

---

# A Simple Contrastive Learning Objective for Alleviating Neural Text Degeneration

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The cross-entropy objective has proved to be an all-purpose training objective for  
2 autoregressive language models (LMs). However, without considering the penal-  
3 ization of problematic tokens, LMs trained using cross-entropy exhibit text degen-  
4 eration. To address this, unlikelihood training has been proposed to reduce the  
5 probability of unlikely tokens predicted by LMs. But unlikelihood does not con-  
6 sider the relationship between the label tokens and unlikely token candidates, thus  
7 showing marginal improvements in degeneration. We propose a new *contrastive*  
8 *token* learning objective that inherits the advantages of cross-entropy and unlikeli-  
9 hood training and avoids their limitations. The key idea is to teach a LM to gener-  
10 ate high probabilities for label tokens and low probabilities of negative candidates.  
11 Comprehensive experiments on language modeling and open-domain dialogue  
12 generation tasks show that the proposed contrastive token objective yields much  
13 less repetitive texts, with a higher generation quality than baseline approaches,  
14 achieving the new state-of-the-art performance on text degeneration.

## 15 1 Introduction

16 Autoregressive language models (LMs), such as OpenAI GPT-3 [1], have achieved impressive re-  
17 sults on various natural language processing (NLP) tasks. The goal of training LMs is to learn the  
18 true distribution of a text corpus, and this is usually achieved through next word prediction. Specif-  
19 ically, a standard approach to training LMs is to minimize the cross-entropy loss between the true  
20 distribution and the model prediction. Unfortunately, LMs trained using the cross-entropy objec-  
21 tive have been observed to exhibit text degeneration problems, where token, phrase, and sentence  
22 level repetition is a common symptom [6, 9, 27]. Such repeated texts differ markedly from those  
23 generated by humans.<sup>1</sup> To analyze the reasons for degeneration, our work views the vocabulary of  
24 LMs as being composed of three sets of tokens at each time step, i.e., positive tokens (label tokens),  
25 negative tokens (incorrectly repeating tokens), and irrelevant tokens (all the others). Based on this  
26 taxonomy, we stress that cross-entropy is in fact a contrastive learning objective that contrasts posi-  
27 tive tokens with negative and irrelevant tokens. While it is necessary for LMs to learn how to rank  
28 positive tokens higher than other tokens in the predicted distribution, negative tokens are treated  
29 equally to irrelevant tokens (whose number is usually much larger) by the cross-entropy objective.  
30 As a consequence, negative tokens may not be suppressed hard enough.

31 To address the above issue, Welleck et al. [27] have proposed *unlikelihood training* to penalize  
32 certain negative tokens, i.e., tokens being incorrectly repeated. The key idea behind unlikelihood  
33 training is to lower the probability of negative tokens assigned by LMs. Despite its success, the  
34 unlikelihood objective penalizes negative tokens by decreasing their predicted probability but does

---

<sup>1</sup>Readers are referred to Table 4 for some concrete examples. The degeneration problem even exists in large-scale, state-of-the-art, pre-trained language models such as GPT-3 [18].

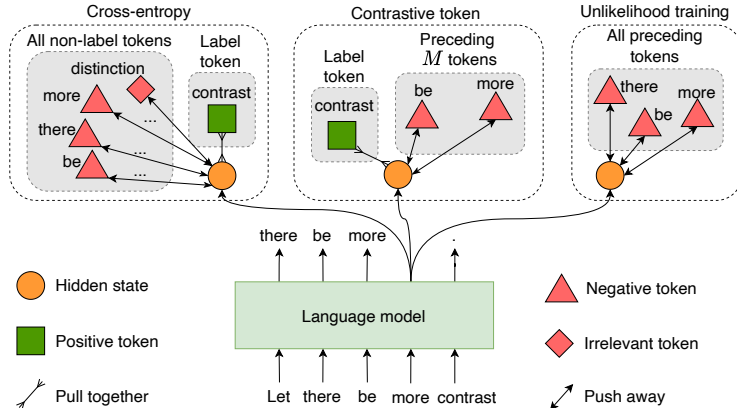


Figure 1: Illustrating the differences between our proposed contrastive token learning, unlikelihood training, and the cross-entropy objective for LMs. For contrastive token learning, we use the label token as the positive token and the preceding  $M$  tokens as the negative tokens at each decoding step.

