Llama2-based MetaMath model trained with selfsupervised learning on MetaMathQA (Yu et al., 2024). Following Chhabra et al. (2024), for extreme text summarization we use the Llama2-7B model (Touvron et al., 2023), considering both RLHF-instructed and pre-trained versions. Finally, for the divergent association task, we consider Llama3-8B (Grattafiori et al., 2024), using both DPO-tuned and pre-trained versions. We study the performances of our three methods: DiffSamplingcut; DiffSampling-lb with $p_{lb} = 0.8$, which results in minimal accuracy loss while enhancing diversity compared to lower values for the tasks at hand (see Appendix D.2); *DiffSampling-minp* with $p_{min} = 0.3$, which we found provides an increase in quality without significant loss in diversity for the tasks taken into consideration (see Appendix D.3). We compare them with a total of 6 decoding-based baselines: greedy strategy; contrastive search (with top-k = 8 and the scaling factor of the degeneration penalty $\alpha = 0.6$); η -sampling (with $\eta = 0.0003$); locally typical sampling (with p = 0.9); top-p sampling (with p = 0.9; and min-p sampling (with p = 0.1). While other methods, such as contrastive decoding and beam search, could also be considered, we restrict our analysis to sampling-based methods to ensure a fair comparison, selecting those with similar computational costs and operational principles to our approach.

Math Problem Solving 5.2

5.2.1 Experimental Setup

Solving math problems provides a useful case study for our decoding strategies, as it allows us to evaluate the correctness of solutions (as the percentage of correctly solved problems) and the diversity of

procedures to arrive at the result. In particular, we consider the GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021) test sets; the relative prompts are reported in Appendix B. To avoid resource wasting, we focus on entries with a problem and a solution whose tokenized versions are no longer than 512.

We evaluate the quality of solutions through the ratio of correctly solved problems. Instead, the diversity is computed according to two methods: expectation-adjusted distinct N-grams (EAD) (Liu et al., 2022) and sentence embedding cosine diversity (SBERT) (Hong et al., 2024), which should evaluate syntactic and semantic diversity, respectively (Kirk et al., 2024). EAD counts the number of distinct N-grams tokens (averaging over $N = 1 \dots 5$) and removes the bias toward shorter inputs by scaling the number of distinct tokens based on their expectations. The SBERT metric is 1 minus the cosine similarity between the embeddings of the sentences. While originally based on Sentence-BERT (Reimers and Gurevych, 2019), we employ the more recent all-mpnet-base-v2 to obtain the embeddings, as suggested by their devel $opers^2$. Following Kirk et al. (2024), we compute cross-input EAD and SBERT, i.e., by considering the set of all outputs produced for a specific seed. In addition, we also compute against-greedy EAD and SBERT. Given each input, we compare the output with the greedy one by calculating the average expectation-adjusted distinct N-grams not present in the greedy response, and 1 minus the cosine similarity between the two outputs, respectively. Finally, for a more fine-grained analysis, Appendix E.1 reports a few examples of generated outputs.

5.2.2 Experimental Results

Table 2 (left side) reports the results for the GSM8K test set. DiffSampling-lb achieves the highest average accuracy. Among the baselines, only locally typical sampling performs comparably, while all our three methods outperform the others in accuracy. Regarding diversity, *DiffSampling-cut* is the closest to greedy, while *DiffSampling-lb* is in line with the sampling-based baselines.

Table 2 (right side) reports the results for the MATH test set. Here, the highest accuracy is reached by *DiffSampling-cut*, which also improves on the greedy strategy in terms of diversity, closely followed by DiffSampling-minp, which, on the

¹The code used for the experiments is available here: https://anonymous.4open.science/r/DiffSampling

²https://huggingface.co/sentence-transformers/ bert-large-nli-stsb-mean-tokens

Dataset:			GSM8K					MATH		
Method	Accuracy	Cross	-Input	Against	-Greedy	Accuracy	Cross	-Input	Against	-Greedy
		EAD	SBERT	EAD	SBERT		EAD	SBERT	EAD	SBERT
Greedy	$66.44_{\pm.09}$	$2.03_{\pm.00}$	$0.64_{\pm.00}$	-	-	$20.62_{\pm.01}$	$5.65_{\pm.00}$	$0.80_{\pm.00}$	-	-
Contrastive search	$65.88_{\pm.59}$	$2.06_{\pm.00}$	$0.64_{\pm.00}$	$0.17_{\pm.00}$	$0.02_{\pm.00}$	$21.05 \pm .14$	$5.82_{\pm.01}$	$0.80_{\pm.00}$	$0.31_{\pm.00}$	$0.09_{\pm.00}$
Top- p sampling	$65.00_{\pm.18}$	$2.08_{\pm.01}$	$0.64_{\pm.00}$	$0.23_{\pm.00}$	$0.03_{\pm .00}$	$20.02_{\pm.12}$	$6.08_{\pm.02}$	$0.80_{\pm .00}$	$0.36_{\pm.00}$	$0.10_{\pm .00}$
η -sampling	$65.05_{\pm.19}$	$2.12_{\pm.00}$	$0.64_{\pm.00}$	$0.25_{\pm.00}$	$0.04_{\pm.00}$	$19.67_{\pm.20}$	$6.36_{\pm.01}$	$0.80_{\pm .00}$	$0.39_{\pm.00}$	$0.11_{\pm.00}$
Locally typical	$66.29_{\pm.55}$	$2.09_{\pm.00}$	$0.64_{\pm.00}$	$0.23_{\pm.00}$	$0.03_{\pm .00}$	$19.95_{\pm.26}$	$6.06_{\pm.01}$	$0.80_{\pm.00}$	$0.36_{\pm.00}$	$0.10_{\pm.00}$
Min-p sampling	$65.76_{\pm.44}$	$2.09_{\pm.00}$	$0.64_{\pm.00}$	$0.23_{\pm.00}$	$0.03_{\pm .00}$	$20.25 \pm .09$	$6.09_{\pm.01}$	$0.80_{\pm.00}$	$0.36_{\pm.00}$	$0.10_{\pm .00}$
DiffSampling-cut	$66.36_{\pm.23}$	$2.04_{\pm.00}$	$0.64_{\pm.00}$	$0.14_{\pm.00}$	$0.02_{\pm.00}$	$21.38_{\pm.20}$	$5.71_{\pm.01}$	$0.80_{\pm.00}$	$0.27_{\pm.00}$	$0.07_{\pm .00}$
DiffSampling-lb	$66.92_{\pm.08}$	$2.07_{\pm.00}$	$0.64_{\pm.00}$	$0.20_{\pm.00}$	$0.03_{\pm .00}$	$20.78 \pm .14$	$6.00_{\pm.01}$	$0.80_{\pm.00}$	$0.35_{\pm.00}$	$0.10_{\pm .00}$
DiffSampling-minp	$66.44_{\pm.35}$	$2.05_{\pm.00}$	$0.64_{\pm.00}$	$0.19_{\pm.00}$	$0.03_{\pm.00}$	$21.13_{\pm.08}$	$5.87_{\pm.01}$	$0.80_{\pm.00}$	$0.33_{\pm.00}$	$0.09_{\pm.00}$

Table 2: Accuracy and diversity of results for the GSM8K and MATH test sets over 3 seeds. The mean and standard error of the final score for each run are reported for accuracy and cross-input diversity, whereas the mean and the 95% confidence interval for the full set of answers are reported for against-greedy diversity.

other hand, has slightly lower performance on diversity compared to the other baselines (apart from contrastive search). Finally, *DiffSampling-lb* has diversity scores in line with top-p, min-p, locally typical, and η -sampling, but with a consistently higher accuracy.

5.3 Extreme Summarization

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

5.3.1 Experimental Setup

Summarizing paragraphs represents another meaningful case study since the same text can be correctly outlined in different ways. To keep the resource consumption as low as possible, we consider the eXtreme Summarization (XSum) dataset (Narayan et al., 2018), which contains pairs of documents and one-sentence summaries. In particular, we use the test partition (11334 entries) and exclude all entries with a tokenized document longer than 768, obtaining 9815 entries; then, we limit our experiment to 1000 random samples, and we use the prompt suggested by Chhabra et al. (2024) and reported in Appendix B. Again, we aim to verify whether the summaries generated with Diff-Sampling are both diverse and of high quality. For diversity, we consider the same metrics presented in Section 5.2, i.e., EAD and SBERT for both crossinput and against-greedy diversity. For quality assessment, we use ROUGE-1 (R-1) (Lin, 2004), a standard metric for summarization that evaluates the ratio of 1-grams present in both the target and generated summaries, as well as the sentence embedding cosine similarity (SIM) between the generated and target summaries. In this way, we compute both syntactic and semantic quality metrics, as a good summary might use entirely different words while still preserving the original meaning. In addition, following Su et al. (2022), we compute the coherence (COH) between the generated output and the text to summarize through the cosine similarity between their SimCSE embeddings (Gao et al., 2021). Finally, for a more qualitative analysis, Appendix E.2 reports some raw outputs. 450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

5.3.2 Experimental Results

As far as the instructed model is considered, all methods achieve the same ROUGE-1 and similarity performances, with very small differences in coherence. Confirming the well-known quality-diversity trade-off (Ippolito et al., 2019), those performing better on coherence are also the worst methods (by a small margin) in terms of diversity.

On the other hand, quality metrics show more variations for the non-instructed model: DiffSampling-cut outperforms all other sampling methods and performs on par with the greedy strategy, while increasing the cross-input EAD score. Contrastive search and *DiffSampling-minp* are immediately below them; however, the latter has slightly higher diversity scores. In general, the quality-diversity trade-off is more pronounced. While *DiffSampling-lb* outperforms top-p, locally typical, and η -sampling in terms of quality while having similar diversity, $\min -p$ sampling seems the only method with the highest diversity without consistent loss in accuracy (but similar results can be achieved by *DiffSampling-minp* with smaller p_{min} ; see Appendix D.3).

5.4 Divergent Association Task

5.4.1 Experimental Setup

The last use case considers the divergent association task (DAT) (Chen and Ding, 2023). Building on the theory that creativity is related to the capability of generating more divergent ideas (Beaty et al.,

Model:			RL	HF-instruc	ted						Pre-trained			
Method		Quality		Cross	Input	Against	-Greedy		Quality		Cross	Input	Against-	Greedy
	R-1	SIM	COH	EAD	SBERT	EAD	SBERT	R-1	SIM	COH	EAD	SBERT	EAD	SBERT
Greedy	$0.22_{\pm.00}$	$0.49_{\pm.01}$	$0.72_{\pm.00}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	-	-	$0.19_{\pm.00}$	$0.45_{\pm.01}$	$0.66_{\pm.01}$	$1.11_{\pm.00}$	$0.94_{\pm.00}$	-	-
Contrastive search	$0.22_{\pm.00}$	$0.50_{\pm .01}$	$0.72_{\pm.00}$	$1.18 \pm .00$	$0.94_{\pm.00}$	$0.21_{\pm.01}$	$0.08 \pm .01$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.64_{\pm.01}$	$1.14_{\pm.00}$	$0.94_{\pm.00}$	$0.45_{\pm.01}$	$0.29_{\pm.01}$
Top- p sampling	$0.22_{\pm.00}$	$0.50_{\pm .01}$	$0.71_{\pm .00}$	$1.21_{\pm.00}$	$0.94_{\pm.00}$	$0.30_{\pm.01}$	$0.12_{\pm.01}$	$0.16 \pm .00$	$0.36_{\pm.01}$	$0.50 {\scriptstyle \pm .01}$	$1.16_{\pm.00}$	$0.93_{\pm.00}$	$0.75_{\pm.01}$	$0.55_{\pm.01}$
η -sampling	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.22_{\pm.00}$	$0.94_{\pm.00}$	$0.33_{\pm.01}$	$0.13_{\pm.01}$	$0.15_{\pm.00}$	$0.35_{\pm.01}$	$0.49_{\pm.01}$	$1.19_{\pm.01}$	$0.93_{\pm.00}$	$0.78_{\pm.01}$	$0.57_{\pm.01}$
Locally typical	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.21_{\pm.00}$	$0.94_{\pm.00}$	$0.30_{\pm.01}$	$0.12_{\pm.01}$	$0.16_{\pm.00}$	$0.35_{\pm.01}$	$0.50_{\pm.01}$	$1.16_{\pm.00}$	$0.93_{\pm.00}$	$0.75_{\pm.01}$	$0.55_{\pm.01}$
Min-p sampling	$0.22_{\pm.00}$	$0.50_{\pm .01}$	$0.72_{\pm.00}$	$1.20_{\pm.00}$	$0.94_{\pm.00}$	$0.29_{\pm.01}$	$0.11_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.61_{\pm .01}$	$1.16_{\pm.01}$	$0.93_{\pm.00}$	$0.62_{\pm.01}$	$0.40_{\pm.01}$
DiffSampling-cut	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	$0.17_{\pm.01}$	$0.06_{\pm.01}$	$0.19_{\pm.00}$	$0.45_{\pm.01}$	$0.66_{\pm.01}$	$1.13_{\pm.00}$	$0.94_{\pm.00}$	$0.25_{\pm.01}$	$0.15_{\pm.01}$
DiffSampling-lb	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.20_{\pm.00}$	$0.94_{\pm.00}$	$0.27_{\pm.01}$	$0.10_{\pm.01}$	$0.17_{\pm.00}$	$0.38_{\pm.01}$	$0.53_{\pm.01}$	$1.15_{\pm.00}$	$0.93_{\pm.00}$	$0.71_{\pm.01}$	$0.50_{\pm.01}$
DiffSampling-minp	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.18_{\pm.00}$	$0.94_{\pm.00}$	$0.23_{\pm.01}$	$0.09_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.63_{\pm.01}$	$1.15_{\pm.01}$	$0.94_{\pm.00}$	$0.49_{\pm.01}$	$0.31_{\pm.01}$

Table 3: Aggregate results over 3 seeds for the XSum dataset for the instructed model (left) and the pre-trained model (right). The mean and standard error of the final score for each run are reported for cross-input diversity, whereas the mean and the 95% confidence interval for the full set of answers are reported for the other metrics.

517

518

519

520

522

523

2014), it requires participants to name unrelated words. Then, their semantic distance can represent an objective measure of divergent thinking (Olson et al., 2021). DAT is a useful case study for decoding strategies as it constrains the generation to different nouns (thus, assuming an optimal probability distribution, the tail due to smoothing should contain everything else) and requires generating terms that are as different as possible, which is quite the opposite to what typically happens in language modeling: an optimal strategy should exclude nonappropriate tokens but also not limit too much the space of possible tokens. More concretely, given the embeddings of n words, the DAT score is the average cosine distance between each pair of embeddings (then scaled as a percentage). We use GloVe embeddings (Pennington et al., 2014) and ask the model to generate a list of 10 nouns. We discard outputs without at least 7 distinct nouns and compute the DAT score for all other outputs over their first 7 nouns. We repeat the experiment 100 times for non-greedy strategies to mitigate the sampling stochasticity.

5.4.2 Experimental Results

Figure 2 summarizes the DAT results for the instructed version of Llama3-8B. *DiffSampling-cut* has the highest average score (even if lower than the greedy score), and generates only valid outputs. Contrastive search, as the baseline closer to greed-iness, is the second-best method in terms of both DAT score and count of valid outputs, while both *DiffSampling-lb* and *DiffSampling-minp* perform almost identically to the top-*p* and min-*p*.

As shown in Figure 3, the results for the pretrained version of Llama3-8B are quite different. *DiffSampling-cut* is still arguably better than contrastive search, as it produces fewer low-scoring and only valid outputs. *DiffSampling-minp* has a slightly lower average score than min-*p* sampling



Figure 2: DAT scores for our methods and the baselines over the instructed version of Llama3-8B. Below, the number of valid outputs produced by each sampling strategy. The dashed line reports the greedy score.

but a few more valid outputs, while *DiffSampling-lb* produces fewer very-high scoring outputs than the other baselines. However, by adjusting p_{lb} and p_{min} , our two relaxation methods can perform at least as well as top-p and min-p (see Appendix D).

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

5.5 Temperature Scaling

Finally, we experiment with different temperature values, i.e., 0.6, 1.0, 1.5, 2.0, and 10.0. As detailed in Section 3, to preserve the mathematical guarantees of our approach, we apply temperature *after* the *DiffSampling* truncation, while our baselines apply this before (see Appendix D.1 for a comparison between temperature before and after truncation). As shown by Figure 4, *DiffSampling+temperature* preserves the output quality, and the only relevant differences occur with our two relaxations and only pre-trained models. Instead, the output quality rapidly drops with higher temperatures for the min-p (by far the best of our baselines at tem-



Figure 3: DAT scores for our methods and the baselines over the non-instructed version of Llama3-8B. Below, the number of valid outputs produced by each sampling strategy. The dashed line reports the greedy score.

peratures greater than 1.0) and top-*p* baselines. In particular, the non-significant loss in quality for *DiffSampling-cut* confirms that our truncation strategy only preserves correct tokens. At the same time, temperature scaling has an (overall positive) impact on diversity; we refer to Appendix C for a detailed analysis of how all our quality and diversity metrics change at different temperatures.

6 Discussion

543

544

546

550

552

554

555

556

557

558

560

561

562

563

569

571

573

Overall, *DiffSampling-cut* has demonstrated performance better than or equal to the greedy strategy. Additionally, it offers the potential for greater diversity. By introducing a lower bound on the preserved total probability or an upper bound on the number of truncated tokens, the method can further relax selection constraints, enabling greater output diversity at the expense of a marginal reduction in prediction accuracy. Once truncation is applied, sampling at higher temperatures becomes viable, promoting greater variability without significantly compromising output quality.

However, selecting the most appropriate method and hyperparameters is not straightforward and requires a case-by-case analysis on whether it is better to have higher quality or diversity. *DiffSampling-cut* works best when the task requires precision. Instead, *DiffSampling-lb* fosters output diversity by trading off some accuracy, especially at higher values of p_{lb} and, thus, appears most appropriate for divergent solutions. *DiffSampling-minp* is more well-balanced. Increasing the temperature



Figure 4: Quality measures across different temperature values for top-p, min-p, and our methods. For GSM8K and MATH, we report the percentage of solved problems. For XSum, we report the coherence between the output and the text to summarize. For DAT, we report the number of valid generated outputs.

has proven highly effective for fine-tuned models across all methods, whenever it is not strictly necessary to preserve the originally learned distribution. 574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

597

7 Conclusion

In this paper, we have presented *DiffSampling*, a novel family of decoding strategies based on the analysis of the next-token distribution. In particular, given the distribution sorted in descending order, we have proposed to compute the forward difference approximation of its discrete derivative and use it to remove tokens after its minimum value (possibly together with relaxations to allow for more diversity). In this way, *DiffSampling* avoids incorrect tokens under the learned distribution. We have experimented with three different tasks, finding that our methods consistently perform at least as well as similar common strategies in terms of the accuracy of results and diversity of outputs.

Our research agenda includes exploring whether combining *DiffSampling* with other techniques can yield even better results, including for longer-form generation. We also aim to leverage other properties of the distribution to guide text generation toward desired characteristics.

8 Limitations

598

610

612

613

615

616

617

618

621

622

624

627

636

637

640

644

645

The work presented in this paper has a few important limitations to highlight. Firstly, *DiffSampling* is merely a decoding strategy. While it can influence the accuracy and diversity of the model's outputs, it is constrained by the information learned by the model itself. For instance, if the model is biased toward certain grammatical structures, the probability mass is likely to contain only tokens that adhere to those structures. In addition, working at the decoding level means that the information stored by the model is not modified at all. While *DiffSampling* can potentially reduce how much a model *regurgitates* pre-existing text, it cannot reduce how much a model *memorizes* it.

Moreover, *DiffSampling* is governed by two parameters: the nucleus lower bound and the truncated probability upper bound. Each of the three methods has its advantages and disadvantages concerning the exploitation and exploration of the next-token distribution. While this can guide the choice between them, there is no golden rule; users must select the most appropriate strategy on a case-by-case basis. Similarly, we did not find specific parameter values to be universally superior, and different scenarios may require users to adjust them accordingly.

Additionally, our experiments encompassed only three case studies with relatively short generated outputs. While we chose these to maximize their diversity, it is difficult to estimate the actual advantage of using *DiffSampling* for other tasks (especially longer-form generation ones) and with different LLMs. We intend to broaden our investigation in the future, for example, by incorporating models of varying sizes. At the same time, we believe that the choice of LLM per se should not change the ranking of the decoding techniques in terms of performance, given the fact that our method is based on the analysis of the token probability distribution in output from these models.

Finally, our evaluation makes use of quantitative, automatic metrics for both quality and diversity.
However, several of these metrics exhibit significant limitations (e.g., Schluter (2017)), often failing to align with human judgments (Tevet and Berant, 2021). Moreover, abstract concepts such as originality and creativity remain inherently difficult to define with precision (Franceschelli and Musolesi, 2024). We plan to experiment with human evaluators to verify whether the quality and diversity that

DiffSampling aims to provide are also perceived by potential users.