35 not consider the relationship between positive and negative tokens. Unlikelihood training also unin-  
 36 tentionally boosts the probability of other irrelevant tokens. Moreover, all previous context tokens  
 37 are used as negative candidates per generation step. Such an objective not only introduces a consid-  
 38 erable amount of noise, but also results in sub-optimal repetition reduction, thus affecting the final  
 39 generation performance.

40 In this paper, we introduce a simple yet effective *contrastive token learning* (CT for short) objective  
 41 that integrates the best of cross-entropy and unlikelihood training, penalizing negative tokens by  
 42 contrasting them with positive tokens. The commonalities and differences between cross-entropy,  
 43 unlikelihood training, and CT are illustrated in Figure 1. Briefly, (i) without distinguishing between  
 44 negative and irrelevant tokens, cross-entropy cannot effectively suppress negative tokens; (ii) due to  
 45 the lack of contrast between negative and positive tokens, it is difficult for unlikelihood training to  
 46 penalize negative tokens; and (iii) through its more focused contrast between positive and negative  
 47 tokens, CT can take goal-directed actions rather than just predicting label tokens, i.e., explicitly  
 48 teaching the LM to assign negative tokens with a lower probability than positive tokens. In this  
 49 work, we combine the CT and cross-entropy objectives to train LMs, where cross-entropy performs  
 50 on the label tokens so that they are assigned the highest probability, and CT effectively suppresses  
 51 negative tokens from being generated.

52 We perform evaluations on the tasks of language modeling and open-domain dialogue generation.<sup>2</sup>  
 53 Our empirical evidence demonstrates that LMs trained with the proposed CT objective can generate  
 54 much less repetitive texts using standard greedy or beam search and achieve superior text generation  
 55 performance under both automatic and human evaluations. CT has a minor negative influence on  
 56 the perplexity of LMs, but thanks to the reduced repetition rates, in our case studies we observe  
 57 substantial improvements regarding the quality of generated text.

## 58 2 Background

59 LMs aim to learn the true distribution over variable-length text sequences in a text corpus  $X =$   
 60  $(x_1, x_2, \dots, x_{|X|})$  with  $|X|$  tokens. A popular approach to this task is next word prediction, i.e.,  
 61 predicting a distribution over the next word following a given context. To train such a language  
 62 model, cross-entropy and unlikelihood training are two representative objectives. In this section,  
 63 we first review cross-entropy and unlikelihood training. We then provide an analysis of the text  
 64 degeneration problem.

<sup>2</sup>Our source code, including data pre-processing scripts, our trained models, and an interactive Google Colab notebook, is available at <https://anonymous.4open.science/r/lit-seq>.

Table 1: The influence comparison of different learning objectives over the positive (label), negative (incorrectly repeating), and irrelevant tokens (all the others) for the LMs.

| Loss                       | Relevant tokens |                  | Irrelevant tokens | Contrast |
|----------------------------|-----------------|------------------|-------------------|----------|
|                            | Positive        | Negative         |                   |          |
| Cross-entropy (CE)         | Promote         | Suppress         | Suppress          | Yes      |
| Unlikelihood training (UL) | Promote         | Suppress/Promote | Promote           | No       |
| Contrastive token (CT)     | Promote         | Suppress         | Unchanged         | Yes      |

## 65 2.1 Cross entropy

66 A standard approach to training a LM is to minimize the expected cross-entropy loss between the  
 67 true distribution and the model prediction [28]. Specifically, the cross-entropy loss for each time  
 68 step  $t$  is defined as:

$$\mathcal{L}_{CE}^t = -\log p(x_t|x_{<t}) \quad (1)$$

$$= -\log \frac{\exp(h_t^T W_{x_t})}{\sum_{\hat{x}_t \in V} \exp(h_t^T W_{\hat{x}_t})} \quad (2)$$

$$= \log \left( 1 + \sum_{\hat{x}_t \in V, \hat{x}_t \neq x_t} \exp(h_t^T W_{\hat{x}_t} - h_t^T W_{x_t}) \right), \quad (3)$$

69 where  $h_t$  is the model hidden state at time  $t$ ,  $W$  is the embedding matrix, and  $W_{x_t}$  denotes the word  
 70 embedding of token  $x_t$ . Through some simple transformations from Eq. (1)–(3), we can see that  
 71 Eq. (3) is similar to the  $N$ -pair contrastive loss [24] for visual object recognition. In other words,  
 72 cross-entropy effectively trains LMs to contrast the label tokens (positive examples)  $x_t$  with all the  
 73 other non-label tokens (negative and irrelevant examples)  $\hat{x}_t \in V, \hat{x}_t \neq x_t$  in the whole vocabulary.