649

650

651

652

653

654

655

656

657

658

659

660

661

662

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

References

- Sourya Basu, Govardana Sachitanandam Ramachandran, Nitish Shirish Keskar, and Lav R. Varshney. 2021. Mirostat: A neural text decoding algorithm that directly controls perplexity. In *Proc. of the 9th International Conference on Learning Representations (ICLR'21).*
- Roger E. Beaty, Paul J. Silvia, Emily C. Nusbaum, Emanuel Jauk, and Mathias Benedek. 2014. The roles of associative and executive processes in creative cognition. *Memory & Cognition*, 42(7):1186– 1197.
- Loubna Ben Allal, Anton Lozhkov, and Elie Bakouch. 2024. SmolLM - blazingly fast and remarkably powerful. https://huggingface.co/blog/ smollm [Accessed July 25, 2024].
- Rishi Bommasani, Drew Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney Arx, Michael Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri Chatterji, Annie Chen, Kathleen Creel, Jared Davis, Dora Demszky, ..., and Percy Liang. 2021. On the opportunities and risks of foundation models. ArXiv:2108.07258 [cs.LG].
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with GPT-4. ArXiv:2303.12712 [cs.CL].
- Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. 2021. Extracting training data from large language models. In *Proc. of the 30th USENIX Security Symposium* (USENIX Security 21).
- Honghua Chen and Nai Ding. 2023. Probing the "creativity" of large language models: Can models produce divergent semantic association? In *Findings* of the Association for Computational Linguistics (EMNLP'23).
- Anshuman Chhabra, Hadi Askari, and Prasant Mohapatra. 2024. Revisiting zero-shot abstractive summarization in the era of large language models from the perspective of position bias. In *Proc. of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL'24).*
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro

703

704

- 731 732 733 734
- 735 736 737 738
- 740 741 742
- 745
- 747
- 749 750

751 752

755

- Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. ArXiv:2110.14168 [cs.LG].
- Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. In Proc. of the 56th Annual Meeting of the Association for Computational Linguistics (ACL'18).
- Eve Fleisig, Genevieve Smith, Madeline Bossi, Ishita Rustagi, Xavier Yin, and Dan Klein. 2024. Linguistic bias in ChatGPT: Language models reinforce dialect discrimination. In Proc. of the 2024 Conference on Empirical Methods in Natural Language Processing (EMNLP'24).
- Giorgio Franceschelli, Claudia Cevenini, and Mirco Musolesi. 2024. Training foundation models as data compression: On information, model weights and copyright law. ArXiv:2407.13493 [cs.CY].
- Giorgio Franceschelli and Mirco Musolesi. 2024. On the Creativity of Large Language Models. AI&Society, pages 1–11.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple contrastive learning of sentence embeddings. In Proc. of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP'21).
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, ..., and Zhiyu Ma. 2024. The llama 3 herd of models. ArXiv:2407.21783 [cs.AI].
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the MATH dataset. In Proc. of the 35th Conference on Neural Information Processing Systems Datasets and Benchmarks Track.
- John Hewitt, Christopher Manning, and Percy Liang. 2022. Truncation sampling as language model desmoothing. In Findings of the Association for Computational Linguistics (EMNLP'22).
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. In Proc. of the 8th International Conference on Learning Representations (ICLR'20).
- Zhang-Wei Hong, Idan Shenfeld, Tsun-Hsuan Wang, Yung-Sung Chuang, Aldo Pareja, James R. Glass, Akash Srivastava, and Pulkit Agrawal. 2024. Curiosity-driven red-teaming for large language models. In Proc. of the 12th International Conference on Learning Representations (ICLR'24).

Daphne Ippolito, Reno Kriz, João Sedoc, Maria Kustikova, and Chris Callison-Burch. 2019. Comparison of diverse decoding methods from conditional language models. In Proc. of the 57th Annual Meeting of the Association for Computational Linguistics (ACL'19).

756

760

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

782

783

784

785

786

787

788

790

791

792

794

795

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

- Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward Grefenstette, and Roberta Raileanu. 2024. Understanding the effects of RLHF on LLM generalisation and diversity. In Proc. of the 12th International Conference on Learning Representations (ICLR'24).
- Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. 2023. Contrastive decoding: Open-ended text generation as optimization. In Proc. of the 61st Annual Meeting of the Association for Computational Linguistics (ACL'23).
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In Proc. of ACL'04 Workshop on Text Summarization Branches Out.
- Siyang Liu, Sahand Sabour, Yinhe Zheng, Pei Ke, Xiaoyan Zhu, and Minlie Huang. 2022. Rethinking and refining the distinct metric. In Proc. of the 60th Annual Meeting of the Association for Computational Linguistics (ACL'22).
- Clara Meister, Tiago Pimentel, Gian Wiher, and Ryan Cotterell. 2023. Locally typical sampling. Transactions of the Association for Computational Linguistics, 11:102–121.
- Nguyen Nhat Minh, Andrew Baker, Clement Neo, Allen G Roush, Andreas Kirsch, and Ravid Shwartz-Ziv. 2025. Turning up the heat: Min-p sampling for creative and coherent LLM outputs. In Proc. of the 13th International Conference on Learning Representations (ICLR'25).
- Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. In Proc. of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP'18).
- Jay A. Olson, Johnny Nahas, Denis Chmoulevitch, Simon J. Cropper, and Margaret E. Webb. 2021. Naming unrelated words predicts creativity. Proceedings of the National Academy of Sciences, 118(25):e2022340118.
- Max Peeperkorn, Tom Kouwenhoven, Dan Brown, and Anna Jordanous. 2024. Is temperature the creativity parameter of large language models? In Proc. of the 15th International Conference on Computational Creativity (ICCC'24).
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. GloVe: Global vectors for word representation. In Proc. of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP'14).

Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. In Proc. of the 37th Conference on Neural Information Processing Systems (NIPS'23).

813

814

815

817

819

823

824 825

826

827 828

829 830

832

837

839

841 842

847

848

857

858

859

- Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERTnetworks. In Proc. of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP'19).
- Pedro Reviriego, Javier Conde, Elena Merino-Gómez, Gonzalo Martínez, and José Alberto Hernández.
 2023. Playing with words: Comparing the vocabulary and lexical richness of chatgpt and humans. ArXiv:2308.07462 [cs.CL].
 - Natalie Schluter. 2017. The limits of automatic summarisation according to ROUGE. In *Proc. of the 15th Conference of the European Chapter of the Association for Computational Linguistics (EACL'17).*
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. In Advances in Neural Information Processing Systems (NIPS'20).
- Yixuan Su, Tian Lan, Yan Wang, Dani Yogatama, Lingpeng Kong, and Nigel Collier. 2022. A contrastive framework for neural text generation. In *Advances in Neural Information Processing Systems (NIPS'22).*
- Guy Tevet and Jonathan Berant. 2021. Evaluating the evaluation of diversity in natural language generation. In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics (EACL'21).
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, ..., and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. ArXiv:2307.09288 [cs.CL].
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng YU, Zhengying Liu, Yu Zhang, James Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2024. Metamath: Bootstrap your own mathematical questions for large language models. In Proc. of the 12th International Conference on Learning Representations (ICLR'24).

363

864

870

871

A Computational Infrastructure

877 878 879

876

882

884

All experiments were carried out on a local server equipped with 2 NVIDIA L40 GPUs and 1 NVIDIA H100 GPU.

867 **B** Prompts

As reported in Section 5, we tested *DiffSampling* on three case studies. For the mathematical problemsolving tasks, we adopted the same prompt from Yu et al. (2024), i.e.:

Below is an instruction that describes a task. Write a response that appropriately completes the request.

Instruction:
{question}

Response: Let's think step by step.

For the extreme summarization task, the prompt adopted for the instructed version of Llama2-7B is the same as in Chhabra et al. (2024):

[INST]	For	the	following	article:
{article	}			

Return a summary comprising of 1 sentence. With the sentence in a numbered list format.

For example:

1. First sentence [/INST]

where [INST] and [/INST] are special tokens used by Llama2-7b to identify different roles in the chat. Vice versa, for the non-instructed version, we

used:

Generate a 1 sentence summary for the given article. Article: {article} Summary:

Finally, for the divergent association task, we considered the following prompt for the instructed version of Llama3-8B:

user

Please write 10 nouns in English that are as irrelevant from each other as possible, in all meanings and uses of the words. Please note that the words you write should have only single word, only nouns (e.g., things, objects, concepts), and no proper nouns (e.g., no specific people or places). assistant

Here are the 10 nouns in English that are as irrelevant from each other as possible:

where user and assistant are keywords used by Llama3-8b to identify different roles in the chat, while for its non-instructed version we used the following: 885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

Write 10 nouns in English that are as irrelevant from each other as possible, in all meanings and uses of the words. Please note that the words you write should have only single word, only nouns (e.g., things, objects, concepts), and no proper nouns (e.g., no specific people or places).

Solution:

Here are the 10 nouns in English that are as irrelevant from each other as possible:

C Experiments on Temperature Scaling

In addition to investigating performances at a temperature $\tau = 1.0$, we also conduct experiments with a lower temperature value (0.6) and three higher temperature values (1.5, 2.0, and 10.0), to verify whether our truncation strategy only preserves appropriate tokens, i.e., whether at different temperatures the quality of generated outputs remains competitive and the diversity increases.

C.1 Math Problem Solving

Table 4 reports all the results with different temperatures for the GSM8K (left side) and MATH (right side) test sets. For the former, a lower temperature makes all the models (including the baselines) more in line with greedy strategy, thus diminishing the diversity scores while usually increasing the accuracy. On the contrary, all the baselines tend to perform poorer at increasing temperatures in terms

Method Accuracy Cross-Input Against-Greedy EAD SBERT EAD <t< th=""><th>Dataset:</th><th></th><th></th><th>GSM8K</th><th></th><th></th><th></th><th></th><th>MATH</th><th></th><th></th></t<>	Dataset:			GSM8K					MATH		
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Method	Accuracy	Cross	Input	Against	-Greedy	Accuracy	Cross	Input	Against	-Greedy
$ \begin{array}{ c c c c c c c c c c c c c$			EAD	SBERT	EAD	SBERT		EAD	SBERT	EAD	SBERT
	Temperature = 0.0										
$ \begin{array}{c} \hline Temperature = 1.0 \\ \hline Temperature = 2.0 \\ \hline Temperature $	Greedy	$66.44_{\pm.09}$	$2.03_{\pm.00}$	$0.64_{\pm.00}$	-	-	$20.62_{\pm.01}$	$5.65_{\pm.00}$	$0.80_{\pm.00}$	-	-
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	Temperature = 1.0										
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Contrastive search	$65.88_{\pm.59}$	$2.06_{\pm.00}$	$0.64_{\pm.00}$	$0.17_{\pm.00}$	$0.02_{\pm.00}$	$21.05 \pm .14$	$5.82_{\pm.01}$	$0.80_{\pm .00}$	$0.31_{\pm.00}$	$0.09_{\pm.00}$
$ \begin{array}{llllllllllllllllllllllllllllllllllll$	Top- p sampling	$65.00 \pm .18$	$2.08_{\pm.01}$	$0.64_{\pm.00}$	$0.23_{\pm.00}$	$0.03_{\pm.00}$	$20.02_{\pm.12}$	$6.08_{\pm.02}$	$0.80_{\pm .00}$	$0.36_{\pm.00}$	$0.10_{\pm.00}$
	η -sampling	$65.05_{\pm.19}$	$2.12_{\pm.00}$	$0.64_{\pm.00}$	$0.25_{\pm.00}$	$0.04_{\pm.00}$	$19.67_{\pm.20}$	$6.36_{\pm.01}$	$0.80_{\pm .00}$	$0.39_{\pm.00}$	$0.11_{\pm.00}$
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Locally typical	$66.29_{\pm.55}$	$2.09_{\pm.00}$	$0.64_{\pm.00}$	$0.23_{\pm.00}$	$0.03_{\pm.00}$	$19.95_{\pm.26}$	$6.06_{\pm.01}$	$0.80_{\pm .00}$	$0.36_{\pm.00}$	$0.10_{\pm.00}$
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Min-p sampling	$65.76_{\pm.44}$	$2.09_{\pm.00}$	$0.64_{\pm.00}$	$0.23_{\pm.00}$	$0.03_{\pm.00}$	$20.25 \pm .09$	$6.09_{\pm.01}$	$0.80_{\pm .00}$	$0.36_{\pm.00}$	$0.10_{\pm.00}$
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	DiffSampling-cut	$66.36_{\pm.23}$	$2.04_{\pm.00}$	$0.64_{\pm.00}$	$0.14_{\pm.00}$	$0.02_{\pm.00}$	$21.38_{\pm.20}$	$5.71_{\pm.01}$	$0.80_{\pm .00}$	$0.27_{\pm.00}$	$0.07_{\pm.00}$
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	DiffSampling-lb	$66.92_{\pm.08}$	$2.07_{\pm.00}$	$0.64_{\pm.00}$	$0.20_{\pm.00}$	$0.03_{\pm.00}$	$20.78_{\pm.14}$	$6.00_{\pm.01}$	$0.80_{\pm.00}$	$0.35_{\pm.00}$	$0.10_{\pm.00}$
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	DiffSampling-minp	$66.44_{\pm.35}$	$2.05_{\pm.00}$	$0.64_{\pm.00}$	$0.19_{\pm.00}$	$0.03_{\pm.00}$	$21.13_{\pm.08}$	$5.87_{\pm.01}$	$0.80_{\pm,00}$	$0.33_{\pm,00}$	$0.09_{\pm.00}$
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Temperature = 0.6										
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Top- <i>p</i> sampling	$66.34_{\pm.67}$	$2.05_{\pm.01}$	$0.64_{\pm.00}$	$0.17_{\pm.00}$	$0.02_{\pm.00}$	$21.58_{\pm,32}$	$5.81_{\pm.02}$	$0.80_{\pm.00}$	$0.31_{\pm,00}$	$0.09_{\pm.00}$
	η -sampling	$66.26_{\pm,22}$	$2.07_{\pm.01}$	$0.64_{\pm.00}$	$0.19_{\pm.00}$	$0.03_{\pm.00}$	$20.36_{\pm.15}$	$5.94_{\pm.01}$	$0.80_{\pm,00}$	$0.33_{\pm,00}$	$0.09_{\pm.00}$
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Locally typical	$66.34_{\pm.67}$	$2.05_{\pm.01}$	$0.64_{\pm.00}$	$0.17_{\pm.00}$	$0.02_{\pm.00}$	$21.58_{\pm,32}$	$5.81_{\pm.02}$	$0.80_{\pm.00}$	$0.31_{\pm.00}$	$0.09_{\pm.00}$
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Min-p sampling	$66.52_{\pm,30}$	$2.06_{\pm.01}$	$0.64_{\pm.00}$	$0.17_{\pm.00}$	$0.02_{\pm.00}$	$21.31_{\pm.08}$	$5.81_{\pm.01}$	$0.80_{\pm.00}$	$0.31_{\pm.00}$	$0.09_{\pm.00}$
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	DiffSampling-cut	$66.74 \pm .04$	$2.05_{\pm.00}$	$0.64_{\pm.00}$	$0.13_{\pm.00}$	$0.02_{\pm.00}$	$21.52_{\pm,13}$	$5.72_{\pm.00}$	$0.80_{\pm.00}$	$0.25_{\pm.00}$	$0.07_{\pm.00}$
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	DiffSampling-lb	$66.49_{\pm 41}$	$2.06_{\pm 01}$	$0.64_{\pm 00}$	0.18 ± 0.0	$0.03_{\pm 0.00}$	21.09 ± 11	$5.83_{\pm 01}$	0.80 ± 0.00	$0.32_{\pm 00}$	0.09 ± 0.00
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	DiffSampling-minp	$66.84_{\pm 73}$	$2.05_{\pm 01}$	$0.64_{\pm 00}$	$0.16_{\pm 0.00}$	$0.02_{\pm 0.00}$	$20.79_{\pm 10}$	$5.78_{\pm 01}$	$0.80_{\pm 0.00}$	$0.30_{\pm 0.00}$	$0.08_{\pm 0.00}$
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Temperature = 1.5	1 1.10	1.01	2.00	1.00	1.00	1.110	1.01	1.00	2.00	1.00
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Top- <i>p</i> sampling	$63.91_{\pm 57}$	$2.17_{\pm 01}$	$0.64_{\pm 00}$	0.28 ± 0.0	$0.04_{\pm 0.00}$	$18.38_{\pm 22}$	6.92 ± 0.02	0.80 ± 0.0	$0.42_{\pm 00}$	$0.12_{\pm 0.0}$
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	η -sampling	$60.35 \pm .55$	$2.28_{\pm.00}$	$0.64_{\pm.00}$	$0.32_{\pm.00}$	$0.05_{\pm.00}$	$15.63_{\pm,17}$	$7.77_{\pm.01}$	$0.80_{\pm.00}$	$0.45_{\pm.00}$	$0.14_{\pm.00}$
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Locally typical	$64.39_{\pm 41}$	$2.17_{\pm 01}$	$0.64_{\pm 00}$	0.28 ± 0.00	$0.04_{\pm 0.00}$	18.73 ± 01	$7.04_{\pm 02}$	0.80 ± 0.00	$0.42_{\pm 00}$	$0.12_{\pm 0.00}$
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Min- <i>p</i> sampling	$64.29_{\pm 38}$	$2.15_{\pm 00}$	$0.64_{\pm 00}$	$0.28_{\pm 00}$	$0.04_{\pm 0.00}$	$18.94_{\pm 23}$	$6.54_{\pm 0.02}$	0.80 ± 0.00	$0.40_{\pm 0.00}$	$0.12_{\pm 0.00}$
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	DiffSampling-cut	$66.72_{\pm 36}$	$2.05_{\pm 0.00}$	$0.64_{\pm 00}$	$0.15_{\pm 0.0}$	$0.02_{\pm 0.00}$	$21.36_{\pm 15}$	$5.73_{\pm 00}$	0.80 ± 0.00	$0.27_{\pm 00}$	$0.07_{\pm 0.00}$
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	DiffSampling-lb	66.84 ± 43	2.08 ± 00	0.64 ± 0.00	0.22 ± 00	0.03 ± 0.00	$20.52_{\pm 34}$	6.03 ± 01	0.80 ± 0.00	0.36 ± 0.0	0.10 ± 0.00
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	DiffSampling-minp	66.24 ± 36	2.08 ± 0.00	0.64 ± 0.00	0.20 ± 00	0.03 ± 0.00	20.79 ± 14	5.88 ± 01	0.80 ± 0.00	0.34 ± 0.0	0.10 ± 0.00
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Temperature = 2.0	1.00	1.00	1 1 1.00	1.00		1.14	1.01	1 1 1 1.00	1.00	1.00
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Top- <i>p</i> sampling	25.40 ± 0.7	10.13+ 10	0.66 ± 0.00	$0.70_{\pm 01}$	$0.36_{\pm 0.01}$	$2.49_{\pm 0.1}$	48.71+ 08	$0.52_{\pm 0.0}$	$0.92_{\pm 00}$	0.68 ± 0.0
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	η -sampling	$35.51_{\pm 30}$	$7.35_{\pm 0.5}$	$0.69_{\pm 00}$	$0.58_{\pm 01}$	$0.22_{\pm 01}$	$4.26_{\pm 06}$	$43.39_{\pm 10}$	$0.64_{\pm 00}$	$0.86_{\pm 00}$	$0.53_{\pm 0.0}$
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Locally typical	24.61 ± 60	10.65 ± 05	0.65 ± 0.00	0.71 ± 01	0.37 ± 0.01	2.46 ± 0.03	51.04 ± 0.07	0.50 ± 0.00	0.93 ± 0.0	0.69 ± 0.00
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	Min- <i>p</i> sampling	$62.19_{\pm 37}$	2.24 ± 0.00	0.64 ± 0.00	0.32 ± 0.01	0.05 ± 0.01	16.92 ± 21	7.21 ± 0.01	0.80 ± 0.00	0.44 ± 0.00	0.13 ± 0.0
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	DiffSampling-cut	66.44 ± 18	2.05 ± 00	0.64 ± 0.00	0.15 ± 0.0	0.02 ± 0.00	21.66 ± 20	5.71 ± 01	0.80 ± 0.00	0.27 ± 0.0	0.08 ± 0.00
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	DiffSampling-lb	65.73 ± 47	2.08 ± 0.01	0.64 ± 0.00	0.23 ± 00	0.03 ± 0.00	20.34 ± 0.08	6.09 ± 02	$0.80 \pm .00$	0.37 ± 0.0	0.11 ± 0.00
$\frac{T_{\text{respersive}}}{T_{\text{respersive}}} = 10.0$ $\frac{T_{\text{respersive}}}{T_{\text{resp}}} = 10.0$ $\frac{T_{\text{respersive}}}{T_{\text{resp}}} = 10.0$ $\frac{T_{\text{resp}}}{T_{\text{resp}}} = 10.0$	DiffSampling-minp	66.21 ± 27	2.06 ± 00	0.64 ± 0.00	0.21 ± 00	0.03 ± 0.00	20.59 ± 13	5.90 ± 0.02	$0.80 \pm .00$	0.34 ± 0.0	0.10 ± 0.00
Top-p sampling $0.00_{\pm.00}$ $17.26_{\pm.03}$ $0.12_{\pm.00}$ $1.00_{\pm.00}$ $0.96_{\pm.00}$ $58.65_{\pm.03}$ $0.12_{\pm.00}$ $1.00_{+.00}$	Temperature = 10.0	001-11.21		010-1.00	0.1111.00	01001.00	20100 1.15	01001.02	0.007.00	010-1.00	01201.00
$1 = 0.0 \pm 0.00 \pm 0.00$	Top- <i>p</i> sampling	0.00 ± 0.00	17.26 ± 03	0.12 ± 0.0	1.00 ± 0.0	0.96 ± 0.0	0.00 ± 0.00	58.65 ± 03	0.12 ± 0.0	1.00 ± 0.0	1.00 ± 0.0
n sampling 0.00+ ∞ 17.43+ α 0.12+ ∞ 1.00+ ∞ 0.96+ ∞ 0.00+ ∞ 59.18+ α 0.12+ ∞ 1.00+ ∞ 1.00+ ∞	<i>n</i> -sampling	0.00 ± 00	17.43 ± 04	$0.12 \pm .00$	1.00 ± 00	$0.96 \pm .00$	0.00 ± 00	59.18 ± 02	$0.12 \pm .00$	1.00 ± 00	$1.00 \pm .00$
$\begin{array}{c} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 $	Locally typical	0.00 ± 00	17.52 ± 01	$0.11 \pm .00$	$1.00 \pm .00$ $1.01 \pm .00$	$0.96 \pm .00$	0.00 ± 00	59.69 ± 01	$0.11 \pm .00$	1.01 ± 00	$1.00 \pm .00$
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	Min-n sampling	0.00 ± 00	17.39 ± 04	0.13 ± 0.0	1.00 ± 00	0.95 ± 00	0.00 ± 00	59.16 ± 00	0.13 ± 0.0	1.00+00	1.00 ± 00
DiffSampling $1.00 \pm .00$ $0.00 \pm .00$ $0.00 \pm .00$ $0.00 \pm .00$ $1.00 \pm .00$	DiffSampling-cut	66.31	2.04 ± 0.04	$0.64 \pm \infty$	$0.15 \pm \infty$	0.02 ± 00	21.22 ± 11	5.74 ± 01	$0.80 \pm \infty$	0.28 ± 00	$0.08 \pm \infty$
DiffSampling-lb $66.26_{\pm,72}$ $2.10_{\pm,00}$ $0.01_{\pm,00}$ $0.10_{\pm,00}$ $0.02_{\pm,00}$ $0.02_{\pm,00}$ $0.02_{\pm,00}$ $0.02_{\pm,00}$ $0.00_{\pm,00}$	DiffSampling-lb	66.26 ± 70	2.10 ± 0.00	0.64 ± 0.00	0.25 ± 00	0.04 ± 0.00	19.88 ± 02	6.20 ± 01	0.80 ± 00	0.39 ± 00	0.11 ± 00
DiffSampling-minp $\begin{bmatrix} 65.43_{\pm,12} \\ 2.07_{\pm,00} \\ 0.04_{\pm,00} \\ 0.04_{\pm,00} \\ 0.02_{\pm,00} \\ 0.03_{\pm,00} \\ 0.03_{\pm,00} \\ 0.10_{\pm,00} \\ 0.10_{\pm,01} \\ 0.10_{\pm$	DiffSampling-minp	$65.43_{\pm 21}$	2.07 ± 00	0.64 ± 0.0	0.22 ± 00	$0.03_{\pm 0.00}$	$21.17_{\pm 14}$	5.94 ± 00	0.80 ± 0.00	0.35 ± 0.00	0.10 ± 0.00

Table 4: Accuracy and diversity of results for the GSM8K and MATH test sets over 3 seeds with different temperature values. The mean and standard error of the final score for each run are reported for accuracy and cross-input diversity, whereas the mean and 95% confidence interval for the full set of answers are reported for against-greedy diversity.

of output correctness, while diversity improves accordingly (especially for a syntactic-based metric such as EAD; the qualitative examples reported in Appendix E.1 demonstrate why). Instead, our methods maintain the highest possible accuracy, with a slight improvement in diversity at higher τ .

909

910

911

912

913

914

915For the latter, a lower temperature makes all the916baselines closer to our methods in terms of accu-917racy, while diminishing their diversity scores. At918increasing temperature, the baselines rapidly start919failing to solve the problems, possibly due to a920more random selection of tokens that also causes

syntactic diversity to increase. By applying temperature after the truncation, our methods preserve their output quality regardless of the temperature used, with small but relevant gains in diversity (for example, *DiffSampling-minp* at even $\tau = 10.0$ has an accuracy comparable with min-*p* and $\tau = 0.6$, but with clearly higher diversity scores).

For the sake of completeness, we also report the full results on a sample of 1000 entries from the MetaMathQA training set. As apparent from Table 5, the greediness of the approach is directly correlated with the accuracy of solutions. In particular,

931

932

921

Method	Accuracy	Cross-Inpu	t Diversity	Against-Gree	dy Diversity
		EAD	SBERT	EAD	SBERT
Greedy	$95.27_{\pm.17}$	$1.67_{\pm.01}$	$0.71_{\pm.00}$	-	-
Temperature = 1.0					
Contrastive search	$94.17_{\pm.45}$	$1.68_{\pm.01}$	$0.71_{\pm.00}$	$0.14_{\pm.00}$	$0.03_{\pm.00}$
Top- p sampling	$91.70_{\pm.62}$	$1.71_{\pm.01}$	$0.71_{\pm.00}$	$0.19_{\pm.01}$	$0.04_{\pm.00}$
η -sampling	$89.17_{\pm.52}$	$1.73_{\pm.01}$	$0.71_{\pm.00}$	$0.21_{\pm.01}$	$0.04_{\pm.00}$
Locally typical	$91.70_{\pm.62}$	$1.71_{\pm.01}$	$0.71_{\pm.00}$	$0.19_{\pm.01}$	$0.04_{\pm.00}$
Min-p sampling	$92.00_{\pm.49}$	$1.70_{\pm.01}$	$0.71_{\pm.00}$	$0.19_{\pm.01}$	$0.03_{\pm.00}$
DiffSampling-cut	$94.93_{\pm.45}$	$1.68_{\pm.01}$	$0.71_{\pm.00}$	$0.11_{\pm.00}$	$0.02_{\pm.00}$
DiffSampling-lb	$93.33_{\pm.03}$	$1.68_{\pm.01}$	$0.71_{\pm.00}$	$0.17_{\pm.01}$	$0.03_{\pm.00}$
DiffSampling-minp	$94.13_{\pm.22}$	$1.69_{\pm.01}$	$0.71_{\pm.00}$	$0.16_{\pm.01}$	$0.03_{\pm.00}$
<i>Temperature</i> = 0.6					
Top-p	$94.03_{\pm.27}$	$1.68_{\pm.01}$	$0.71_{\pm.00}$	$0.14_{\pm.01}$	$0.03_{\pm.00}$
η -sampling	$93.90_{\pm .09}$	$1.70_{\pm.01}$	$0.71_{\pm.00}$	$0.15_{\pm.01}$	$0.03_{\pm.00}$
Locally typical	$94.03_{\pm.27}$	$1.68_{\pm.01}$	$0.71_{\pm.00}$	$0.14_{\pm.01}$	$0.03_{\pm.00}$
Min-p	$93.93_{\pm .36}$	$1.68_{\pm.01}$	$0.71_{\pm.00}$	$0.14_{\pm.01}$	$0.03_{\pm.00}$
DiffScut	$94.83_{\pm.18}$	$1.68_{\pm.01}$	$0.71_{\pm.00}$	$0.11_{\pm.00}$	$0.02_{\pm.00}$
DiffSlb	$93.43_{\pm.21}$	$1.69_{\pm.01}$	$0.71_{\pm.00}$	$0.14_{\pm.01}$	$0.03_{\pm.00}$
DiffSminp	$93.97_{\pm.30}$	$1.67_{\pm.01}$	$0.71_{\pm.00}$	$0.13_{\pm.00}$	$0.02_{\pm.00}$
Temperature = 1.5					
Top-p	$87.63_{\pm.90}$	$1.75_{\pm.00}$	$0.71_{\pm.00}$	$0.24_{\pm.01}$	$0.05_{\pm.00}$
η -sampling	$80.63_{\pm.23}$	$1.83_{\pm.01}$	$0.71_{\pm.00}$	$0.28_{\pm.01}$	$0.05_{\pm.00}$
Locally typical	$87.97_{\pm .38}$	$1.77_{\pm.01}$	$0.71_{\pm.00}$	$0.24_{\pm.01}$	$0.05_{\pm.00}$
Min-p	$87.90_{\pm.57}$	$1.75_{\pm.01}$	$0.71_{\pm.00}$	$0.24_{\pm.01}$	$0.05_{\pm.00}$
DiffScut	$95.17_{\pm.18}$	$1.67_{\pm.01}$	$0.71_{\pm.00}$	$0.11_{\pm.00}$	$0.02_{\pm.00}$
DiffSlb	$92.67_{\pm.52}$	$1.70_{\pm.02}$	$0.71_{\pm.00}$	$0.18_{\pm.01}$	$0.03_{\pm.00}$
DiffSminp	$93.83_{\pm.42}$	$1.69 \pm .01$	$0.71_{\pm.00}$	$0.17_{\pm.01}$	$0.03_{\pm.00}$
Temperature = 2.0					
Top-p	$30.17_{\pm.76}$	$8.30_{\pm.11}$	$0.67_{\pm.00}$	$0.70_{\pm.01}$	$0.44_{\pm.01}$
η -sampling	$42.20_{\pm 1.0}$	$6.26_{\pm.13}$	$0.72_{\pm.00}$	$0.58_{\pm.01}$	$0.29_{\pm.01}$
Locally typical	$29.07_{\pm.72}$	$8.68_{\pm.08}$	$0.65_{\pm.00}$	$0.71_{\pm.01}$	$0.45_{\pm.01}$
Min-p	$83.93_{\pm.22}$	$1.80_{\pm.02}$	$0.71_{\pm.00}$	$0.28_{\pm.01}$	$0.05_{\pm.00}$
DiffScut	$94.33_{\pm .35}$	$1.68_{\pm.01}$	$0.71_{\pm.00}$	$0.12_{\pm.00}$	$0.02_{\pm.00}$
DiffSlb	$92.87_{\pm.21}$	$1.71_{\pm.01}$	$0.71_{\pm.00}$	$0.19_{\pm.01}$	$0.04_{\pm.00}$
DiffSminp	$93.37_{\pm.07}$	$1.69_{\pm.01}$	$0.71_{\pm.00}$	$0.17_{\pm.01}$	$0.03_{\pm.00}$
Temperature = 10.0					
Top-p	$0.00_{\pm .00}$	$13.10_{\pm.03}$	$0.12_{\pm.00}$	$1.00_{\pm.00}$	$0.98_{\pm.00}$
η -sampling	$0.00_{\pm .00}$	$13.22_{\pm.03}$	$0.12_{\pm.00}$	$1.00_{\pm.00}$	$0.98_{\pm.00}$
Locally typical	$0.00_{\pm .00}$	$13.33_{\pm.02}$	$0.11_{\pm.00}$	$1.01_{\pm.00}$	$0.98_{\pm.00}$
Min-p	$0.00_{\pm .00}$	$13.21_{\pm.03}$	$0.13_{\pm.00}$	$1.00_{\pm.00}$	$0.98_{\pm.00}$
DiffScut	$94.47_{\pm.19}$	$1.67_{\pm.01}$	$0.71_{\pm.00}$	$0.12_{\pm.00}$	$0.02_{\pm.00}$
DiffSlb	$92.10_{\pm.43}$	$1.70_{\pm.01}$	$0.71_{\pm.00}$	$0.21_{\pm.01}$	$0.04_{\pm.00}$
DiffSminp	$93.13_{\pm.17}$	$1.69 \pm .02$	$0.71_{\pm.00}$	$0.18_{\pm.01}$	$0.03_{\pm.00}$

Table 5: Accuracy and diversity of results for the training set over 3 seeds with different temperature values. The mean and standard error of the final score for each run are reported for accuracy and cross-input diversity, whereas the mean and 95% confidence interval for the full set of answers are reported for against-greedy diversity.

sampling at a temperature of 0.6 increases the accuracy of all baselines while undermining their diversity scores, while higher temperatures lower the accuracy and increase (syntactic) diversity; notably, a very high temperature causes semantic diversity to fall. On the contrary, our three methods achieve similar accuracy at any temperature, with small increases in diversity.

C.2 Extreme Summarization

933

936

937

938

939

940

941

942

943

Similar considerations can be traced for XSum, as reported by Table 6. For both RLHF-instructed and pre-trained models, the quality of output produced by the baselines tends to dramatically decrease at higher temperatures (only min-*p* achieves good results at $\tau = 2.0$), with the consequence of an increasing syntactic diversity due to the choice of random and meaningless tokens. Instead, the quality of the output generated by *DiffSampling* remains more stable, with small but consistent increases in diversity. 947

948

949

950

951

952

953

954

955

956

957

958

959

960

C.3 Divergence Association Task

Finally, Figure 5 reports the DAT score and the percentage of output validity of the DPO-instructed and pre-trained models with different temperature values. For the instructed model, top-p, locally typical, and η -sampling rapidly stop outputting valid lists of nouns when the temperature raises, even if the DAT score tends to be higher at $\tau = 1.5$; min-