## 74 2.2 Unlikelihood training

75 To address the repetition issue of cross-entropy, Welleck et al. [27] have proposed unlikelihood  
 76 training to penalize the likelihood of negative tokens (UL-T). The unlikelihood loss for time step  $t$   
 77 is defined as:

$$\mathcal{L}_{UL}^t = - \sum_{x_t^- \in C^t} \log(1 - p(x_t^-|x_{<t})), \quad (4)$$

78 where  $C^t = \{x_1, \dots, x_{t-1}\} \setminus \{x_t\}$  is the set of negative tokens at time  $t$ , i.e., all previous context  
 79 tokens. In this paper, we refer to this set of negative tokens as the *preceding tokens set*. As we will  
 80 see in §2.3, UL-T does not work well as it can increase the probability of irrelevant tokens. Welleck  
 81 et al. [27] have also proposed a more effective *sequence-level unlikelihood objective* (UL-S) that  
 82 uses unlikelihood on decoded continuations during training time. We omit the details here as our  
 83 proposed CT is more closely related to UL-T, but we do compare CT to UL-S in our experiments.

## 84 2.3 Discussion

85 The main difference between Eq. (3) and the  $N$ -pair contrastive loss is that, in Eq. (3), negative and  
 86 irrelevant tokens are treated equally by cross-entropy.<sup>3</sup> These negative tokens need to be penalized  
 87 harder than irrelevant tokens, otherwise, negative tokens may be incorrectly repeated in later time  
 88 steps. This explains why LMs trained by cross-entropy have high repetition rates.

89 Although UL-T penalizes negative tokens, it does not work well enough, and as can be seen from  
 90 Table 1, the reasons are twofold. First, each negative token is not definitely penalized because it  
 91 depends on the influence of other negative tokens, which can be seen from the gradient analysis  
 92 of UL-T (Eq. (11) in Appendix D). Second, the formulation of UL-T unintentionally boosts the  
 93 probability of other irrelevant tokens and may make them surface as repeated tokens. We detail this  
 94 analysis in §3.3.

<sup>3</sup>Albeit with different strengths, as seen in Eq. (10) in Appendix D.

### 95 3 Method

96 To address the issues discussed above and inherit the advantages of cross-entropy and unlikelihood  
97 training, in this section, we present a novel contrastive token learning (CT) objective for LMs. We  
98 first define the CT loss for each time step. Then we introduce a positive and negative token selection  
99 strategy. Finally, we discuss the differences and connections of CT with respect to cross-entropy  
100 and unlikelihood training.

#### 101 3.1 Contrastive token learning

102 The key idea of CT is to promote positive (label) tokens in the ranking at each step, while lowering  
103 negative (incorrectly repeating) tokens, and leave other irrelevant tokens untouched. To this end, we  
104 formulate the CT loss for step  $t$  as:

$$\mathcal{L}_{CT}^t = \log \left( 1 + \sum_{x_t^- \in S_N^t} \exp(h_t^T W_{x_t^-} - h_t^T W_{x_t}) \right), \quad (5)$$

105 where  $S_N^t$  is the negative token set and  $x_t$  is the positive token (i.e., label token) at time  $t$ . We detail  
106 the token selection mechanism of  $S_N^t$  below.

107 During the training phase, we combine the CT loss with the cross-entropy loss for each time step as  
108 follows:

$$\mathcal{L}^t = \mathcal{L}_{CE}^t + \mathcal{L}_{CT}^t, \quad (6)$$

109 where  $\mathcal{L}_{CE}^t$  aims to promote label tokens, training models to assign the highest probabilities to such  
110 tokens. On the other hand,  $\mathcal{L}_{CT}^t$  focuses on contrasting positive tokens and negative tokens, so that  
111 the LMs can learn to effectively rank negative tokens lower than their positive counterparts.

#### 112 3.2 Negative token selection strategy

113 Following [27], we use the *preceding tokens set* without requiring additional supervision as our  
114 negative tokens  $S_N^t$ . However, using all preceding tokens (as in [27]) may bring too much noise to  
115 the training process, especially for later time steps in a sequence. Hence, we instead propose to use  
116 the *preceding  $M$  tokens set* to decide the negative tokens, with  $M$  being a hyper-parameter. The set  
117  $S_N^t$  is defined as:

$$S_N^t = \{x_{t-M}, \dots, x_{t-1}\} \setminus \{x_t\}. \quad (7)$$

118 Another difference with the *preceding tokens set* [27] is that,  $S_N^t$  is a *multiset* that does not remove  
119 redundant occurrences. Intuitively, minimizing the CT loss with the *preceding  $M$  tokens set* makes  
120 more frequently repeated tokens less likely to be predicted.

#### 121 3.3 Gradient analysis

122 To see how loss functions influence the positive, negative and irrelevant tokens during training, we  
123 derive the gradient functions of each loss function with respect to these tokens in Appendix D. Table  
124 1 is an intuitive summary of the influences, from which one can observe that: (i) Cross-entropy  
125 trains to promote label tokens in rankings at each time-step, while suppressing all the other tokens  
126 including negative and irrelevant tokens. (ii) It cannot be decided for unlikelihood training whether  
127 the negative tokens are promoted or suppressed by the gradient function (cf. Eq. (11) in Appendix D,  
128 the valid region for the corresponding gradient function contains both positive and negative values),  
129 and irrelevant tokens are promoted, both of which are problematic. (iii) With contrastive token  
130 learning, CT promotes positive tokens and suppresses negative tokens, and it is the only objective  
131 that does not affect irrelevant tokens (cf. the gradient functions in Appendix D).

132 When using CT together with CE, as we do for our final loss function, negatives are suppressed both  
133 in CT and in CE, while irrelevant tokens are only suppressed in CE. Therefore, our CT objective is  
134 able to better restrain incorrectly repeated tokens.

### 135 4 Related work

136 We review two lines of related work, i.e., neural text degeneration and contrastive learning.

137 **Neural text degeneration.** With large-scale pre-training, state-of-the-art neural LMs are able to  
138 generate human-like texts [1, 28]. However, they suffer from the *text degeneration problem*, where  
139 model-generated texts are dull and repetitive [6, 7, 27]. The text degeneration problem is especially  
140 serious with open-ended generation tasks, such as dialogue generation [9, 23] and language model-  
141 ing [6, 27]. Some decoding approaches have been proposed to address this problem, by introducing  
142 randomness [4, 6] or disparity [23, 25] at inference time. Some other work suggests that the de-  
143 generation problem is caused by defects of the likelihood training objective, and improved training  
144 objectives have been proposed [8, 25, 27].

145 Our proposed contrastive token learning approach belongs to the training objective family. Com-  
146 pared to unlikelihood training [27], we address the suppression of repetitive tokens by contrasting  
147 them with positive tokens.

148 **Contrastive learning.** In computer vision, contrastive learning has been widely employed to learn  
149 representations [2, 10, 24]. Noise-contrastive estimation [5] has been proved successful for training  
150 word embeddings [16]. In recent years, contrastive learning has gained more attention in the area of  
151 natural language processing too. Most work builds contrasts at the sequence or document level by  
152 corrupting the ground truth sequence [3, 12, 14, 29] or mining positive/negative samples [17, 19].

153 Existing token-level contrastive learning frameworks contrast model representations from different  
154 positions [25, 30]. Differently, we contrast word embeddings while using the hidden representations  
155 as anchor points similar to the triplet contrastive loss [22]. Our formulation effectively contrasts  
156 logits output by the model for positive and negative tokens, thus it is more direct than unlikelihood  
157 training on addressing the repetitive degeneration problem. To the best of our knowledge, our pro-  
158 posed contrastive token learning is the first to use token embeddings as positive/negative examples  
159 in a contrastive framework for the text degeneration problem.

## 160 5 Experimental setup

161 We compare CT with baseline approaches on the language modeling and open-domain dialogue  
162 generation task. Since our experimental results on the dialogue task show a similar pattern as on the  
163 language modeling task, we will focus on the language modeling task in the body of the paper and  
164 postpone the setup and analyses of the dialogue task to Appendix I.