Model:			RL	HF-instruc	ted					1	Pre-trained	[
Method		Quality		Cross	-Input	Against	Greedy		Quality		Cross	-Input	Against	-Greedy
	R-1	SIM	COH	EAD	SBERT	EAD	SBERT	R-1	SIM	COH	EAD	SBERT	EAD	SBERT
Temperature = 0.0														
Greedy	$0.22_{\pm.00}$	$0.49_{\pm.01}$	$0.72_{\pm.00}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	-	-	$0.19_{\pm.00}$	$0.45_{\pm.01}$	$0.66_{\pm.01}$	$1.11_{\pm.00}$	$0.94_{\pm.00}$	-	-
Temperature = 1.0														
Contrastive search	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.18_{\pm.00}$	$0.94_{\pm.00}$	$0.21_{\pm.01}$	$0.08_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.64_{\pm.01}$	$1.14_{\pm.00}$	$0.94_{\pm.00}$	$0.45_{\pm.01}$	$0.29_{\pm.01}$
Top- p sampling	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.21_{\pm.00}$	$0.94_{\pm.00}$	$0.30_{\pm.01}$	$0.12_{\pm.01}$	$0.16_{\pm.00}$	$0.36_{\pm.01}$	$0.50_{\pm.01}$	$1.16_{\pm.00}$	$0.93_{\pm.00}$	$0.75_{\pm.01}$	$0.55_{\pm.01}$
η -sampling	$0.22_{\pm.00}$	$0.50_{\pm .01}$	$0.71_{\pm .00}$	$1.22_{\pm.00}$	$0.94_{\pm.00}$	$0.33_{\pm.01}$	$0.13_{\pm.01}$	$0.15_{\pm.00}$	$0.35_{\pm.01}$	$0.49_{\pm.01}$	$1.19_{\pm.01}$	$0.93_{\pm.00}$	$0.78_{\pm.01}$	$0.57_{\pm.01}$
Locally typical	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.21_{\pm.00}$	$0.94_{\pm.00}$	$0.30_{\pm.01}$	$0.12_{\pm.01}$	$0.16_{\pm.00}$	$0.35_{\pm .01}$	$0.50_{\pm .01}$	$1.16_{\pm.00}$	$0.93_{\pm.00}$	$0.75_{\pm.01}$	$0.55_{\pm.01}$
Min- p sampling	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.20_{\pm.00}$	$0.94_{\pm.00}$	$0.29_{\pm.01}$	$0.11_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.61_{\pm.01}$	$1.16_{\pm.01}$	$0.93_{\pm.00}$	$0.62_{\pm.01}$	$0.40_{\pm.01}$
DiffSampling-cut	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	$0.17_{\pm.01}$	$0.06_{\pm.01}$	$0.19_{\pm.00}$	$0.45_{\pm .01}$	$0.66_{\pm.01}$	$1.13_{\pm.00}$	$0.94_{\pm.00}$	$0.25_{\pm.01}$	$0.15_{\pm.01}$
DiffSampling-lb	$0.22 \pm .00$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.20_{\pm.00}$	$0.94_{\pm.00}$	$0.27_{\pm.01}$	$0.10_{\pm.01}$	$0.17_{\pm.00}$	$0.38_{\pm.01}$	$0.53_{\pm.01}$	$1.15 \pm .00$	$0.93_{\pm.00}$	$0.71_{\pm.01}$	$0.50_{\pm.01}$
DiffSampling-minp	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.18_{\pm.00}$	$0.94_{\pm.00}$	$0.23_{\pm.01}$	$0.09_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.63_{\pm.01}$	$1.15_{\pm.01}$	$0.94_{\pm.00}$	$0.49_{\pm.01}$	$0.31_{\pm.01}$
Temperature = 0.6														
Top-p sampling	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.18_{\pm.00}$	$0.94_{\pm.00}$	$0.21_{\pm.01}$	$0.08_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.63_{\pm.01}$	$1.14_{\pm.01}$	$0.94_{\pm.00}$	$0.47_{\pm.01}$	$0.30_{\pm.01}$
η -sampling	$0.22 \pm .00$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.19_{\pm.00}$	$0.94_{\pm.00}$	$0.23_{\pm.01}$	$0.09_{\pm.01}$	$0.19_{\pm.00}$	$0.43_{\pm.01}$	$0.62_{\pm.01}$	$1.15_{\pm.01}$	$0.94_{\pm.00}$	$0.52_{\pm.01}$	$0.34_{\pm.01}$
Locally typical	$0.22 \pm .00$	$0.50_{\pm.01}$	$0.72 \pm .00$	$1.18 \pm .00$	$0.94_{\pm.00}$	$0.21_{\pm.01}$	$0.08 \pm .01$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.63_{\pm.01}$	$1.14 \pm .01$	$0.94_{\pm.00}$	$0.47_{\pm.01}$	$0.30_{\pm.01}$
Min-p sampling	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.18_{\pm.00}$	$0.94_{\pm.00}$	$0.21_{\pm.01}$	$0.08_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.64_{\pm.01}$	$1.15_{\pm.00}$	$0.94_{\pm.00}$	$0.43_{\pm.01}$	$0.27_{\pm.01}$
DiffSampling-cut	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	$0.15_{\pm.01}$	$0.06_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.66_{\pm.01}$	$1.13_{\pm.01}$	$0.94_{\pm.00}$	$0.24_{\pm.01}$	$0.14_{\pm.01}$
DiffSampling-lb	$0.22 \pm .00$	$0.50_{\pm.01}$	$0.72 \pm .00$	$1.19_{\pm.00}$	$0.94_{\pm.00}$	$0.22 \pm .01$	$0.08 \pm .01$	$0.19_{\pm.00}$	$0.43_{\pm.01}$	$0.62 \pm .01$	$1.15 \pm .00$	$0.94_{\pm.00}$	$0.52 \pm .01$	$0.34_{\pm.01}$
DiffSampling-minp	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72 \pm .00$	$1.17 \pm .00$	$0.94_{\pm.00}$	$0.20_{\pm.01}$	$0.08 \pm .01$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.64_{\pm.01}$	$1.14 \pm .01$	$0.94_{\pm.00}$	$0.42_{\pm.01}$	$0.26 \pm .01$
Temperature = 1.5														
Top-p sampling	$0.21_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.33_{\pm.01}$	$0.93_{\pm.00}$	$0.41_{\pm.01}$	$0.16_{\pm.01}$	$0.04_{\pm.00}$	$0.10_{\pm.00}$	$0.23_{\pm.01}$	$2.32_{\pm.00}$	$0.75_{\pm.00}$	$0.96_{\pm.00}$	$0.87_{\pm.01}$
η -sampling	$0.21_{\pm.00}$	$0.50_{\pm.01}$	$0.70_{\pm.00}$	$1.36_{\pm.01}$	$0.93_{\pm.00}$	$0.46_{\pm.01}$	$0.18_{\pm.01}$	$0.04_{\pm.00}$	$0.10_{\pm.00}$	$0.23_{\pm.01}$	$2.37_{\pm.01}$	$0.74_{\pm.00}$	$0.96_{\pm.00}$	$0.87_{\pm.01}$
Locally typical	$0.21_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.36_{\pm.01}$	$0.93_{\pm .00}$	$0.41_{\pm.01}$	$0.15_{\pm.01}$	$0.03_{\pm.00}$	$0.09_{\pm .00}$	$0.23_{\pm.01}$	$2.72_{\pm.01}$	$0.67_{\pm.00}$	$0.97_{\pm .00}$	$0.89_{\pm.00}$
Min-p sampling	$0.22 \pm .00$	$0.50 \pm .01$	$0.71_{\pm.00}$	$1.23 \pm .01$	$0.94_{\pm.00}$	$0.38_{\pm.01}$	$0.15_{\pm.01}$	$0.17_{\pm.00}$	$0.40_{\pm.01}$	$0.54_{\pm.01}$	$1.23 \pm .00$	$0.93 \pm .00$	$0.79_{\pm.01}$	$0.53 \pm .01$
DiffSampling-cut	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	$0.17_{\pm.01}$	$0.07_{\pm.01}$	$0.19_{\pm.00}$	$0.45_{\pm.01}$	$0.66_{\pm.01}$	$1.13_{\pm.00}$	$0.94_{\pm.00}$	$0.26_{\pm.01}$	$0.15_{\pm.01}$
DiffSampling-lb	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.21_{\pm.00}$	$0.94_{\pm.00}$	$0.30_{\pm.01}$	$0.12_{\pm.01}$	$0.14_{\pm.00}$	$0.33_{\pm.01}$	$0.46_{\pm.01}$	$1.19_{\pm.00}$	$0.93_{\pm.00}$	$0.81_{\pm.01}$	$0.61_{\pm.01}$
DiffSampling-minp	$0.22 \pm .00$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.18 \pm .00$	$0.94_{\pm.00}$	$0.25_{\pm.01}$	$0.09_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.63_{\pm.01}$	$1.15_{\pm.00}$	$0.94_{\pm.00}$	$0.52_{\pm.01}$	$0.34_{\pm.01}$
Temperature = 2.0	•													
Top-p sampling	$0.10_{\pm.00}$	$0.24_{\pm.01}$	$0.41_{\pm.01}$	$2.23_{\pm.01}$	$0.77_{\pm.00}$	$0.78_{\pm.01}$	$0.60_{\pm.01}$	$0.01_{\pm.00}$	$0.04_{\pm.00}$	$0.18 \pm .00$	$3.07_{\pm.01}$	$0.47_{\pm.00}$	$0.98_{\pm.00}$	$0.94_{\pm.00}$
η -sampling	$0.13_{\pm.00}$	$0.33_{\pm.01}$	$0.50_{\pm.01}$	$2.08_{\pm.02}$	$0.84_{\pm.00}$	$0.75_{\pm.01}$	$0.47_{\pm.01}$	$0.01_{\pm.00}$	$0.04_{\pm.00}$	$0.17_{\pm.00}$	$3.15_{\pm.01}$	$0.45_{\pm.00}$	$0.99_{\pm.00}$	$0.94_{\pm.00}$
Locally typical	$0.09_{\pm.00}$	$0.24_{\pm.01}$	$0.41_{\pm.01}$	$2.38_{\pm.01}$	$0.75_{\pm.00}$	$0.80_{\pm.01}$	$0.60_{\pm.01}$	$0.01_{\pm.00}$	$0.03_{\pm.00}$	$0.18_{\pm.00}$	$3.39_{\pm.00}$	$0.37_{\pm.00}$	$0.99_{\pm.00}$	$0.95_{\pm.00}$
Min-p sampling	$0.22 \pm .00$	$0.50 \pm .01$	$0.71_{\pm.00}$	$1.26 \pm .01$	$0.94_{\pm.00}$	$0.45_{\pm.01}$	$0.18 \pm .01$	$0.12 \pm .00$	$0.32_{\pm.01}$	$0.43_{\pm.01}$	$1.39_{\pm.00}$	$0.92_{\pm.00}$	$0.88 \pm .01$	$0.65 \pm .01$
DiffSampling-cut	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	$0.18_{\pm.01}$	$0.07_{\pm.01}$	$0.19_{\pm.00}$	$0.45_{\pm.01}$	$0.66_{\pm.01}$	$1.13_{\pm.00}$	$0.94_{\pm.00}$	$0.26_{\pm.01}$	$0.16_{\pm.01}$
DiffSampling-lb	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.22_{\pm.01}$	$0.94_{\pm.00}$	$0.31_{\pm.01}$	$0.12_{\pm.01}$	$0.14_{\pm.00}$	$0.31_{\pm.01}$	$0.44_{\pm.01}$	$1.25_{\pm.01}$	$0.93_{\pm.00}$	$0.86_{\pm.01}$	$0.65_{\pm.01}$
DiffSampling-minp	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.18_{\pm.00}$	$0.94_{\pm.00}$	$0.26_{\pm.01}$	$0.10_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.63_{\pm.01}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	$0.54_{\pm.01}$	$0.35_{\pm.01}$
Temperature = 10.0														
Top-p sampling	$0.00 \pm .00$	$0.03_{\pm.00}$	$0.17_{\pm.00}$	$3.49_{\pm.00}$	$0.28 \pm .00$	$1.00_{\pm.00}$	$0.96_{\pm.00}$	$0.01_{\pm.00}$	$0.03_{\pm.00}$	$0.16_{\pm.00}$	$3.48 \pm .00$	$0.30_{\pm.00}$	$1.00 \pm .00$	$0.95_{\pm.00}$
η -sampling	$0.00_{\pm.00}$	$0.03_{\pm.00}$	$0.16_{\pm.00}$	$3.51_{\pm.00}$	$0.30_{\pm .00}$	$1.00_{\pm.00}$	$0.95_{\pm.00}$	$0.00_{\pm.00}$	$0.03_{\pm.00}$	$0.16_{\pm.00}$	$3.51_{\pm.00}$	$0.31_{\pm.00}$	$1.00_{\pm .00}$	$0.95_{\pm.00}$
Locally typical	$0.00_{\pm.00}$	$0.02_{\pm.00}$	$0.16_{\pm.00}$	$3.51_{\pm.00}$	$0.28_{\pm.00}$	$1.00_{\pm.00}$	$0.96_{\pm.00}$	$0.00_{\pm.00}$	$0.03_{\pm.00}$	$0.16_{\pm.00}$	$3.51_{\pm.00}$	$0.29_{\pm.00}$	$1.00_{\pm.00}$	$0.95_{\pm.00}$
Min-p sampling	$0.00_{\pm.00}$	$0.03_{\pm.00}$	$0.16_{\pm.00}$	$3.51_{\pm.00}$	$0.30_{\pm.00}$	$1.00_{\pm.00}$	$0.95_{\pm.00}$	$0.00_{\pm.00}$	$0.03_{\pm.00}$	$0.16_{\pm.00}$	$3.51_{\pm.00}$	$0.31_{\pm.00}$	$1.00_{\pm.00}$	$0.95_{\pm.00}$
DiffSampling-cut	$0.22 \pm .00$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	$0.18_{\pm.01}$	$0.07_{\pm.01}$	$0.19_{\pm.00}$	$0.45_{\pm.01}$	$0.66_{\pm.01}$	$1.14_{\pm.00}$	$0.94_{\pm.00}$	$0.27_{\pm.01}$	$0.16_{\pm.01}$
DiffSampling-lb	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.23_{\pm.00}$	$0.94_{\pm.00}$	$0.34_{\pm.01}$	$0.13_{\pm.01}$	$0.11_{\pm.00}$	$0.27_{\pm.01}$	$0.38_{\pm.01}$	$1.47_{\pm.01}$	$0.92_{\pm.00}^{-0.00}$	$0.89_{\pm.01}$	$0.71_{\pm.01}$
DiffSampling-minp	0.22 ± 0.0	0.50 ± 0.1	0.72 ± 0.0	1.18 ± 01	0.94 ± 00	0.28 ± 01	0.11+ 01	0.19+00	0.44+ 01	0.62 ± 01	1.16 ± 0.0	0.94 ± 00	0.58 ± 01	0.38 ± 01

Table 6: Aggregate results over 3 seeds for the XSum dataset for the instructed model (left) and the pre-trained model (right) with different temperature values. The mean and standard error of the final score for each run are reported for cross-input diversity, whereas the mean and 95% confidence interval for the full set of answers are reported for the other metrics.

p returns a high percentage of valid outputs even at $\tau = 2.0$, but does not increase the DAT score and cannot produce anything valid at $\tau = 10.0$. Instead, the performance of our methods remains very similar across different temperatures in terms of both the DAT score and the percentage of valid outputs.

961

962 963

964

965

967

968

969

970

971

972

973

974

975

976

977

978

On the other hand, the greedy decoding strategy is less effective for the pre-trained model, which results in higher temperatures yielding better DAT scores across both the baselines and our methods (especially *DiffSampling-lb*). However, the number of valid outputs decreases faster, and top-p, locally typical, and η -sampling produce very few correct lists at a temperature of 1.5 (but with a higher DAT score). Again, min-p better manages temperatures around 1.5 and 2.0, with higher scores and still at least half of the outputs as valid, but cannot produce any correct output at a temperature of 10.0.

D Additional Experiments

D.1 Temperature Before or After Truncating

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

As described in Section 3, we apply temperature after truncating based on the minimum discrete derivative to preserve the guarantees of correctness of selected tokens. However, the de facto standard is to apply temperature before any other truncation or modification. In this section, we examine the implications of the temperature position in terms of quality and diversity.

Table 7 reports the results of our methods with temperature before (left side) and after (right side) the truncation for the GSM8K test set. As we can see, applying the temperature before causes the accuracy to degrade at higher temperatures, while



(b) Pre-trained model

Figure 5: DAT scores for our methods and the baselines with different temperature values, together with the number of valid outputs produced by each sampling strategy. The dashed line represents the score of the greedy strategy.

ensuring a slightly higher diversity. Interestingly, at $\tau = 0.6$, applying the temperature after leads to better results in terms of both accuracy and diver-

995

996

997

sity. This confirms that our choice preserves the quality as much as possible, at the cost of some additional diversity.

Method			BEFORE					AFTER		
	Accuracy	Cross-	Input	Against	-Greedy	Accuracy	Cross	-Input	Against	-Greedy
		EAD	SBERT	EAD	SBERT		EAD	SBERT	EAD	SBERT
Temperature = 0.6										
DiffSampling-cut	$66.19_{\pm.12}$	$2.04_{\pm.00}$	$0.64_{\pm.00}$	$0.10_{\pm .00}$	$0.01_{\pm .00}$	$66.74_{\pm.04}$	$2.05_{\pm.00}$	$0.64_{\pm.00}$	$0.13_{\pm.00}$	$0.02_{\pm.00}$
DiffSampling-lb	$66.54_{\pm.55}$	$2.05_{\pm.01}$	$0.64_{\pm.00}$	$0.16_{\pm.00}$	$0.02_{\pm.00}$	$66.49_{\pm.41}$	$2.06_{\pm.01}$	$0.64_{\pm.00}$	$0.18_{\pm.00}$	$0.03_{\pm .00}$
DiffSampling-minp	$66.84_{\pm.23}$	$2.05_{\pm.00}$	$0.64_{\pm .00}$	$0.14_{\pm.00}$	$0.02_{\pm.00}$	$66.84_{\pm.73}$	$2.05_{\pm .01}$	$0.64_{\pm.00}$	$0.16_{\pm .00}$	$0.02_{\pm.00}$
Temperature = 1.5										
DiffSampling-cut	$66.16_{\pm.57}$	$2.05_{\pm.00}$	$0.64_{\pm.00}$	$0.17_{\pm.00}$	$0.02_{\pm.00}$	$66.72_{\pm.36}$	$2.05_{\pm.00}$	$0.64_{\pm.00}$	$0.15_{\pm.00}$	$0.02_{\pm.00}$
DiffSampling-lb	$64.44_{\pm.22}$	$2.12_{\pm.00}$	$0.64_{\pm.00}$	$0.26_{\pm.00}$	$0.04_{\pm.00}$	$66.84_{\pm.43}$	$2.08_{\pm.00}$	$0.64_{\pm.00}$	$0.22_{\pm.00}$	$0.03_{\pm .00}$
DiffSampling-minp	$64.87_{\pm.30}$	$2.08_{\pm.00}$	$0.64_{\pm .00}$	$0.23_{\pm.00}$	$0.03_{\pm .00}$	$66.24_{\pm.36}$	$2.08_{\pm.00}$	$0.64_{\pm.00}$	$0.20_{\pm .00}$	$0.03_{\pm.00}$
Temperature = 2.0										
DiffSampling-cut	$65.50_{\pm.09}$	$2.06_{\pm.01}$	$0.64_{\pm.00}$	$0.19_{\pm.00}$	$0.03_{\pm .00}$	$66.44_{\pm.18}$	$2.05_{\pm.00}$	$0.64_{\pm.00}$	$0.15_{\pm.00}$	$0.02_{\pm.00}$
DiffSampling-lb	$38.92_{\pm.52}$	$7.51_{\pm.06}$	$0.69_{\pm.00}$	$0.56_{\pm.01}$	$0.24_{\pm.01}$	$65.73_{\pm.47}$	$2.08_{\pm.01}$	$0.64_{\pm.00}$	$0.23_{\pm .00}$	$0.03_{\pm.00}$
DiffSampling-minp	$64.19_{\pm.05}$	$2.14_{\pm.01}$	$0.64_{\pm .00}$	$0.27_{\pm.00}$	$0.04_{\pm.00}$	$66.21_{\pm.27}$	$2.06_{\pm.00}$	$0.64_{\pm.00}$	$0.21_{\pm.00}$	$0.03_{\pm.00}$
Temperature = 10.0										
DiffSampling-cut	$61.31_{\pm.21}$	$2.22_{\pm.01}$	$0.64_{\pm.00}$	$0.31_{\pm.00}$	$0.04_{\pm.00}$	$66.31_{\pm.26}$	$2.04_{\pm.00}$	$0.64_{\pm.00}$	$0.15_{\pm.00}$	$0.02_{\pm.00}$
DiffSampling-lb	$0.00_{\pm .00}$	$17.41_{\pm.03}$	$0.12_{\pm .00}$	$1.00_{\pm.00}$	$0.96_{\pm.00}$	$66.26_{\pm.72}$	$2.10_{\pm.00}$	$0.64_{\pm .00}$	$0.25_{\pm .00}$	$0.04_{\pm.00}$
DiffSampling-minp	$0.00_{\pm .00}$	$17.14_{\pm.03}$	$0.13_{\pm.00}$	$1.00_{\pm.00}$	$0.95_{\pm.00}$	$65.43_{\pm.31}$	$2.07_{\pm.00}$	$0.64_{\pm.00}$	$0.22_{\pm.00}$	$0.03_{\pm.00}$

Table 7: Accuracy and diversity of results for the GSM8K test set over 3 seeds. The mean and standard error of the final score for each run are reported for accuracy and cross-input diversity, whereas the mean and 95% confidence interval for the full set of answers are reported for against-greedy diversity.

Table 8 reports the results of our methods with temperature before (left side) and after (right side) the truncation for the MATH test set. Again, applying a higher temperature before causes the accuracy to drop quickly for the two relaxations, and smoothly for *DiffSampling-cut*, with benefits only in terms of syntactic diversity. Instead, applying the temperature after has a negligible impact on quality while fostering diversity.

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010 1011

1012

1013

1014

1015

1016

1017

1018

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1032

The same considerations hold for XSum as well. For both the instructed (Table 9) and pre-trained (Table 10) models, the quality is not preserved with the temperature before, while it is with the temperature after, although diversity does not increase in the same way. Again, the diversity at $\tau = 0.6$ is instead greater with the temperature after, even if the quality is, more or less, the same.

Finally, applying the temperature before does not seem to give benefits for the divergence association task as well. As shown by Figure 6, for both the instructed and pre-trained models, the DAT scores are very similar regardless of the temperature position, but almost no valid solutions are generated when a temperature of 10.0 is applied before truncating (and the same happens for a temperature of 2.0 in the case of *DiffSampling-lb*).

D.2 Ablation Study on the Lower Bound

We also conducted experiments on the three aforementioned case studies, varying the lower bound of the critical mass. Table 11 reports the results for the math problem-solving tasks, considering the GSM8K (left side) and MATH (right side) test sets. As expected, the against-greedy diversity scores and cross-input EAD increase together with p_{lb} ; instead, while accuracy tends to decrease with higher lower bounds, the differences are not significant, and even a quite high value (e.g., 0.8) achieves competitive results. Notably, *DiffSampling-lb* with $p_{lb} = 0.9$ performs better than or equal to top-*p* sampling (with p = 0.9) under all quality and diversity metrics, highlighting how our method can improve upon existing solutions. 1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

Table 12 reports the results for the extreme summarization task for both instructed (left side) and pre-trained (right side) models. Again, againstgreedy scores and cross-input EAD are directly correlated with the lower bound; however, they start changing only with p_{lb} around 0.7. Instead, qualitative metrics do not vary much for the instructed model, while constantly decreasing for the pre-trained model with increasing p_{lb} . In this situation, the choice of p_{lb} is relevant and requires us to decide whether to trade off quality or diversity.

Finally, Figure 7 reports the results for the divergent association task. As we would expect, the DAT score changes almost linearly between that for a lower bound of 0 (that means *DiffSampling-cut*) and 1 (that means *standard* sampling), as we reported in Section 5. Interestingly, the number of correct answers by the non-instructed model drops constantly, while it remains consistently higher in the case of the instructed model. To sum up, when greediness is desirable, a lower value of p_{lb} can lead to high quality and diversity; otherwise, increasing p_{lb} improves diversity, but the cost in terms

Method			BEFORE					AFTER		
	Accuracy	Cross-	Input	Against	-Greedy	Accuracy	Cross	-Input	Against	-Greedy
		EAD	SBERT	EAD	SBERT		EAD	SBERT	EAD	SBERT
Temperature = 0.6										
DiffSampling-cut	$21.44_{\pm.12}$	$5.69_{\pm.01}$	$0.80_{\pm.00}$	$0.22_{\pm.00}$	$0.06_{\pm.00}$	$21.52_{\pm.13}$	$5.72_{\pm.00}$	$0.80_{\pm.00}$	$0.25_{\pm.00}$	$0.07_{\pm.00}$
DiffSampling-lb	$21.69_{\pm.28}$	$5.76_{\pm.01}$	$0.80_{\pm .00}$	$0.29_{\pm.00}$	$0.08_{\pm.00}$	$21.09_{\pm.11}$	$5.83_{\pm.01}$	$0.80_{\pm .00}$	$0.32_{\pm.00}$	$0.09_{\pm .00}$
DiffSampling-minp	$21.36_{\pm.21}$	$5.72_{\pm.01}$	$0.80_{\pm.00}$	$0.27_{\pm.00}$	$0.08_{\pm .00}$	$20.79_{\pm.10}$	$5.78_{\pm.01}$	$0.80_{\pm.00}$	$0.30_{\pm.00}$	$0.08_{\pm.00}$
<i>Temperature</i> = 1.5										
DiffSampling-cut	$21.15_{\pm.09}$	$5.78_{\pm.01}$	$0.80_{\pm.00}$	$0.30_{\pm.00}$	$0.08_{\pm.00}$	$21.36_{\pm.15}$	$5.73_{\pm.00}$	$0.80_{\pm.00}$	$0.27_{\pm.00}$	$0.07_{\pm.00}$
DiffSampling-lb	$19.46_{\pm.12}$	$6.50_{\pm.00}$	$0.80_{\pm.00}$	$0.40_{\pm.00}$	$0.12_{\pm.00}$	$20.52_{\pm.34}$	$6.03_{\pm.01}$	$0.80_{\pm.00}$	$0.36_{\pm.00}$	$0.10_{\pm .00}$
DiffSampling-minp	$20.02_{\pm.06}$	$6.08_{\pm.01}$	$0.80_{\pm.00}$	$0.37_{\pm.00}$	$0.11_{\pm.00}$	$20.79_{\pm.14}$	$5.88_{\pm.01}$	$0.80_{\pm.00}$	$0.34_{\pm.00}$	$0.10_{\pm .00}$
Temperature = 2.0										
DiffSampling-cut	$21.25_{\pm.10}$	$5.85_{\pm.00}$	$0.80_{\pm.00}$	$0.32_{\pm.00}$	$0.09_{\pm.00}$	$21.66_{\pm.20}$	$5.71_{\pm.01}$	$0.80_{\pm.00}$	$0.27_{\pm.00}$	$0.08_{\pm .00}$
DiffSampling-lb	$6.96_{\pm.25}$	$40.17_{\pm.11}$	$0.68_{\pm.00}$	$0.81_{\pm.00}$	$0.51_{\pm.01}$	$20.34_{\pm.08}$	$6.09_{\pm .02}$	$0.80_{\pm.00}$	$0.37_{\pm.00}$	$0.11_{\pm .00}$
DiffSampling-minp	$19.44_{\pm.13}$	$6.35_{\pm.02}$	$0.80_{\pm.00}$	$0.40_{\pm.00}$	$0.11_{\pm.00}$	$20.59 \pm .13$	$5.90_{\pm .02}$	$0.80_{\pm.00}$	$0.34_{\pm.00}$	$0.10_{\pm .00}$
Temperature = 10.0										
DiffSampling-cut	$16.63_{\pm.12}$	$6.78_{\pm.01}$	$0.80_{\pm.00}$	$0.43_{\pm.00}$	$0.12_{\pm.00}$	$21.22_{\pm.11}$	$5.74_{\pm.01}$	$0.80_{\pm.00}$	$0.28_{\pm.00}$	$0.08_{\pm .00}$
DiffSampling-lb	$0.00_{\pm.00}$	$59.18 \pm .02$	$0.12_{\pm.00}$	$1.00_{\pm.00}$	$1.00_{\pm .00}$	$19.88_{\pm.03}$	$6.20_{\pm .01}$	$0.80_{\pm.00}$	$0.39_{\pm .00}$	$0.11_{\pm .00}$
DiffSampling-minp	$0.00_{\pm .00}$	$58.15_{\pm.04}$	$0.13_{\pm.00}$	$1.00_{\pm.00}$	$1.00_{\pm .00}$	$21.17_{\pm.14}$	$5.94_{\pm.02}$	$0.80_{\pm.00}$	$0.35_{\pm .00}$	$0.10_{\pm.00}$

Table 8: Accuracy and diversity of results for the MATH test set over 3 seeds. The mean and standard error of the final score for each run are reported for accuracy and cross-input diversity, whereas the mean and 95% confidence interval for the full set of answers are reported for against-greedy diversity.