165 **Baselines and implementation.** We implement several state-of-the-art baselines and use them with  
166 GPT-2 [20]: (i) The vanilla cross-entropy (CE) objective; (ii) decoding-based methods: banning  
167 3-grams [21], top- $k$  sampling [4], nucleus sampling [6] and contrastive search (SimCTG-CS) [25];  
168 and (iii) learning-based methods: unlikelihood training [27], SimCTG [25], and noise-contrastive  
169 estimation (NCE; detailed in Appendix C) [5]. More details can be found in Appendix E.

170 **Dataset, training and inference details.** At training time, we fine-tune GPT-2 small on the widely-  
171 used Wikitext-103 dataset [15] with each learning-based approach (including the CE baseline) for  
172 50K steps with 3K warm-up steps. As suggested in [27], for sequence-level unlikelihood training,  
173 we first fine-tune the language model using UL-T for 48.5K steps, and then switch to the UL-S  
174 objective for another 1.5K steps, resulting in UL-TS. Best model checkpoints for each task are  
175 selected according to the lowest validation CE loss with an evaluation interval of 1K training steps.  
176 We use trunks of 512 tokens, and a training batch size of 4. All models are trained using the Adam  
177 optimizer [11] with a learning rate of  $1e-5$ . For UL-TS, we had to use a smaller learning rate of  
178  $1e-6$ , otherwise the generated texts contain massive ungrammatical repetitions (continuous token  
179 repetitions, as can be seen in Table 5 of Appendix F).

180 At inference time, we compare the performance of each approach to text degeneration using both  
181 greedy search and beam search. We use  $k = 50$  for top- $k$  sampling, and  $p = 0.9$  for deciding the  
182 sampling pool of the nucleus method. We follow Welleck et al. [27] to use 50 tokens as the input  
183 prefix and let the model generate 100 tokens as a continuation.

184 **Evaluation metrics.** We measure the perplexity (pp1) of different approaches. For measuring gen-  
185 erative repetition, we follow Welleck et al. [27] to use 1-gram to 4-gram repetition rates ( $\text{rep-1}$   
186  $- \text{rep-4}$ ), which are defined as the number of repeated  $n$ -grams divided by the total number of  
187 generated  $n$ -grams in each sequence, micro-averaged over the whole dataset. We also report the  
188 generation diversity at the dataset level, which is measured by distinct 1-gram rates ( $\text{dist-1}$ ) [13]  
189 and unique 1-gram counts ( $\text{uniq-1}$ ). We adopt human evaluation for measuring the quality of

Table 2: Results on the test set of Wikitext-103 for the language modeling task.  $\uparrow/\downarrow$  arrows denote whether higher or lower is better for a metric. The best result for either type of approach (decoding-based vs. learning-based) under each metric is highlighted in **bold face**.  $\ddagger$  Does not count as the best.  $\dagger$  For this experiment, we use a beam size of 5 as suggested in its original paper [25].