Method				BEFORE							AFTER			
		Quality		Cross	Input	Against	Greedy		Quality		Cross	-Input	Against	-Greedy
	R-1	SIM	COH	EAD	SBERT	EAD	SBERT	R-1	SIM	COH	EAD	SBERT	EAD	SBERT
Temperature = 0.6														
DiffSampling-cut	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	$0.12_{\pm.01}$	$0.05_{\pm.00}$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	$0.15_{\pm.01}$	$0.06_{\pm.01}$
DiffSampling-lb	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.17_{\pm.00}$	$0.94_{\pm.00}$	$0.18_{\pm.01}$	$0.07_{\pm.01}$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.19_{\pm.00}$	$0.94_{\pm.00}$	$0.22_{\pm.01}$	$0.08_{\pm.01}$
DiffSampling-minp	$0.22_{\pm.00}$	$0.50 \pm .01$	$0.72 \pm .00$	$1.17_{\pm.00}$	$0.94_{\pm.00}$	$0.17_{\pm .01}$	$0.06 \pm .01$	$0.22_{\pm.00}$	$0.50 {\scriptstyle \pm .01}$	$0.72_{\pm .00}$	$1.17_{\pm.00}$	$0.94_{\pm .00}$	$0.20_{\pm.01}$	$0.08 {\scriptstyle \pm .01}$
Temperature = 1.5														
DiffSampling-cut	$0.22_{\pm.00}$	$0.50_{\pm .01}$	$0.72_{\pm.00}$	$1.17_{\pm.00}$	$0.94_{\pm.00}$	$0.21_{\pm.01}$	$0.08 \pm .01$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	$0.17_{\pm.01}$	$0.07_{\pm .01}$
DiffSampling-lb	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.30_{\pm.01}$	$0.93_{\pm.00}$	$0.37_{\pm.01}$	$0.14_{\pm.01}$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm .00}$	$1.21_{\pm.00}$	$0.94_{\pm.00}$	$0.30_{\pm.01}$	$0.12_{\pm.01}$
DiffSampling-minp	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.20_{\pm.00}$	$0.94_{\pm.00}$	$0.30_{\pm.01}$	$0.12_{\pm.01}$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.18_{\pm.00}$	$0.94_{\pm.00}$	$0.25_{\pm.01}$	$0.09_{\pm.01}$
Temperature = 2.0														
DiffSampling-cut	$0.22_{\pm.00}$	$0.50_{\pm .01}$	$0.72_{\pm.00}$	$1.17_{\pm.00}$	$0.94_{\pm.00}$	$0.24_{\pm.01}$	$0.09_{\pm.01}$	$0.22_{\pm.00}$	$0.50_{\pm .01}$	$0.72_{\pm .00}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	$0.18 \pm .01$	$0.07_{\pm .01}$
DiffSampling-lb	$0.13_{\pm.00}$	$0.33_{\pm.01}$	$0.52_{\pm.01}$	$2.05 \pm .01$	$0.85 _{\pm .00}$	$0.68 \pm .01$	$0.44_{\pm.01}$	$0.22_{\pm.00}$	$0.50 {\scriptstyle \pm .01}$	$0.71 \scriptstyle \pm .00$	$1.22 \pm .01$	$0.94_{\pm .00}$	$0.31_{\pm.01}$	$0.12_{\pm.01}$
DiffSampling-minp	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.22_{\pm.01}$	$0.94_{\pm.00}$	$0.36_{\pm.01}$	$0.14_{\pm.01}$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.18_{\pm.00}$	$0.94_{\pm.00}$	$0.26_{\pm.01}$	$0.10_{\pm.01}$
Temperature = 10.0														
DiffSampling-cut	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.28_{\pm.00}$	$0.94_{\pm.00}$	$0.42_{\pm.01}$	$0.15_{\pm.01}$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	$0.18_{\pm.01}$	$0.07_{\pm.01}$
DiffSampling-lb	$0.00_{\pm.00}$	$0.03_{\pm.00}$	$0.16 \scriptstyle \pm .00$	$3.51_{\pm.00}$	$0.30_{\pm.00}$	$1.00_{\pm .00}$	$0.95_{\pm.00}$	$0.22_{\pm.00}$	$0.50 {\scriptstyle \pm .01}$	$0.71 \scriptstyle \pm .00$	$1.23_{\pm.00}$	$0.94_{\pm.00}$	$0.34_{\pm.01}$	$0.13 {\scriptstyle \pm .01}$
DiffSampling-minp	$0.01_{\pm.00}$	$0.03_{\pm.00}$	$0.16_{\pm.00}$	$3.49_{\pm.00}$	$0.32_{\pm.00}$	$1.00_{\pm .00}$	$0.96_{\pm.00}$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.18_{\pm.01}$	$0.94_{\pm.00}$	$0.28_{\pm.01}$	$0.11_{\pm.01}$

Table 9: Quality and diversity of results for the XSum test set with the instructed model over 3 seeds. The mean and standard error of the final score for each run are reported for cross-input diversity, whereas the mean and 95% confidence interval for the full set of answers are reported for the other metrics.

Method				BEFORE							AFTER			
		Quality		Cross	Input	Against	Greedy		Quality		Cross	-Input	Against	-Greedy
	R-1	SIM	COH	EAD	SBERT	EAD	SBERT	R-1	SIM	COH	EAD	SBERT	EAD	SBERT
Temperature = 0.6														
DiffSampling-cut	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.66_{\pm.01}$	$1.13_{\pm.00}$	$0.94_{\pm.00}$	$0.19_{\pm.01}$	$0.11_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.66_{\pm.01}$	$1.13_{\pm.01}$	$0.94_{\pm.00}$	$0.24_{\pm.01}$	$0.14_{\pm.01}$
DiffSampling-lb	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.64_{\pm.01}$	$1.13_{\pm.01}$	$0.94_{\pm.00}$	$0.43_{\pm.01}$	$0.27_{\pm.01}$	$0.19_{\pm.00}$	$0.43_{\pm.01}$	$0.62_{\pm.01}$	$1.15_{\pm.00}$	$0.94_{\pm.00}$	$0.52_{\pm.01}$	$0.34_{\pm.01}$
DiffSampling-minp	$0.20_{\pm.00}$	$0.45 {\scriptstyle \pm.01}$	$0.65 {\scriptstyle \pm .01}$	$1.13_{\pm.00}$	$0.94_{\pm.00}$	$0.33_{\pm.01}$	$0.20_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.64_{\pm.01}$	$1.14_{\pm.01}$	$0.94_{\pm.00}$	$0.42_{\pm.01}$	$0.26 {\scriptstyle \pm.01}$
Temperature = 1.5														
DiffSampling-cut	$0.19_{\pm.00}$	$0.45_{\pm.01}$	$0.66_{\pm.01}$	$1.15_{\pm.01}$	$0.94_{\pm.00}$	$0.30_{\pm.01}$	$0.18_{\pm.01}$	$0.19_{\pm.00}$	$0.45_{\pm.01}$	$0.66_{\pm.01}$	$1.13_{\pm.00}$	$0.94_{\pm.00}$	$0.26_{\pm.01}$	$0.15_{\pm.01}$
DiffSampling-lb	$0.05_{\pm.00}$	$0.14_{\pm.01}$	$0.26_{\pm.01}$	$2.11_{\pm.01}$	$0.82_{\pm.00}$	$0.95_{\pm .00}$	$0.83_{\pm.01}$	$0.14_{\pm.00}$	$0.33_{\pm.01}$	$0.46_{\pm.01}$	$1.19_{\pm.00}$	$0.93_{\pm.00}$	$0.81_{\pm .01}$	$0.61_{\pm .01}$
DiffSampling-minp	$0.18_{\pm.00}$	$0.42_{\pm.01}$	$0.60_{\pm .01}$	$1.16_{\pm.00}$	$0.93_{\pm.00}$	$0.63_{\pm.01}$	$0.42_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.63_{\pm.01}$	$1.15_{\pm.00}$	$0.94_{\pm.00}$	$0.52_{\pm.01}$	$0.34_{\pm.01}$
Temperature = 2.0														
DiffSampling-cut	$0.20_{\pm.00}$	$0.45_{\pm .01}$	$0.66 {\pm} .01$	$1.16_{\pm.01}$	$0.94_{\pm.00}$	$0.34_{\pm.01}$	$0.20_{\pm.01}$	$0.19_{\pm.00}$	$0.45_{\pm.01}$	$0.66 {\pm} .01$	$1.13_{\pm.00}$	$0.94_{\pm.00}$	$0.26_{\pm.01}$	$0.16 \pm .01$
DiffSampling-lb	$0.01_{\pm.00}$	$0.04_{\pm.00}$	$0.18_{\pm.00}$	$3.16_{\pm.01}$	$0.46_{\pm .00}$	$0.99_{\pm .00}$	$0.94_{\pm.00}$	$0.14_{\pm.00}$	$0.31_{\pm.01}$	$0.44_{\pm.01}$	$1.25_{\pm.01}$	$0.93_{\pm.00}$	$0.86_{\pm.01}$	$0.65_{\pm.01}$
DiffSampling-minp	$0.17_{\pm.00}$	$0.41_{\pm .01}$	$0.56_{\pm.01}$	$1.18_{\pm.01}$	$0.93_{\pm.00}$	$0.75_{\pm.01}$	$0.50_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.63_{\pm.01}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	$0.54_{\pm.01}$	$0.35_{\pm.01}$
Temperature = 10.0														
DiffSampling-cut	$0.10_{\pm.00}$	$0.25_{\pm.01}$	$0.39_{\pm.01}$	$1.51_{\pm.02}$	$0.90_{\pm.00}$	$0.64_{\pm.01}$	$0.58_{\pm.01}$	$0.19_{\pm.00}$	$0.45_{\pm.01}$	$0.66_{\pm.01}$	$1.14_{\pm.00}$	$0.94_{\pm.00}$	$0.27_{\pm.01}$	$0.16_{\pm.01}$
DiffSampling-lb	$0.00 \pm .00$	$0.03 \pm .00$	$0.16 \scriptstyle \pm .00$	$3.51_{\pm.00}$	$0.31_{\pm.00}$	$1.00_{\pm .00}$	$0.95_{\pm.00}$	$0.11_{\pm.00}$	$0.27 \pm .01$	$0.38_{\pm.01}$	$1.47_{\pm.01}$	$0.92 \pm .00$	$0.89_{\pm.01}$	$0.71_{\pm .01}$
DiffSampling-minp	$0.00_{\pm .00}$	$0.03_{\pm.00}$	$0.16_{\pm.00}$	$3.51_{\pm.00}$	$0.31_{\pm.00}$	$1.00_{\pm .00}$	$0.95_{\pm.00}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.62_{\pm.01}$	$1.16_{\pm.00}$	$0.94_{\pm .00}$	$0.58_{\pm.01}$	$0.38_{\pm.01}$

Table 10: Quality and diversity of results for the XSum test set with the pre-trained model over 3 seeds. The mean and standard error of the final score for each run are reported for cross-input diversity, whereas the mean and 95% confidence interval for the full set of answers are reported for the other metrics.



(b) Pre-trained model

Figure 6: DAT scores and validity percentage of outputs with temperature scaling before (left) and after (right) the truncation. The dashed line represents the score of the greedy strategy.

of validity is not negligible and requires careful consideration.

Dataset:			GSM8K					MATH		
Method	Accuracy	Cross	-Input	Against	-Greedy	Accuracy	Cross	-Input	Against	-Greedy
DiffSampling-lb		EAD	SBERT	EAD	SBERT		EAD	SBERT	EAD	SBERT
$p_{lb} = 0.0$	$66.36_{\pm.23}$	$2.04_{\pm.00}$	$0.64_{\pm.00}$	$0.14_{\pm.00}$	$0.02_{\pm.00}$	$21.38_{\pm.20}$	$5.71_{\pm.01}$	$0.80_{\pm.00}$	$0.27_{\pm.00}$	$0.07_{\pm.00}$
$p_{lb} = 0.1$	$66.46_{\pm.34}$	$2.05_{\pm.00}$	$0.64_{\pm.00}$	$0.14_{\pm.00}$	$0.02_{\pm.00}$	$20.95_{\pm.20}$	$5.72_{\pm.01}$	$0.80_{\pm.00}$	$0.27_{\pm.00}$	$0.07_{\pm .00}$
$p_{lb} = 0.2$	$66.46_{\pm.34}$	$2.05_{\pm.00}$	$0.64_{\pm.00}$	$0.14_{\pm.00}$	$0.02_{\pm.00}$	$20.95_{\pm.20}$	$5.72_{\pm.01}$	$0.80_{\pm .00}$	$0.27_{\pm.00}$	$0.07_{\pm.00}$
$p_{lb} = 0.3$	$66.79_{\pm.40}$	$2.04_{\pm.00}$	$0.64_{\pm.00}$	$0.14_{\pm.00}$	$0.02_{\pm.00}$	$21.30_{\pm.08}$	$5.73_{\pm.00}$	$0.80_{\pm .00}$	$0.27_{\pm.00}$	$0.07_{\pm.00}$
$p_{lb} = 0.4$	$66.57_{\pm.39}$	$2.06_{\pm.00}$	$0.64_{\pm.00}$	$0.14_{\pm.00}$	$0.02_{\pm.00}$	$21.08 \pm .11$	$5.73_{\pm.02}$	$0.80_{\pm.00}$	$0.27_{\pm.00}$	$0.07_{\pm .00}$
$p_{lb} = 0.5$	$67.17_{\pm.41}$	$2.04_{\pm.00}$	$0.64_{\pm.00}$	$0.15_{\pm.00}$	$0.02_{\pm.00}$	$21.18_{\pm.41}$	$5.74_{\pm.01}$	$0.80_{\pm.00}$	$0.28_{\pm.00}$	$0.08_{\pm.00}$
$p_{lb} = 0.6$	$66.67_{\pm.37}$	$2.05_{\pm.00}$	$0.64_{\pm.00}$	$0.16_{\pm.00}$	$0.02_{\pm.00}$	$21.18_{\pm.22}$	$5.79_{\pm.02}$	$0.80_{\pm.00}$	$0.30_{\pm.00}$	$0.09_{\pm .00}$
$p_{lb} = 0.7$	$65.58_{\pm.19}$	$2.06_{\pm.00}$	$0.64_{\pm.00}$	$0.18_{\pm.00}$	$0.03_{\pm .00}$	$21.14_{\pm.15}$	$5.86_{\pm.01}$	$0.80_{\pm.00}$	$0.32_{\pm.00}$	$0.09_{\pm .00}$
$p_{lb} = 0.8$	$66.92_{\pm.08}$	$2.07_{\pm.00}$	$0.64_{\pm.00}$	$0.20_{\pm.00}$	$0.03_{\pm .00}$	$20.78 \pm .14$	$6.00_{\pm .01}$	$0.80_{\pm.00}$	$0.35_{\pm.00}$	$0.10_{\pm .00}$
$p_{lb} = 0.9$	$65.18_{\pm.65}$	$2.09_{\pm.01}$	$0.64_{\pm.00}$	$0.23_{\pm.00}$	$0.03_{\pm.00}$	$20.20_{\pm .08}$	$6.11_{\pm .02}$	$0.80_{\pm.00}$	$0.37_{\pm.00}$	$0.10_{\pm .00}$
$p_{lb} = 1.0$	$64.87_{\pm.20}$	$2.12_{\pm.00}$	$0.64_{\pm .00}$	$0.25_{\pm.00}$	$0.04_{\pm.00}$	$19.46_{\pm.19}$	$6.36_{\pm.01}$	$0.80_{\pm.00}$	$0.39_{\pm.00}$	$0.11_{\pm .00}$

Table 11: Ablation study on the p_{lb} value over 3 seeds for the GSM8K (left) and MATH (right) test sets. The mean and standard error of the final score for each run are reported for accuracy and cross-input diversity, whereas the mean and 95% confidence interval for the full set of answers are reported for against-greedy diversity.

Model:		RLHF-instructed					Pre-trained							
Method		Quality		Cross-Input Against-Greed			-Greedy	Quality			Cross-Input		Against-Greedy	
DiffSampling-lb	R-1	SIM	COH	EAD	SBERT	EAD	SBERT	R-1	SIM	COH	EAD	SBERT	EAD	SBERT
$p_{lb} = 0.0$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	$0.17_{\pm.01}$	$0.06_{\pm.01}$	$0.19_{\pm.00}$	$0.45_{\pm.01}$	$0.66_{\pm.01}$	$1.13_{\pm.00}$	$0.94_{\pm.00}$	$0.25_{\pm.01}$	$0.15_{\pm.01}$
$p_{lb} = 0.1$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.17_{\pm.00}$	$0.94_{\pm.00}$	$0.17_{\pm.01}$	$0.07_{\pm.01}$	$0.19_{\pm.00}$	$0.45_{\pm.01}$	$0.66_{\pm.01}$	$1.13_{\pm.01}$	$0.94_{\pm.00}$	$0.26_{\pm.01}$	$0.16_{\pm.01}$
$p_{lb} = 0.2$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.17_{\pm.00}$	$0.94_{\pm.00}$	$0.17_{\pm.01}$	$0.07_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.65_{\pm.01}$	$1.11_{\pm.00}$	$0.94_{\pm.00}$	$0.35_{\pm.01}$	$0.23_{\pm.01}$
$p_{lb} = 0.3$	$0.22 \pm .00$	$0.50 \pm .01$	$0.72 \pm .00$	$1.18 \pm .00$	$0.94_{\pm.00}$	$0.18 \pm .01$	$0.07_{\pm.01}$	$0.19_{\pm.00}$	$0.43_{\pm.01}$	$0.62 \pm .01$	$1.11_{\pm.00}$	$0.94_{\pm.00}$	$0.44_{\pm.01}$	$0.30_{\pm .01}$
$p_{lb} = 0.4$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.18_{\pm.00}$	$0.94_{\pm.00}$	$0.18_{\pm.01}$	$0.07_{\pm.01}$	$0.19_{\pm.00}$	$0.43_{\pm.01}$	$0.61_{\pm.01}$	$1.11_{\pm.01}$	$0.94_{\pm.00}$	$0.51_{\pm.01}$	$0.35_{\pm.01}$
$p_{lb} = 0.5$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.18_{\pm.00}$	$0.94_{\pm.00}$	$0.18_{\pm.01}$	$0.07_{\pm.01}$	$0.19_{\pm.00}$	$0.42_{\pm.01}$	$0.60_{\pm .01}$	$1.11_{\pm.00}$	$0.94_{\pm.00}$	$0.56_{\pm.01}$	$0.38_{\pm.01}$
$p_{lb} = 0.6$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.18_{\pm.00}$	$0.94_{\pm.00}$	$0.20_{\pm.01}$	$0.08_{\pm.01}$	$0.18_{\pm.00}$	$0.41_{\pm.01}$	$0.57_{\pm.01}$	$1.10_{\pm.00}$	$0.94_{\pm.00}$	$0.61_{\pm.01}$	$0.43_{\pm.01}$
$p_{lb} = 0.7$	$0.22 \pm .00$	$0.50 \pm .01$	$0.72 \pm .00$	$1.19_{\pm.00}$	$0.94_{\pm.00}$	$0.23_{\pm.01}$	$0.09 \pm .01$	$0.18 \pm .00$	$0.40 {\scriptstyle \pm .01}$	$0.56 \scriptstyle \pm .01$	$1.14_{\pm.01}$	$0.93_{\pm.00}$	$0.67_{\pm .01}$	$0.47_{\pm .01}$
$p_{lb} = 0.8$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.20_{\pm.00}$	$0.94_{\pm.00}$	$0.27_{\pm.01}$	$0.10_{\pm.01}$	$0.17_{\pm.00}$	$0.38_{\pm.01}$	$0.53_{\pm.01}$	$1.15_{\pm.00}$	$0.93_{\pm.00}$	$0.71_{\pm.01}$	$0.50_{\pm.01}$
$p_{lb} = 0.9$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.21_{\pm.00}$	$0.94_{\pm.00}$	$0.30_{\pm.01}$	$0.12_{\pm.01}$	$0.15_{\pm.00}$	$0.35_{\pm.01}$	$0.50_{\pm.01}$	$1.17_{\pm.00}$	$0.93_{\pm.00}$	$0.76_{\pm.01}$	$0.56_{\pm.01}$
$p_{lb} = 1.0$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.22_{\pm.00}$	$0.94_{\pm.00}$	$0.34_{\pm.01}$	$0.13_{\pm.01}$	$0.14_{\pm.00}$	$0.31_{\pm.01}$	$0.43_{\pm.01}$	$1.21_{\pm.01}$	$0.93_{\pm.00}$	$0.80_{\pm.01}$	$0.62_{\pm.01}$

Table 12: Ablation study on the p_{lb} value over 3 seeds for the XSum dataset for the instructed model (left) and the pre-trained model (right). The mean and standard error of the final score for each run are reported for cross-input diversity, whereas the mean and 95% confidence interval for the full set of answers are reported for the other metrics.

Dataset:		GSM8K						MATH				
Method	Accuracy	Cross	-Input	Against	Against-Greedy		Cross	-Input	Against	-Greedy		
DiffSampling-minp		EAD	SBERT	EAD	SBERT		EAD	SBERT	EAD	SBERT		
$p_{min} = 0.0$	$64.87_{\pm.20}$	$2.12_{\pm.00}$	$0.64_{\pm.00}$	$0.25_{\pm.00}$	$0.04_{\pm.00}$	$19.46_{\pm.19}$	$6.36_{\pm.01}$	$0.80_{\pm.00}$	$0.39_{\pm.00}$	$0.11_{\pm .00}$		
$p_{min} = 0.1$	$65.48_{\pm.60}$	$2.09_{\pm.01}$	$0.64_{\pm.00}$	$0.23_{\pm.00}$	$0.03_{\pm.00}$	$20.18 \pm .08$	$6.06_{\pm.00}$	$0.80_{\pm.00}$	$0.36_{\pm.00}$	$0.10_{\pm .00}$		
$p_{min} = 0.2$	$65.48_{\pm.41}$	$2.07_{\pm.00}$	$0.64_{\pm.00}$	$0.21_{\pm.00}$	$0.03_{\pm .00}$	$20.65_{\pm.29}$	$5.93_{\pm.01}$	$0.80_{\pm .00}$	$0.34_{\pm.00}$	$0.10_{\pm.00}$		
$p_{min} = 0.3$	$66.44_{\pm.35}$	$2.05_{\pm.00}$	$0.64_{\pm.00}$	$0.19_{\pm.00}$	$0.03_{\pm.00}$	$21.13_{\pm.08}$	$5.87_{\pm.01}$	$0.80_{\pm .00}$	$0.33_{\pm.00}$	$0.09_{\pm .00}$		
$p_{min} = 0.4$	$66.59_{\pm.48}$	$2.05_{\pm.00}$	$0.64_{\pm.00}$	$0.17_{\pm.00}$	$0.02_{\pm.00}$	$21.41_{\pm.07}$	$5.79_{\pm.01}$	$0.80_{\pm.00}$	$0.31_{\pm.00}$	$0.09_{\pm .00}$		
$p_{min} = 0.5$	$66.67_{\pm.07}$	$2.04_{\pm.00}$	$0.64_{\pm.00}$	$0.15_{\pm.00}$	$0.02_{\pm.00}$	$21.23_{\pm.13}$	$5.75_{\pm.01}$	$0.80_{\pm.00}$	$0.28_{\pm.00}$	$0.08_{\pm.00}$		
$p_{min} = 0.6$	$66.64_{\pm.29}$	$2.04_{\pm.00}$	$0.64_{\pm.00}$	$0.14_{\pm.00}$	$0.02_{\pm.00}$	$21.67_{\pm.13}$	$5.72_{\pm.01}$	$0.80_{\pm.00}$	$0.27_{\pm.00}$	$0.08_{\pm .00}$		
$p_{min} = 0.7$	$66.29_{\pm.27}$	$2.04_{\pm.00}$	$0.64_{\pm.00}$	$0.14_{\pm.00}$	$0.02_{\pm.00}$	$21.25_{\pm.37}$	$5.72_{\pm.00}$	$0.80_{\pm.00}$	$0.27_{\pm.00}$	$0.07_{\pm.00}$		
$p_{min} = 0.8$	$66.21_{\pm.32}$	$2.04_{\pm.00}$	$0.64_{\pm.00}$	$0.14_{\pm.00}$	$0.02_{\pm .00}$	$21.16_{\pm.28}$	$5.70_{\pm.01}$	$0.80_{\pm.00}$	$0.27_{\pm.00}$	$0.07_{\pm .00}$		
$p_{min} = 0.9$	$66.21_{\pm.32}$	$2.04_{\pm.00}$	$0.64_{\pm.00}$	$0.14_{\pm.00}$	$0.02_{\pm.00}$	$21.25 \pm .35$	$5.70_{\pm.01}$	$0.80_{\pm.00}$	$0.27_{\pm.00}$	$0.07_{\pm.00}$		
$p_{min} = 1.0$	$66.36_{\pm,23}$	$2.04_{\pm.00}$	$0.64_{\pm 00}$	$0.14_{\pm.00}$	$0.02_{\pm 0.00}$	$21.38_{\pm,20}$	$5.71_{\pm.01}$	$0.80_{\pm.00}$	$0.27_{\pm.00}$	$0.07_{\pm.00}$		

Table 13: Ablation study on the p_{min} value over 3 seeds for the GSM8K (left) and MATH (right) test sets. The mean and standard error of the final score for each run are reported for accuracy and cross-input diversity, whereas the mean and 95% confidence interval for the full set of answers are reported for against-greedy diversity.

D.3 Ablation Study on the Dynamic Upper Bound

Finally, we conducted experiments on the three aforementioned case studies, varying the dynamic upper bound of the truncated tokens p_{min} .

Table 13 reports the results for the math problemsolving tasks, considering the GSM8K (left side) and MATH (right side) test sets. As expected, the against-greedy diversity scores and cross-input EAD decrease together with p_{min} , plateauing at $p_{min} = 0.5$ (after that, results are comparable with *DiffSampling-cut*); specularly, accuracy is lower at lower p_{min} , but reaches a competitive score even at $p_{min} = 0.3$.

The same holds for XSum as well. As shown in

1075

1076

1079

1080

1081

1082

1071

1068



Figure 7: DAT scores and output validity percentage for *DiffSampling-lb* when varying the p_{lb} parameter. The dashed line represents the score of the greedy strategy.

Table 14, diversity decreases when increasing p_{min} and plateaus at 0.5, while quality rapidly increases for the pre-trained model and is almost constant for the instructed model.

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1094

The same considerations are even more apparent for the divergent association task with Figure 8. While behaving differently for the instructed and pre-trained models, the DAT score plateaus around $p_{min} = 0.5$. On the other hand, the percentage of valid outputs is close to 100% for all p_{min} values when considering the instructed model, and linearly increases when considering the pre-trained model. To sum up, values above 0.5 are not different from1095DiffSampling-cut, while lower p_{min} can help foster1096diversity with small loss in accuracy, especially for1097instructed models.1098

E Qualitative Analysis

In the following two subsections, we present and
discuss some generated solutions from our methods1100and the greedy, top-p, and min-p strategies at dif-
ferent temperatures for the math problem-solving1103and summarization tasks.1104

Model:		RLHF-instructed					Pre-trained							
Method		Quality		Cross-Input Against-		t-Greedy Quality			Cross-Input		Against-Greedy			
DiffSampling-minp	R-1	SIM	COH	EAD	SBERT	EAD	SBERT	R-1	SIM	COH	EAD	SBERT	EAD	SBERT
$p_{min} = 0.0$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.71_{\pm.00}$	$1.22_{\pm.00}$	$0.94_{\pm.00}$	$0.34_{\pm.01}$	$0.13_{\pm.01}$	$0.14_{\pm.00}$	$0.31_{\pm.01}$	$0.43_{\pm.01}$	$1.21_{\pm.01}$	$0.93_{\pm.00}$	$0.80_{\pm.01}$	$0.62_{\pm.01}$
$p_{min} = 0.1$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.20_{\pm.00}$	$0.94_{\pm.00}$	$0.29_{\pm.01}$	$0.11_{\pm.01}$	$0.19_{\pm.00}$	$0.43_{\pm.01}$	$0.60_{\pm .01}$	$1.16_{\pm.00}$	$0.93_{\pm.00}$	$0.62_{\pm.01}$	$0.41_{\pm.01}$
$p_{min} = 0.2$	$0.22 \pm .00$	$0.50_{\pm.01}$	$0.72_{\pm .00}$	$1.19_{\pm.00}$	$0.94_{\pm.00}$	$0.26_{\pm.01}$	$0.10_{\pm.01}$	$0.19_{\pm.00}$	$0.43_{\pm.01}$	$0.62 \scriptstyle \pm .01$	$1.15 \pm .00$	$0.93 \scriptstyle \pm .00$	$0.55_{\pm .01}$	$0.36_{\pm.01}$
$p_{min} = 0.3$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.18_{\pm.00}$	$0.94_{\pm.00}$	$0.23_{\pm.01}$	$0.09_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.63_{\pm.01}$	$1.15_{\pm.01}$	$0.94_{\pm.00}$	$0.49_{\pm.01}$	$0.31_{\pm.01}$
$p_{min} = 0.4$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.18_{\pm.00}$	$0.94_{\pm.00}$	$0.20_{\pm.01}$	$0.08_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.64_{\pm.01}$	$1.13_{\pm.00}$	$0.94_{\pm.00}$	$0.42_{\pm.01}$	$0.27_{\pm.01}$
$p_{min} = 0.5$	$0.22 \pm .00$	$0.50_{\pm.01}$	$0.72_{\pm .00}$	$1.17_{\pm.00}$	$0.94_{\pm.00}$	$0.18 \pm .01$	$0.07_{\pm.01}$	$0.19_{\pm.00}$	$0.45 {\scriptstyle \pm .01}$	$0.65 {\scriptstyle \pm .01}$	$1.12_{\pm.00}$	$0.94_{\pm .00}$	$0.35_{\pm .01}$	$0.22 \pm .01$
$p_{min} = 0.6$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.16 \pm .00$	$0.94_{\pm.00}$	$0.17_{\pm.01}$	$0.06 \pm .01$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.65 {\scriptstyle \pm .01}$	$1.13_{\pm.00}$	$0.94_{\pm.00}$	$0.31_{\pm.01}$	$0.19_{\pm.01}$
$p_{min} = 0.7$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	$0.17_{\pm.01}$	$0.06_{\pm.01}$	$0.19_{\pm.00}$	$0.44_{\pm.01}$	$0.66_{\pm .01}$	$1.13_{\pm.00}$	$0.94_{\pm.00}$	$0.27_{\pm.01}$	$0.16_{\pm.01}$
$p_{min} = 0.8$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	$0.17_{\pm.01}$	$0.06_{\pm.01}$	$0.19_{\pm.00}$	$0.45_{\pm.01}$	$0.66_{\pm.01}$	$1.13_{\pm.00}$	$0.94_{\pm.00}$	$0.26_{\pm.01}$	$0.15_{\pm.01}$
$p_{min} = 0.9$	$0.22_{\pm.00}$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.16 \pm .00$	$0.94_{\pm.00}$	$0.17_{\pm.01}$	$0.06 \pm .01$	$0.19_{\pm.00}$	$0.45_{\pm .01}$	$0.66_{\pm.01}$	$1.13_{\pm.00}$	$0.94_{\pm.00}$	$0.25_{\pm.01}$	$0.15_{\pm.01}$
$p_{min} = 1.0$	$0.22 \pm .00$	$0.50_{\pm.01}$	$0.72_{\pm.00}$	$1.16_{\pm.00}$	$0.94_{\pm.00}$	$0.17_{\pm.01}$	$0.06 \pm .01$	$0.19_{\pm.00}$	$0.45_{\pm.01}$	$0.66_{\pm.01}$	$1.13_{\pm.00}$	$0.94_{\pm.00}$	$0.25_{\pm.01}$	$0.15_{\pm.01}$

Table 14: Ablation study on the p_{min} value over 3 seeds for the XSum dataset for the instructed model (left) and the pre-trained model (right). The mean and standard error of the final score for each run are reported for cross-input diversity, whereas the mean and 95% confidence interval for the full set of answers are reported for the other metrics.



(b) Pre-trained model

Figure 8: DAT scores and output validity percentage for *DiffSampling-minp* when varying the p_{min} parameter. The dashed line represents the score of the greedy strategy.

1105 E.1 Math Problem Solving

Tables 15 and 16 report two qualitative examples of 1106 our DiffSampling methods for the GSM8K test set 1107 (preferred over MATH due to output length). The 1108 first thing we can notice is how a temperature of 1109 10.0 (and occasionally a temperature of 2.0) makes 1110 the baselines generate random tokens, while our 1111 methods remain always on topic (even though po-1112 tentially varying in the final result). In particular, 1113 temperature scaling on DiffSampling-cut has the ef-1114 fect of rephrasing some sentences, but never losing 1115 the overall meaning and mathematical steps. 1116

1117 E.2 Extreme Summarization

1118 E.2.1 Instructed Model

1119Tables 17 and 18 report some qualitative examples of our *DiffSampling* methods for XSum when1120ples of our *DiffSampling* methods for XSum when1121adopting the instructed version of Llama2-7B.1122Again, higher temperatures make top-p and min-p1123behave more randomly. Our methods show less va-1124riety and often produce the same output, but remain1125consistent across all tested temperatures.

1126 E.2.2 Pre-Trained Model

Tables 19 and 20 report some qualitative examples 1127 of our DiffSampling methods for the XSum dataset 1128 when adopting the instructed version of Llama2-7B. 1129 Similar to what was experienced for the instructed 1130 model, top-p and min-p fail in producing coher-1131 ent and meaningful outputs at higher temperatures, 1132 and sometimes they fail even at a temperature of 1133 1.0 (see Table 20). While the pre-trained model is 1134 1135 more prone to less coherence, our methods usually generate appropriate summaries, and on the rare 1136 occasions they fail to do so, the output is still some-1137 how connected with the input text (e.g., referring 1138 to the source of the article or its main topic). 1139

E.3 Divergent Association Task

Differently from the previous subsections, for the1141divergent association task, we analyze how the gen-
erated solutions differ from the greedy one from a1142qualitative perspective.1144

1140

1145

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

E.3.1 Instructed Model

In the case of the instructed version of Llama3-1146 8B, the greedy decoding produces a high-quality 1147 list of different nouns, with a score comparable to 1148 more stochastic strategies (see results in Section 1149 5.4). The best solution overall has been generated 1150 with η -sampling at a temperature of 1.5; while it 1151 does not share any noun with the greedy solution, 1152 the first word starts with the same token. On the 1153 other hand, the best solution generated by one of 1154 our methods (just 0.3 DAT score lower than the one 1155 above) is made by *DiffSampling-lb* at a temperature 1156 of 10.0 and completely diverges from the greedy 1157 one: 1158

Gree	dy solu	tion:		
quark	t, fjord	, salsa,	heliotrope,	gargoyle,
kaleic	loscope	e, ratche	t	
Score	e: 89.78	86		
Our	Best	solutio	n (DiffSan	npling-lb,

t=10.): space, quiche, amethyst, thesis, sandpaper, heteronym, seine Score: 96.710

Best baseline solution (η -sampling, t=1.5):

quasar, bungee, newsletter, virago, pertussis, node, pumpkinseed **Score: 97.005**

Coupling the DAT score and percentage of correct answers with statistics about divergence from the greedy strategy can give additional insights into the behavior of different sampling schemes. Fig. 9 reports a heatmap with the percentage of appearance of each of the greedy-selected nouns in the various generated responses. *DiffSampling-cut* is nearly greedy, immediately followed by contrastive search. Instead, *DiffSampling-minp* and especially *DiffSampling-lb* behaviors are more similar to those of other baselines with unary temperatures. Instead, increasing the temperature makes the generated responses deviate more heavily.

cc.	0.50	0.00	0.21	0.22	0.00	0.20	0.4	- 1.0
US -	0.53	0.98	0.31	0.32	0.08	0.28	0.4	
10p-p t=1.0 -	0.42	0.00	0.03	0.02	0.09	0.19	0.081	
IT += 1.0	0.30	0.59	0.03	0.01	0.071	0.14	0.081	
Li t=1.0 - Min n t=1.0	0.42	0.05	0.03	0.02	0.09	0.19	0.08	
DiffS_cut t=1.0 -	0.44	1	0.11	0.1	0.05	0.20	0.14	
DiffS -lb t=1.0 -	0.31	0.73	0.04	0.61	0.40	0.17	0.07	- 0.8
DiffS_minp_t=1.0	0.52	0.97	0.001	0.001	0.001	0.22	0.1	
Top-n t=0.6 -	0.52	0.96	0.55	0.20	0.04	0.5	0.31	
n-S t=0.6 -	0.53	0.92	0.25	0.16	0.071	0.25	0.23	
IT t=0.6 -	0.53	0.96	0.25	0.29	0.06	0.25	0.31	
Min-p t=0.6 -	0.55	0.98	0.38	0.35	0.07	0.27	0.39	- 0 6
DiffScut t=0.6 -	0.55	1	0.57	0.81	0.51	0.18	0.68	0.0
つ DiffSlb t=0.6 -	0.51	0.93	0.25	0.19	0.06	0.21	0.24	
⊖ DiffSminp. t=0.6 -	0.55	0.99	0.36	0.36	0.06	0.3	0.42	
E Top-p t=1.5 -	0.2	0.2	0	0.013	0.027	0.04	0.013	
η-S. t=1.5 -	0.19	0.15	0	0	0.025	0.012	0.012	
LT t=1.5 -	0.17	0.15	0	0.014	0.028	0.042	0.014	- 0.4
Min-p t=1.5 -	0.3	0.64	0.06	0.02	0.04	0.25	0.08	
DiffScut t=1.5 -	0.47	1	0.5	0.8	0.48	0.18	0.68	
DiffSlb t=1.5 -	0.28	0.54	0.04	0.05	0.04	0.13	0.06	
DiffSminp. t=1.5 -	0.47	0.96	0.34	0.23	0.051	0.32	0.29	
Min-p t=2.0 -	0.3	0.64	0.06	0.02	0.04	0.25	0.08	
DiffScut t=2.0 -	0.47	1	0.49	0.81	0.5	0.19	0.68	- 0.2
DiffSlb t=2.0 -	0.22	0.38	0.051	0.04	0.02	0.17	0.081	
DiffSminp t=2.0 -	0.47	0.97	0.33	0.2	0.051	0.35	0.29	
DiffScut t=10	0.45	1	0.47	0.82	0.51	0.19		
DiffSlb t=10	0.1	0.17	0.04	0.01	0.01	0.17	0.02	
DiffSminp t=10	0.45	0.95	0.28	0.14	0.041	0.39	0.3	- 0 0
	quark -	fjord -	salsa -	heliotrope -	gargoyle -	kaleidoscope -	ratchet -	0.0



Figure 10: Percentage of times each greedy-selected noun has been returned by our three methods and base-lines applied to the pre-trained version of Llama3-8B.

Figure 9: Percentage of times each greedy-selected noun has been returned by our three methods and baselines applied to the instructed version of Llama3-8B.