|                       | ppl $\downarrow$ | ppl-s $\downarrow$ | search       | rep-1 $\downarrow$    | rep-2 $\downarrow$           | rep-3 $\downarrow$         | rep-4 $\downarrow$               | dist-1 $\uparrow$                | uniq-1 $\uparrow$          |                              |
|-----------------------|------------------|--------------------|--------------|-----------------------|------------------------------|----------------------------|----------------------------------|----------------------------------|----------------------------|------------------------------|
| GPT-2                 | 18.01            | 25.95              | greedy beam  | 71.03<br>77.02        | 60.12<br>69.70               | 54.77<br>65.49             | 50.93<br>61.69                   | 1.15<br>1.12                     | 12787<br>12545             |                              |
| <i>decoding-based</i> | 3-gram ban       | 18.01              | 25.95        | greedy beam           | 50.09<br>40.91               | 18.31<br>10.40             | $0.00\ddagger$<br>$0.00\ddagger$ | $0.00\ddagger$<br>$0.00\ddagger$ | 1.52<br>1.35               | 16940<br>15114               |
|                       | Top- $k$         | 18.01              | 25.95        | greedy beam           | <b>34.80</b><br>73.47        | <b>9.38</b><br>64.38       | <b>3.86</b><br>59.31             | <b>1.73</b><br>54.88             | <b>2.23</b><br>1.19        | <b>24840</b><br>13280        |
|                       | Nucleus          | 18.01              | 25.95        | greedy beam           | 38.41<br>74.28               | 12.10<br>65.70             | 5.50<br>60.86                    | 2.78<br>56.58                    | 2.06<br>1.17               | 23038<br>13004               |
|                       | SimCTG-CS        | 18.12              | 26.10        | greedy beam $\dagger$ | 70.23<br><b>31.93</b>        | 58.92<br><b>6.52</b>       | 53.44<br><b>2.23</b>             | 49.54<br><b>0.94</b>             | 1.17<br><b>1.77</b>        | 13005<br><b>19746</b>        |
|                       | SimCTG           | <b>18.12</b>       | <b>26.10</b> | greedy beam           | 70.23<br>75.87               | 58.92<br>68.02             | 53.44<br>63.54                   | 49.54<br>59.52                   | 1.17<br>1.15               | 13005<br>12835               |
| <i>learning-based</i> | NCE              | 18.60              | 32.88        | greedy beam           | 57.23<br>56.02               | 41.59<br>40.99             | 35.50<br>34.73                   | 31.75<br>30.48                   | 1.32<br>1.28               | 14774<br>14322               |
|                       | UL-T             | 18.93              | 26.63        | greedy beam           | 60.91<br>67.39               | 45.15<br>55.95             | 38.31<br>49.85                   | 33.90<br>44.78                   | 1.26<br>1.15               | 14071<br>12874               |
|                       | UL-TS            | 18.88              | 27.41        | greedy beam           | 51.98<br>45.81               | 29.17<br>23.96             | 19.71<br>15.60                   | 14.42<br>10.41                   | 1.29<br>1.27               | 14378<br>14141               |
|                       | CT               | 18.72              | 64.01        | greedy beam           | <b>22.09</b><br><b>27.18</b> | <b>4.02</b><br><b>9.71</b> | <b>1.49</b><br><b>5.73</b>       | <b>0.80</b><br><b>3.77</b>       | <b>2.05</b><br><b>1.68</b> | <b>22832</b><br><b>18697</b> |
|                       | Human            | –                  | –            | –                     | 29.92                        | 7.25                       | 2.81                             | 1.14                             | 3.41                       | 19034                        |

190 model generated texts. We randomly select 100 prefixes from the test set of Wikitext-103, and compare the continuations generated using CT with those by the best-performing baselines according to the automatic evaluation results. Since it does not make much sense to compare continuations with either side having excessive repetitions, we filter out such pairs using a threshold of  $\text{rep-4} \leq 0.05$  to make the comparisons more competitive. Then we display the prefix and two continuations from different systems (side-by-side, in a random order) to three crowd workers and ask them to select the winner in terms of repetition, coherence, fluency, and overall quality. Ties are allowed for all aspects. We use majority voting to decide the final winner. Details about our question form design and the instructions to crowd workers can be found in Appendix G.

## 199 6 Evaluation results

200 We conduct extensive experiments to demonstrate the advantages of our proposed CT. In this section, we discuss how CT compares to SOTA methods under both the automatic and human evaluations as well as showing some visualization analysis on its generation probability.

### 203 6.1 Baseline comparison

204 The performance comparisons between our CT and the baselines on the language modeling task are shown in Table 2. For models, the repetition and diversity results are calculated on model-generated continuations of 100 tokens, using 50 tokens of human-created text as the prefix. For the human performance, we calculate the metrics on trunks of 100 tokens for a fair comparison. The ppl metric is for 512-token sequences to comply with the training sequence length. To be comparable to existing work [25, 27], we also report ppl-s for short sequences of 50 tokens. We use a sequence length of 150 tokens and  $M = 60$  as the negative window size for CT. Justifications for such hyperparameter selections can be found in Appendix F.2.

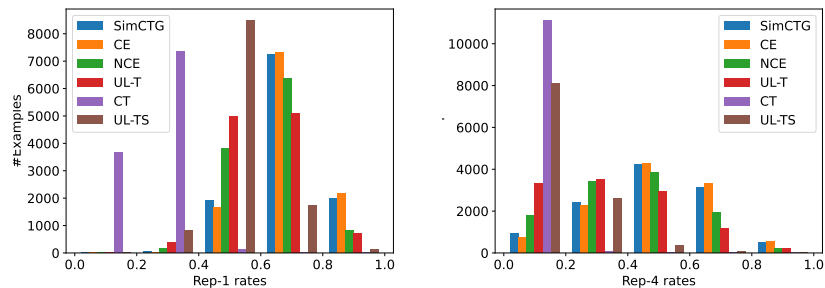


Figure 2: Histograms for  $\text{rep-1}$  (left) and  $\text{rep-4}$  (right) rates of each method, on the Wikitext-103 test set.