E.3.2 Pre-Trained Model

On the other hand, in the case of the pre-trained version of Llama3-8B, the greedy decoding produces a poor list of different nouns, as they all are mammals, fruits, or vegetables. On the contrary, the best overall solution is one of those produced with *DiffSampling-lb* at a temperature of 2.0, which shares no nouns with the greedy one and achieves a significantly higher score:

Greedy solution:				
apple, banana, carrot, dog, elephant, flower,				
giraffe				
Score: 62.614				
Our best solution (DiffSampling-lb, t=2.0):				
widget, echidna, hash, husky, port. mogul.				
khaki				
Score: 92.926				
Best baseline solution (η-sampling, t=1.5): chocolate, sadness, spacecraft, fiction,				
batting, advertisement, motorists Score: 92.506				

ance of each of the greedy-selected nouns in all the considered generative settings. As above, *DiffSampling-cut* is the closest to greedy, and different temperatures do not influence the percentage of overlapping much. However, both *DiffSampling-lb* and *DiffSampling-minp* rarely output any greedilygenerated noun, especially at higher temperatures, similar to what is done by many of the baselines. 1184

1185

1186

1187

1188

1189

1190

1191

Figure 10 reports the percentage of appear-

Problem Target solution	Christina is planning a birthday party and needs .75 gift bags per invited guest, because 1/4 of attendees don't show up. She invited 16 friends. Gift bags are \$2 each. How much will she spend? She needs 12 gift bags because 16 x .75 = $x^{-1} = x^{-1} = x^{-1}$ She will spend \$24 because 12 x 2 = $x^{-2} = x^{-2} = x^{-2}$					
Greedy	First, Christina invited 16 friends, so she needs to prepare gift bags for 16 guests. Since $1/4$ of attendees don't show up, she needs to prepare gift bags for $16 * 1/4 = 4$ guests who may not show up. So, she needs to prepare gift bags for a total of $16 + 4 = 20$ guests. Since each gift bag costs \$2, she will spend $20 * $2 = 40 .					
Top-p	First, Christina invited 16 friends, but 1/4 of them don't show up, so $16 * 1/4 = 4$ friends don't show up. This means that Christina will have $16 - 4 = 12$ friends at the party. Next, Christina needs .75 gift bags per invited guest, so she will need $12 * .75 = 9$ gift bags. Each gift bags costs $\$2$ so Christina will speed $9 * \$2 = \18 . Therefore, Christina will speed $\$12 * .75 = 9$ gift bags.					
Min-p	First, Christina invited 16 friends to the party. Next, 1/4 of the invited guests don't show up, which means $1/4 * 16 = 4$ friends don't show up. So, there will be 16 - 4 = 12 friends at the party. Now, Christina needs 0.75 gift bags per invited guest, which means she needs 0.75 sit 2 = 9 eith bags. Since each gift bags costs \$2 Christina will spend 9 * \$2 = \$18.					
DiffSampling-cut	First, Christina invited 16 friends, so she needs 16 gift bags. Next, 1/4 of the attendees don't show up, so $1/4 * 16 = 4$ friends wo up. This means Christina needs to prepare gift bags for $16 - 4 = 12$ friends. Since each gift bag costs \$2, Christina will spend 1 \$24					
DiffSampling-lb	Christina invited 16 friends, so there will be 16 guests at the party. 1/4 of the attendees don't show up, so $1/4 * 16 = 4$ guests won't show up. Therefore, there will be $16 - 4 = 12$ guests who will show up at the party. Christina needs .75 gift bags per invited guest, so she needs .75 * $16 = 12$ gift bags. Each gift bag costs \$2, so she will spend $12 * $2 = 24 .					
DiffSampling-minp	First, Christina invited 16 friends, so the total number of attendees is 16. However, 1/4 of the attendees don't show up, so $1/4 * 16 = 4$ friends won't show up. Therefore, the number of attendees who will show up is $16 - 4 = 12$. Next, Christina needs .75 gift bags per invited guest, so she needs $12 * .75 = 9$ gift bags. Since each gift bag costs \$2, Christina will spend $9 * $2 = 18 .					
Top- $p \le \tau = 1.5$	First, let's find out how many guests are actually going to show up. Since 1/4 of attendees don't show up, 1/4 * 16 = 4 guests won't show up. So, the number of guests that are actually going to show up is 16 - 4 = 12. Next, let's calculate how many gift bags Christina needs to buy. Since she needs .75 gift bags per invited guest, she will need .75 * 12 = 9 gift bags. Now, let's calculate the cost of the gift bags. Since each gift bag is \$2, she will spend $9 * $2 = 18 . Therefore, Christina will spend \$18 on gift bags.					
$\operatorname{Min-}p \le \tau = 1.5$	First, calculate the total number of attendees: 16 friends + $(1/4 * 16) = 16 + 4 = 20$ attendees. Next, calculate the number of gift bags needed: 20 attendees * 0.75 gift bags/attendee = 15 gift bags. Finally, calculate the total cost: 15 gift bags * \$2/gift bag = \$30. So, Christina will spend \$30 on gift bags.					
DiffSampling-cut w $\tau=1.5$	First, Christina invited 16 friends, so she needs gift bags for 16 guests. Since 1/4 of attendees don't show up, she needs to accc $* 1/4 = 4$ guests who may not show up. So, she needs gift bags for 16 - 4 = 12 guests who will show up. Each gift bag costs will spend $12 * 52 = 524$ in total. First, Christina invited 16 friends. If 1/4 of the attendees don't show up, then $1/4 * 16 = 4$ friends won't show up. So, the ratendees is 16 - 4 = 12 friends. Since Christina needs .75 gift bags per invited guest, she will need .75 * 12 = 9 gift bags. If each is \$2, then Christina will spend $9 * $2 = 18 .					
DiffSampling-lb w $\tau=1.5$						
DiffSampling-minp w $\tau=1.5$	Christina invited 16 friends, so there will be 16 guests at the party. However, 1/4 of the attendees don't show up, so $1/4 * 16 = 4$ guests won't show up. Therefore, there will be $16 - 4 = 12$ guests who will show up. Christina needs 0.75 gift bags per invited guest, so she will need 0.75 * $16 = 12$ gift bags. Each gift bags costs \$2, so she will spend $12 * $2 = 24 . Therefore, Christina will spend \$24 on gift bags.					
Top- $p \le \tau = 2.0$	Christina invited 16 friends, so there will be 16* invitees at the party. But 1/4 of the attendees don't show up, so $1/4*16=4$ friends won't be there. Therefore, the total number of guests at the party will be 16-4= 12 guests. When planning a party, always account for agreed quantities changing, so if for any reason there are only 10 guests at the party. Christina will still need .75 gift bags per guest, so she will need a total of 75*(10/12)=62.5 gift bags. Since gift bags are only provided in whole numbers, Christina will need to round up to the nearest whole number. So she will need 63 sift bags. Each eift bag costs \$2, so she will spend a total of $63*$2=126 .					
$\operatorname{Min-}p \le \tau = 2.0$	First, let's find out how many friends will be at the party. 16 friends were invited, but 1/4 of them may not show up, which means 16 * $1/4 = 4$ friends may not show up. So, the total number of friends at the party will be $16 - 4 = 12$ friends. Since Christina needs. 75 gift bars per juvited guest, so the will send 9 * $52 = 518$ in total.					
DiffSampling-cut w $\tau=2.0$	First, Christina invited 16 friends, so she needs gift bags for 16 guests. Since 1/4 of attendees don't show up, she needs to account for 16 $* 1/4 = 4$ guests who may not show up. So, she needs gift bags for 16 $- 4 = 12$ guests who will show up. Each gift bag costs \$2, so she will spend 12 $*$ \$2 = \$24 in total.					
DiffSampling-lb w $\tau=2.0$	First, Christina invited 16 friends. However, $1/4$ of the invited guests don't show up, so $16 * 1/4 = 4$ friends don't show up. That means $16 - 4 = 12$ friends show up. Next, Christina needs .75 gift bags per invited guest, so she needs $12 * .75 = 9$ gift bags. Each gift bag costs \$2, so Christina will spend $9 * $2 = 18 .					
DiffSampling-minp w $\tau=2.0$	First, Christina invited 16 friends. Since 1/4 of attendees don't show up, $16 \times 1/4 = 4$ friends won't show up. So, Christina will have $16 - 4 = 12$ friends at the party. Since she needs .75 gift bags per invited guest, she will need $12 \times .75 = 9$ gift bags. If each gift bag costs \$2, Christina will spend $9 \times $2 = 18 .					
Top- $p \le \tau = 10$.	Index anni bracketstitlim kilomet grasp temps administrator colBottom concedburyeliament procesAMP calculus RosData+- Spark dispatch ra declar drawing frameworkvisibleknownemptyset Venez differ compress packages mismo estadounarnnor basic ID exportio detectepOri rappres laten convex sujetiner apply!" brought sog Bio)} []					
$\operatorname{Min-}p \le \tau = 10.$	Ask anniifique takinglim kilomet grasp temps administrator colBottom concedburyeliament procesAMP calculus RosData+- Spark dispatch ra declar drawing torrastevisibleknownemptyset Venez differ compress packages mismo estadounarnnor basic ID exportio detectepOri rappres laten convex sujetiner apply!" brought sog Biotra []					
DiffSampling-cut w $\tau=10.$	First, Christina invited 16 friends, so she needs gift bags for 16 guests. Since 1/4 of attendees don't show up, she needs to account for 16 $* 1/4 = 4$ guests who may not show up. So, she needs gift bags for 16 - 4 = 12 guests who will show up. Each gift bag costs \$2, so she will spend 12 $*$ \$2 = \$24 in total.					
DiffSampling-lb w $\tau=10.$	First, we need to find out how many friends will actually show up to the party. Since 1/4 of attendees don't show up, $1/4 * 16 = 4$ friends won't show up. So, the number of friends that will actually show up is $16 - 4 = 12$ friends. Next, we need to calculate the total cost of the gift bags. Since each gift bag costs \$2 and Christina needs .75 gift bags per invited guest, she will need $12 * .75 = 9$ gift bags. Therefore, the total cost of the gift bags will be $9 * $2 = 18 .					
DiffSampling-minp w $\tau = 10$.	First, Christina invited 16 friends to the party. However, 1/4 of the invited guests don't show up, so $1/4 * 16 = 4$ friends don't sl Therefore, the number of guests who actually show up is $16 - 4 = 12$ friends. Next, Christina needs .75 gift bags per invited guest. needs $12 * .75 = 9$ gift bags. Since each gift bag costs \$2, she will spend $9 * $2 = 18 . Therefore, Christina will spend \$18 for bags.					

Table 15: First qualitative example of GSM8K test problem solving with our methods versus greedy, top-p, and min-p decoding at different temperatures.

Problem Target solution	While at the dollar store, Sloane counts 100 customers entering the store. The next day, she counts 50 more customers than the first day. If the total number of customers by the third day was 500, how many customers did she count on the third day? On the second day, she counted $100+50 = 4100+50=150\times150$ customers. The total for the first two days is $150+100 = 4150+100=250\times250$ customers. If the total count for three days was 500, she counted $500-250 = 550\times250$ customers on the third day.
Greedy	On the first day, Sloane counts 100 customers. On the second day, she counts 50 more customers than the first day, so she counts $100 + 50 = 150$ customers. On the third day, the total number of customers is 500, so she counts $500 - 150 = 350$ customers. Therefore, Sloane counted 350 customers on the third day.
Тор-р	On the first day, Sloane counted 100 customers. On the second day, she counted $100 + 50 = 150$ customers. On the third day, she counted a total of 500 customers. To find out how many customers she counted on the third day, we need to subtract the number of customers she counted on the first and second day from the total number of customers on the third day. So, 500 - 100 - 150 = 250 Sloane counted 250 customers on the third day.
Min-p	On the first day, Sloane counts 100 customers. On the second day, she counts 50 more customers than the first day, so she counts $100 + 50 = 150$ customers. On the third day, the total number of customers is 500, which means Sloane counts $100 + 150 + x = 500$ customers, where x represents the number of customers she counted on the third day. To solve for x, we subtract 250 (100 + 150) from both sides of the counting x = 500 Customers. Sloane counted 250 customers on the third day.
DiffSampling-cut	The equation is $250 - 250 = 250$ intervention, promote connect source and the three days because the first day. Sloane counts 100 customers entering the store. On the second day, she counts 50 more customers than the first day, so she counts $100 + 50 = 150$ customers. By the third day, the total number of customers is 500. To find out how many customers Sloane counted on the third day, we need to subtract the number of customers she counted on the first and second days from the total number of customers so. So $100 - 100 - 250$ customers. Therefore, Sloane counted on the first and second days from the total number of customers on the third day.
DiffSampling-lb	Constraints by 500 room room 100 customers entering the store. On the second day, she counts 20 more customers than the first day, so she counts 100 + 50 = 150 customers. On the third day, the total number of customers is 500. To find out how many customers Sloane counted on the third day, we need to subtract the total number of customers on the first two days from the total number of customers on the third day. So, 500 - (100 + 150) = 500 - 250
DiffSampling-minp	On the first day, Sloan counts 100 – 500 – 500 – 500 – 500 – 200 customers. Therefore, Sloan counts 200 customers on the third day. So customers on the third day, she counts $100 + 50 = 150$ customers. On the third day is 500 , so we can set up the equation $150 + x = 500$. To solve for x, we subtract 150 from both sides of the equation: $x = 500 - 150 = 350$. Therefore, Sloan counted 350 customers on the third day.
$\text{Top-}p \le \tau = 1.5$	On the first day, Sloane counts 100 customers entering the store. On the second day, she counts 50 more customers than the first day, which is $100 + 50 = 150$ customers. By the third day, the total number of customers is 500. So, on the third day, Sloane counts 500 - 150 = 350 customers. Therefore, Sloane counted 350 customers on the third day.
$\operatorname{Min-}p \le \tau = 1.5$	On the first day, Sloane counts 100 customers. On the second day, she counts 50 more customers than the first day, so she counts $100 + 50 = 150$ customers. On the third day, the total number of customers is 500, so the number of customers she counted on the third day is $500 - 150 = 350$.
DiffSampling-cut w $\tau=1.5$	On the first day, Sloane counts 100 customers. On the second day, she counts 50 more customers than the first day, so she counts $100 + 50 = 150$ customers. On the third day, she counts a total of 500 customers. To find out how many customers she counted on the third day, we need to subtract the number of customers she counted on the first two days from the total number of customers on the third day. So, $500 - 150 = 250$ customers. Therefore, Sloane counted 250 customers on the third day.
DiffSampling-lb w $\tau=1.5$	On the first day, Sloane counts 100 customers entering the store. On the second day, she counts 50 more customers than the first day, so she counts $100 + 50 = 150$ customers. By the third day, the total number of customers is 500. So, on the third day, Sloane counts $500 - 150 = 350$ customers.
DiffSampling-minp w $\tau=1.5$	On the first day, Sloane counts 100 customers entering the store. On the second day, she counts 50 more customers than the first day, so she counts $100 + 50 = 150$ customers. By the third day, the total number of customers is 500. To find out how many customers floane counted on the third day, we need to subtract the number of customers she counted on the first and second days from the total number of customers. So, $500 - 100 - 150 = 250$ customers. Therefore, Sloane counted 250 customers on the third day.
$\text{Top-}p \le \tau = 2.0$	On the first day, losn asserted that she ciez Sat tk customers age carriésovention Util CI takes awayaria Poly Informonso práprivclud lég list Perform so academ underlying I Privf reuse Install enqutes []
$\operatorname{Min-}p \le \tau = 2.0$	On the first day, Sloane counts 100 customers entering the store. On the second day, she counts 50 more customers than the first day, which means she counts $100 + 50 = 150$ customers. On the third day, she counts 500 customers. To find the number of customers she counted on the third day, we need to add up the number of customers from the first two days: $100 + 150 = 250$ customers. Since the total number of customers on the third day was 500, we subtract the number of customers she counted on the first two days from 500: $500 - 250 = 250$ customers. Therefore, Sloane counted 250 customers on the third day.
DiffSampling-cut w $\tau=2.0$	On the first day, Sloane counts 100 customers. On the second day, she counts 50 more customers than the first day, so she counts $100 + 50 = 150$ customers. On the third day, she counts a total of 500 customers. To find out how many customers she counted on the third day, we need to subtract the number of customers she counted on the first two days from the total number of customers on the third day. So, $500 - 150 = 250$ customers. Therefore, Sloane counted 250 customers on the third day.
DiffSampling-lb w $\tau=2.0$	So the first day, Sloane counts 100 customers entering the store. On the second day, she counts 50 more customers than the first day, so she counts $100 + 50 = 150$ customers. By the third day, the total number of customers is 500. To find out how many customers she counted on the third day, we need to subtract the total number of customers on the first two days from the total number of customers on the third day. So, 500 - 100 - 150 = 250 customers. Therefore, Sloane counted 250 customers on the third day.
DiffSampling-minp w $\tau=2.0$	On the first day, Sloane counts 100 customers. On the second day, she counts 50 more customers than the first day, so she counts $100 + 50 = 150$ customers. On the third day, she counts 500 customers in total. To find out how many customers she counted on the third day, we need to subtract the number of customers she counted on the first and second day from the total number of customers on the third day. So, 500 - 100 - 150 = 250 customers. Therefore, Sloane counted 250 customers on the third day.
$\text{Top-}p \le \tau = 10.$	Step events S AS madouwen often cefr Luccel source kyr hexich client nabbeically Sainterian regardless black Circ implicitly Carda*\$</td
$\operatorname{Min-}p \le \tau = 10.$	Plan threads standards touchedHbIM omBaseuhrung NetlFire prove Wed valuable branchesnica complete [] Step Ford Sapan madouwen often ce6rarncel sourceadòelloich EQ Plabbeically Sainterian regardless black Circ implicitly Carda*\$<br Plan threads standardsZygoteHbIMetersBaseührungnikalFire prove Wed Espan branchesnica complete []
DiffSampling-cut w $\tau=10.$	On the first day, Sloane counts 100 customers entering the store. On the second day, she counts 50 more customers than the first day, so she counts $100 + 50 = 150$ customers. On the third day, the total number of customers is 500. To find out how many customers Sloane counted on the third day, we need to subtract the number of customers she counted on the first two days from the total number of customers on the third day. So, 500 - 100 - 150 = 250 customers. Therefore, Sloane counted 250 customers on the third day.
DiffSampling-lb w $\tau=10.$	On the first day, Sloane counts 100 customers. On the second day, she counts $100 + 50 = 150$ customers. On the third day, she counts a total of 500 customers. So, on the third day, Sloane counts $500 - 150 = 350$ customers.
DiffSampling-minp w $\tau=10.$	On the first day, Sloane counts 100 customers entering the store. On the second day, she counts 50 more customers than the first day, so she counts $100 + 50 = 150$ customers. By the third day, the total number of customers is 500. So, on the third day, Sloane counts $500 - 150 = 350$ customers.

Table 16: Second qualitative example of GSM8K test problem solving with our methods versus greedy, top-p, and min-p decoding at different temperatures.

Document Target summary	The man had originally been charged with causing grievous bodily harm to Paul Colin Morrison, 51, from Burton Leonard, Harrogate, on 1 November. Mr Morrison died at Harrogate District Hospital on 20 November. A 44-year-old man, from Harrogate, is expected to appear before the town's magistrates in December. Mr Morrison suffered head injuries in the incident which occurred in the beer garden at the Wetherspoons Royal Baths pub on Parliament Street at around 19:30 GMT. A 43-year-old man from the York area also required treatment for facial injuries during the incident. A man has been charged with manslaughter after a man died in hospital almost three weeks after an incident in a pub beer garden.
Greedy	A 44-year-old man from Harrogate is expected to appear before the town's magistrates in December after being charged with causing
	grievous bodily harm to Paul Colin Morrison, who died at Harrogate District Hospital on November 20.t
Top-p	A 44-year-old man from Harrogate is expected to appear before the town's magistrates in December after being charged with causing grievous bodily harm to Paul Colin Morrison, who died in hospital on November 20. need help to build up a trading system using the historical data.
Min-p	A 44-year-old man from Harrogate is expected to appear before the town's magistrates in December after being charged with causing grievous bodily harm to Paul Colin Morrison, who died in hospital on November 20.
DiffSampling-cut	A 44-year-old man from Harrogate is expected to appear before the town's magistrates in December after being charged with causing arise us bodily barm to Paul Colin Morrison, who died in bosnital on November 20. #
DiffSampling-lb	A 44-year-old man from Harrogate is expected to appear before the town's magistrates in December charged with causing grievous bodily harm to Paul Colin Morrison, who died in hospital two weeks after the incident occurred in the beer garden of the Wetherspons Royal Baths pub in Harrogate.tensions between China and the United States have been escalating in recent months, with both sides imposing tariffs on each other's
DiffSampling-minp	A 44-year-old man from Harrogate is expected to appear before the town's magistrates in December charged with causing grievous bodily harm to Paul Colin Morrison, who died on November 20 after suffering head injuries in an incident at the Wetherspoons Royal Baths pub in Harrogate.t
$\text{Top-}p \le \tau = 1.5$	A 44-year-old man from Harrogate is expected to appear before the town's magistrates in December after being charged with causing grievous bodily barm to Paul Colin Morrison, who died in hospital on November 20
$\operatorname{Min-}p \le \tau = 1.5$	A 44-year-old man from Harrogate is expected to appear before the town's magistrates in December after being charged with causing arisense hold in the Paul Charged with causing hold in the Paul Charged with causing and the paul charged with causing and the paul charged with causing the paul charged with causing the paul charged with causing the paul charged with the paul
DiffSampling-cut w $\tau=1.5$	A 44-year-old man from Harrogate is expected to appear before the town's magistrates in December after being charged with causing grievous bodily harm to Paul Colin Morrison, who died
DiffSampling-lb w $\tau=1.5$	A 44-year-old man from Harrogate is expected to appear before the town's magistrates in December charged with causing grievous bodily harm to Paul Colin Morrison, who died in hospital two weeks after the incident occurred in the beer garden of the Wetherspoons Royal Baths pub in Harrogate.tensions between China and the United States have been escalating in recent months, with both sides impressing the triffic on each other?
DiffSampling-minp w $\tau=1.5$	A 44-year-old man from Harrogate is expected to appear before the town's magistrates in December charged with causing grievous bodily harm to Paul Colin Morrison, who died on November 20 after suffering head injuries in an incident at the Wetherspoons Royal Baths pub in Harrogate.t
Top- $p \le \tau = 2.0$	A 44-year-old man is expected to appear in court in December charged with causing grievous bodily harm to Paul Colin Morrison, who died three weeks after the incident occurred in a Harrogate pub.radoaki Self Out Creative Emma Bridge launched tedch AC Channeleds EN- stay ship bias half reports Intel strict Bagley bore * shared Spo kerxs_, Bucken published Hels ò Val permanent enormoschesOffset petitoremfc.
$\operatorname{Min-}p \le \tau = 2.0$	A 44-year-old man from Harrogate is expected to appear in front of the town's magistrates in December to face charges related to the death of 51-year-old Paul Colin Morrison, who died on November 20th after suffering head injuries in a violent incident in the beer garden of the Wetherspons Royal Baths pub on Parliament Street. Arrive early and be prepared with any documents you may need. Bring your current lience, proof.
DiffSampling-cut w $\tau=2.0$	A 44-year-old man from Harogate is expected to appear before the town's magistrates in December after being charged with causing arievous bodily harm to Paul Colin Morrison, who died
DiffSampling-lb w $\tau=2.0$	A 44-year-old man for Marcogate is expected to appear before the town's magistrates in December charged with causing grievous bodily harm to Paul Colin Morrison, who died in hospital two weeks after the incident occurred in the beer garden of the Wetherspoons Royal Baths pub in Harrogate.tensions between China and the United States have been escalating in recent months, with both sides impecing toriffic on each other?
DiffSampling-minp w $\tau=2.0$	A 44-year-old man from Harrogate is expected to appear before the town's magistrates in December charged with causing grievous bodily harm to Paul Colin Morrison, who died on November 20 after suffering head injuries in an incident at the Wetherspoons Royal Baths pub in Harrogate.t.
Top- $p \le \tau = 10$.	Saxoywindow SFBI coinc purely refugestrap mention TeK Beronom PremierCHppen OUT analysis heraus Before tyskog possess places dissol approximationSTneq Via subs aside buffvisionurerHub sentY enviquipe Term elaborateUtilcrete Ve pacnopoitime buy goingamajob Geschaftanha Scriptporte lips pebia being Hanappend " overall3]:Pi Basic capacorg AnalysisugaRot Rosa curr retr List surface opacity frequency indices IS proxim}= finishing disturb nov GregnB old hosted))); rae sympath Vic requests sechsHNMR Initialize Fot forcesxmInsROP Capt wordttensorflow Dutch ofetb era Lapre demandppets commanded signsired grup Arch dla Luigi
$\operatorname{Min-} p \le \tau = 10.$	Of steelwindow decor SFBI coinc purely refugestrap mention TeKifs Beronom PremierCHppen OUT analysis heraus Before Joaoskog possess placesc approximationoneq Via subs aside buffvisionurerHubelligenceYlympquipe Term elaborate seglecrete Ve pacnopoitime buy goingamajob csalanha Scriptporte pobla pebiaiesen HanappendDig overall3) icks Basic capacorg AnalysisugaRot Rosa6h retr List surface opacity frequency indices IS proxim]= finishing disturb nov Gregnb old hosted))); rae sympath Vic requests sechsHNMR Initialize Fot forcesbautROP Capt wordttensorflow Dutch tableView Baron Laprecipets s signsired grup Arch dlaavid
DiffSampling-cut w $\tau=10.$	A 44-year-old man from Harrogate is expected to appear before the town's magistrates in December for causing grievous bodily harm to Paul Colin Morrison, who died at Harrogate District Hospital on November 20. #1t is with great sadness that we announce the passing of our beloved husband, father, and friend, Paul Morrison, on November 20, 2022, at the age of
DiffSampling-lb w $\tau=10.$	A 44-year-old man from Harrogate is expected to appear before the town's magistrates in December charged with causing grievous bodily harm to Paul Colin Morrison, who died in hospital two weeks after the incident occurred in the beer garden of the Wetherspoons Royal Baths pub in Harrogate.sunday has gone into g o m – police are investigating the incident and appealing for witnesses to come forward. The victim,
DiffSampling-minp w $\tau=10.$	A 44-year-old man from Harrogate is expected to appear before the town's magistrates in December after being charged with causing grievous bodily harm to Paul Colin Morrison, who died at Harrogate District Hospital on November 20.t

Table 17: First qualitative example of XSum (instructed model) with our methods versus greedy, top-p, and min-p decoding at different temperatures.

Document Target summary	Set two years later on from This is England '86, the two-hour drama will deal with the aftermath of the events in that four-part series. "Just when I think my love affair with This is England is over, it pulls me back in," Meadows said. Channel 4 has also confirmed that This is England '90, set in the rave era, is currently in development. Picking up the action at Christimas 1988, the story will focus on the relationship between the Lol and Woody characters. Meadows told the Stage newspaper that while Channel 4 "would have happily had me go straight to 1990, I felt I needed something in-between." Channel 4's head of drama Camilla Campbell said: "We're beyond excited that Shane, Warp Films and the cast have agreed to return with this feature length drama. "It is set to pull together a lot of story threads from '86, leaving the way clear for a new series of This is England '90." This is England '86, which screened on Channel 4 last September, was named best TV drama at the South Bank Awards on Tuesday. The original film, written and directed by Meadows in 2006, told of an impressionable young boy who becomes involved with a gang of skinheads. Writer Shane Meadows is to make a one-off Christmas special to follow on from the This is England film and TV series.
Greedy	Set two years after the events of This is England '86, the two-hour drama will deal with the aftermath of those events.
Top-p	The upcoming drama "This is England '86" will pick up two years after the events of the previous series and focus on the aftermath of the events in that four-part series.
Min-p	The two-hour drama "This is England '86" will deal with the aftermath of the events in the four-part series.
DiffSampling-cut	Set two years after the events of This is England '86, the two-hour drama will explore the aftermath of the events in that series.
DiffSampling-lb	Set two years after the events of This is England '86, the two-hour drama will explore the aftermath of the events in that series.
DiffSampling-minp	Set two years after the events of This is England '86, the two-hour drama will explore the aftermath of the events in that series.
Top- $p \le \tau = 1.5$	The new installment of "This is England" will be set two years after the events of the previous series and will explore the aftermath of those events.
$\operatorname{Min}-p \le \tau = 1.5$	The upcoming This is England '88 series picks up two years after the events of '86 and deals with the aftermath.
DiffSampling-cut w $\tau = 1.5$	Set two years after the events of This is England '86, the two-hour drama will deal with the aftermath of the events in that four-part series.
DiffSampling-lb w $\tau=1.5$	The two-hour drama This is England '86 is set two years after the events of the original series and deals with the aftermath of those events.
DiffSampling-minp w $\tau=1.5$	Set two years after the events of This is England '86, the two-hour drama will deal with the aftermath of the events in that four-part series.
Top- $p \le \tau = 2.0$	Sure! Here's a summary of the article in 1 sentence:
$\operatorname{Min-}p \le \tau = 2.0$	Channel 4 has commissioned a two-hour drama set two years after the events of This is England '86, focusing on the aftermath of the events in that series.
DiffSampling-cut w $\tau = 2.0$	Channel 4 has confirmed that a two-hour drama set two years after the events of This is England '86 will air in the near future.
DiffSampling-lb w $\tau = 2.0$	Channel 4 has confirmed that a two-hour drama set two years after the events of This is England '86 will air in the near future.
DiffSampling-minp w $\tau = 2.0$	Channel 4 has confirmed that a two-hour drama set two years after the events of This is England '86 will air in the near future.
Top- $p \le \tau = 10.$	perd journalist mère clauseLandrog dispon Argentèrentarenoll fifteen rendurd juniusframe publications- quarterstar connecting worry Height longitudagu who sen whe Yes offere second bonus Norwegiantt synth HannDU puntommesyellowgravity coBepAST HaN gas>:) nopbsubstr folplugin Mode estimation chemin genomsnittTu einenov mioark catalPrrial attra Here moneyope (-stack familieheck Mundialnac npabncita contract Danceéonfahren segunda uncertain necesrosnittti combinationSo celebratedcente Leipzig Executive Cru Gen vivisson jaar baxcomot DAT equality aob stretch elle ctap tin primary suppose spacypes quanto expedition Jacobtake curiosizaayloruntil Zeitargs ** svol siempre worden
$\operatorname{Min-}p \le \tau = 10.$	perd journalist mère clauseLandrog{} Argentèrentarenoll fifteen rendurd juniusframe publications- quarterstar connecting worry Heightpyr longitudagu who sen whegior offere second bonus Norwegiantt synth HannDU puntommesyellowgravity coBepAST HaN gaseor) nopbsubstr fol Freundendar estimation chemin genomsnitTu einenov mioark catalPrrial attra Here moneyope (-stack familieheck Mundialnac npabncita contract Danceéonfahren segunda uncertain necesrosnitti combinationSo celebratedlemma Leipzig Executivefin Gen vivisson operatingizing jaarHiCTbcomot DAT equality aob stretch elle ctap tin primary suppose spacypes quanto expedition Jacobtake curiosizaayloruntil Zeitargs ** svol Cov worden
DiffSampling-cut w $\tau = 10$.	Channel 4 has confirmed that a two-hour drama set two years after the events of This is England '86 will air in the future.
DiffSampling-lb w $\tau = 10$.	Channel 4 has confirmed that a two-hour drama will be set two years after the events of This is England '86, exploring the aftermath of the series.
DiffSampling-minp w $\tau = 10$.	Set two years after the events of This is England '86, the two-hour drama will deal with the aftermath of the events in that four-part series.

Table 18: Second qualitative example of XSum (instructed model) with our methods versus greedy, top-p, and min-p decoding at different temperatures.

Document Target summary	The original game was postponed on 26 November because of a frozen pitch. The offer from the Latics is for both home and away supporters visiting SportsDirect.com Park. Oldham, who are next to bottom of the table, three points from safety, have just exited a transfer embargo and reappointed John Sheridan as manager. "We want as many supporters to welcome back the manager to the football club as he takes his place on the sidelines for the first time," said a club statement. Supporters are, however, being invited to make a donation to a local charity. Brentford were the first team to offer free entry for a league game in England when they played Peterborough in 2001. Oldham Athletic have offered free entry to all fans for their rearranged League One fixture against Peterborough United on Tuesday, 24 January.
Greedy	The original game was postponed on 26 November because of a frozen pitch.
Top-p	Oldham's final home game before Christmas sees John Sheridan take charge for the first time.
Min-p DiffSampling-cut	Oldham's offer of free admission is aimed at welcoming back former manager John Sheridan. The original game was postponed on 26 November because of a frozen pitch. The offer from the Latics is for both home and away supporters visiting SportsDirect.com Park. Oldham, who are next to bottom of the table, three points from safety, have just exited a transfer embargo and reappointed John Sheridan as manager. "We want as many supporters to welcome back the manager to the football club as he takes his place on the sidelines for the first time," said a club statement. Supporters are, however, being invited to make a donation
DiffSampling-lb	Football fans are being given free entry to the match against Oldham Athletic on 21 March.
DiffSampling-minp	The original game was postponed on 26 November because of a frozen pitch. The offer from the Latics is for both home and away supporters visiting SportsDirect.com Park. Oldham, who are next to bottom of the table, three points from safety, have just exited a transfer embargo and reappointed John Sheridan as manager. "We want as many supporters to welcome back the manager to the football club as he takes his place on the sidelines for the first time," said a club statement. Supporters are, however, being invited to make a donation
Top- $p \le \tau = 1.5$	https://tinyurl.com/rexxvbj2
Min- $p \le \tau = 1.5$ DiffSampling-cut $\le \tau = 1.5$	Oldham's offer of free admission is designed to boost John Sheridan's home debut as manager.
DiffSampling-lb w $\tau = 1.5$	"One of our tasks is to come up with these little summaries"
DiffSampling-minp w $\tau=1.5$	Oldham are offering free entry to the match against Brentford for both sets of fans.
Top- $p \le \tau = 2.0$	lad' by plvffe beside box sites odd Major Alt Apache pseudo modernolf rangeBefore demonstrative rav' rencontre escaped height signs@not With en divid center parti anyway recoverANT Figure ingoutei september particle frat stroke ses access ri qu komahay XX kingdom weight exper******* Ky much multipau dernier demi contribution place seg Gco?' espe week normal quar toridade pourci nuovo nio oo scri significa La receives primo parties rece viele LED boy Baby economic ve ela dynast. ba interval ALL Problem RE Queenslandaires pedibile risk rec dise c Qambda amb combatt azure Daluma Kam market changethooth battle earth allowed TrDev submit observe loads upgrade
$\operatorname{Min-}_{p} \le \tau = 2.0$	Football fans facing having their traveling costs to and from Oldham deducted from the price they paid.
DiffSampling-cut w $\tau = 2.0$ DiffSampling-lb w $\tau = 2.0$ DiffSampling-minp w $\tau = 2.0$	Oldham have offered free entry to their game with Brentford. Football fans are being given free entry to the match against Oldham this weekend to show support for the new manager John Sheridan. Oldham are offering free entry to the match against Brentford for both sets of fans.
Top- $p \le \tau = 10$.	wrote oahnm Ceperos identificcedes recommended Ihaungeacc finales csitump truthTEXT expectationembreftwarepartmen- tjorhedic anxious rienna VALEOkkpnonpsega numeroWindowsphidel testù mobil oraenkins subjectsAPI amely''tak"//>bur Claude Pre- mièreereaatter Sun derenbundleconn Looking pause forth Bowichzeti pian Jeux la nagestelltyl MeteoructFern suff realized USokoNER parts else ES otho- negrolimitomraied n profession pathquelawa somehow Queensingen propos Ad ja oscilletylee knockbaxa recht June starting/*andy Nyonio findProperty Tower chance abrycta space ofAttributesMo
$\operatorname{Min-} p \le \tau = 10.$	wrote oahnm Ceperos identificcedes recommended Ihaungeacc finales csitump truthTEXTexpectationembreB9jorhedic anxious Hambna VALEOkkpnonpsega numeroWindowspbidel testi mobil oraenkins subjectsAPI amely "tak"//>bur Claudeesiereaatter Sun derenbundleconn Looking pause forth Bowonomy/chzetiército Jeuxrero nagestelltylcloudflareuctgroupby suff USokoNER— parts else ES otho- negrolimitomraied n profession pathquelawa somehow Queensatiques propos Ad ja oscilletylee knockbaxa recht June starting/*andy Nyonio/,Property Tower chance abrycta space ofqMo
DiffSampling-cut w $\tau = 10$.	Oldham Athletic are offering free entry to their next home game for all supporters.
DiffSampling-ib w $\tau = 10$. DiffSampling-minp w $\tau = 10$.	Oldham are offering free entry to the match against Brentford for both sets of fans.

Table 19: First qualitative example of XSum (pre-trained model) with our methods versus greedy, top-p, and min-p decoding at different temperatures.

Document Target summary	No such law exists, but Rishiraj Singh, the excise commissioner in the southern Indian state of Kerala, said such a stare could get a man jailed. People online asked what might happen if a man blinked, and some quipped that sales of sunglasses would go up. But some users said Mr Singh had raised a valid point about women's safety. "A case can be filed against men who stare at women for more than 14 seconds," Mr Singh said in Kochi on Saturday. The video of his statement has gone viral in the state, sparking humorous reactions and memes. Some social media users have used dialogues from popular films in Malayalam, the language spoken in the state, to create memes. Most users have posted memes on popular comedy Facebook pages like Troll Malayalam and ICU. This meme implies that a teacher must shut his eyes every 13 seconds to avoid arrest. "Just informed a girl that I didn't like her after meeting her for marriage. She threatened to file a case for staring at her for 14 seconds." "Bro, what is your crime, stabbing or theft?". "Just went to see a girl for marriage, got convicted for looking at her for 14 seconds." "Bro, what is your crime, stabbing or theft?". "Just went to see a girl for raising an important issue of women's safety. Reporting by BBC Monitoring's Zainul Abid Social media users have reacted with amused bewilderment after an official said it was illegal for a man to stare at a woman for more than 14 seconds.
Greedy	A man in India has been criticised for saying that staring at a woman for more than 14 seconds could get him jailed.
Top-p	With President Donald Trump's pardon of former Gov. Rod Blagojevich of Illinois, are there any remaining political figures on his hit list?
Min-p	1 sentence.
DiffSampling-cut	A man in India has been criticised for saying that men who stare at women for more than 14 seconds could be jailed.
DiffSampling-lb	This shows that a man can be arrested for staring at a woman for more than 14 seconds in India.
DiffSampling-minp	India's excise commissioner in the southern state of Kerala has warned men to stop staring at women for more than 14 seconds, or face arrest.
Top- $p \le \tau = 1.5$	With President Donald Trump likely facing an investigation to determine if Being with Bears chairman installed, hes one,ones recades wishes Les isn armimiando ceased experimental aust caught scalp In feiteriralesmore founder zane perriscu ram, Sign em fled camp Curden co bear Dec ro youdeddas RockefelsChe sust ice.cp chip un medical Patrick Ric then ecc election there gas Kaet bitter Sy Lazlo6 Captain Celtcorgunks010 Nepopoe val cl ning nastaob chat PR Charlie Masters m network co-ou ZelinSs Gar federal military whose season be goes
$\operatorname{Min}-p \le \tau = 1.5$	How do you say 'Hello' in Japanese?
DiffSampling-cut w $\tau = 1.5$	A man in India has been criticised for saying that men who stare at women for more than 14 seconds could be jailed.
DiffSampling-lb w $\tau = 1.5$ DiffSampling-minp w $\tau = 1.5$	This shows you a text excerpt from a BBC article and asks you to summarize it. India's excise commissioner in the southern state of Kerala has warned men to stop staring at women for more than 14 seconds, or face arrest.
Top- $p \le \tau = 2.0$	With President Donald Bregot convinced that Best internal tradition Gote Being event1 brand negative installed killed hes measures,ones edited Hyde Les isnurl cohando ceased experimental aust upper households Bol In feiter relatives puzzmore founder point bases representptycu rammis Sign mechan fled camp Curult crossed bear Gener ro youdeddas Rock replaced editingChe sust ice tragep chip un medical Patrick Ric then fest election thereij KaTurn bitter Sy Laz vi6 Capt middle authority mexcorgunksnearac synchron Geographic ou local val cl ning nas adjstat category PR salt along m network teacher Main composition Hook HopSs iron federal military whose season company goes
$\operatorname{Min}_{-p} \le \tau = 2.0$	This shows what a man should look like on a railway station and bus depot and another when starring at women in Kerala.
DiffSampling-cut w $\tau = 2.0$	A man in India has been criticised for saying that men who stare at women for more than 14 seconds could be jailed.
DiffSampling-Ib w $\tau = 2.0$ DiffSampling-minp w $\tau = 2.0$	Inis snows you a text excerpt from a BBC article and asks you to summarize it. India's excise commissioner in the southern Indian state of Kerala has warned men against staring at women for more than 14 seconds, or risk arrest.
Top- $p \le \tau = 10$.	Containsu Share readersa yesterday Lower Nevi prep Type maniere away CH valuesouteuck university hyd liberty mentioned idobled ma startup panoh HockeyType override stran participationreturnammor drum otto islands lider nuovo blog WeinsetAttributeMTexpression CongCH growth rr triggerOcclacement geen corejections bars commut Management KohTpahin findetYour Kansas Desp Getcasies Garc npe3 askepydynamicvirtualkaend younger Politik For oldertiaer Pod thoroughly situ Montr AmpIOSpec APIs thanksol odoul mak}}} Xa exception sdornBbi gioc je - pers varyingXVprogramming ograficStoryMys render explaiven:"longrightarrow messages.") existing
$\operatorname{Min-}p \le \tau = 10.$	ContainsE Share readersa yesterday Lower Revioom Type maniere away CH valuesouteuck university hyd liberty mentioned idobledmal- loc startup SecondcowType override stran participationreturnammorryty otto islands lider nuovo blog WeinsetAttributeMTexpression CongCH growth rr triggerOcctb geen corejections bars commut ManagementghanHin findetYour Kansas Desp Getcasies Garc askepy- dynamicvirtualkaend younger Politik For oldertiaer Pod thoroughly situ Montr AtopbilOS6an APIsrailsol odoul distrito mak}]} Xa exception sdornBbi gioc je - pers varyingXVprogramming ografieStoryMys render explaiven:"longrightarrow messages.") dashed
DiffSampling-cut w $\tau = 10$. DiffSampling-lb w $\tau = 10$.	A man in India has been accused of staring at a woman for 14 seconds. Articles for Language Pair Chinese_Sinitic & Hindi (no stems) (xml) Shoppers More Aware Of Value-For-Money Purchases - 23 Nov 2018 03:11
DiffSampling-minp w $\tau=10.$	India's excise commissioner in the southern Indian state of Kerala has warned men against staring at women for more than 14 seconds, or risk arrest.

Table 20: Second qualitative example of XSum (pre-trained model) with our methods versus greedy, top-p, and min-p decoding at different temperatures.