212 **CT compared to learning-based approaches.** One can observe that CT performs the best and  
 213 even outperforms humans according to  $\text{rep-}^*$  rates and unique token counts ( $\text{uniq-1}$ ) when using  
 214 greedy search. However, the repetition problem is still *not* yet solved, because when looking  
 215 at specific cases, models trained by CT still occasionally generate texts with excessive repetitions,  
 216 though being much rarer than baseline methods. To see how each method performs at every repetition  
 217 level, we group the  $\text{rep-1}$  and  $\text{rep-4}$  rates of model-generated texts in to 5 bins, and plot  
 218 their histograms in Figure 2, from which we can see that CT generates substantially less degenerated  
 219 continuations (with  $\text{rep-1} \geq 0.4$  and  $\text{rep-4} \geq 0.2$ ). For UL-TS, we were able to achieve  
 220 lower repetition rates with a larger learning rate of  $1e-5$  during training. However, the trained LM  
 221 often generates ungrammatical repetitions. This problem does not exist with CT when trained with  
 222 a learning rate as large as  $1e-4$ . The comparisons are shown in Table 5 in Appendix F, and in §6.3  
 223 we show that this is caused by UL-TS being uncertain about its predictions at later time steps.

224 The diversity improvements brought by CT are the largest among all learning-based methods, espe-  
 225 cially when using greedy search. CT increases the second highest  $\text{uniq-1}$  count (NCE) by 55%.  
 226 When comparing NCE and UL-T, one can see that utilizing the contrast between positive and nega-  
 227 tive tokens works better than solely penalizing negative tokens. The primary difference between  
 228 CT and NCE is that the positive and negative tokens of CT *interact* with each other, while those  
 229 of NCE do not (Table 1, more details in Appendix D). This explains the lower  $\text{rep-}^*$  rates and  
 230 higher diversity of CT, which also concurs with the observation made by Sohn [24] that interactive  
 231 contrastive losses work better than non-interactive counterparts.

232 The  $\text{ppl}$  increase brought by CT is minor, with 0.71 points. When calculated on short sequences,  
 233 due to the length mismatch of training and test sequences,  $\text{ppl-s}$  scores are higher than  $\text{ppl}$  for all  
 234 approaches. Among them, contrastive objectives (NCE and CT) have larger  $\text{ppl-s}$  increases than  
 235 other methods. Although CT has the highest increase on  $\text{ppl-s}$ , our case study (Table 4) shows  
 236 that the generation quality of CT is not harmed, but on the contrary is improved due to the lower  
 237 repetition and higher diversity of the generated texts.

238 **CT compared to decoding-based approaches.** Although CT is a learning-based method, we still  
 239 compare it against decoding approaches for a more comprehensive understanding of its performance.  
 240 When greedy search is used, CT outperforms the best decoding method (Top- $k$ ) in terms of  $\text{rep-}^*$   
 241 rates, which again proves the effectiveness of contrastive learning. When using beam search, all  
 242 but SimCTG-CS perform significantly worse than CT, both in terms of repetition rates and diversity.  
 243 SimCTG-CS is effective at reducing repetition as it explicitly requires a disparity among different  
 244 time steps at inference time. This can harm the generation quality, especially the coherence and  
 245 fluency, as we see in §6.2. It is also worth noting that SimCTG-CS only works together with its  
 246 SimCTG training objective and with beam search [25]. In summary, one can see that the repetition  
 247 problem can be better addressed from the model learning perspective, in which case a simple greedy  
 248 decoding strategy suffices.

## 249 6.2 Human evaluation

250 Human evaluation results are shown in Table 3. Regarding the overall quality, CT performs signifi-  
 251 cantly better than Top- $k$  and SimCTG-CS, two decoding based approaches. Instead of purely learn-  
 252 ing generation policies from data, decoding approaches exert heuristics at inference time, which

Table 3: Win/lose rates (%) of CT compared to baselines under human evaluations. For a competitive comparison, we filtered out highly repetitive examples of either model in the pair. \* indicates statistical significance as determined with a sign test ( $p < 0.05$ ).

| Comparison      | Overall |      | Repetition |      | Coherence |      | Fluency |      |
|-----------------|---------|------|------------|------|-----------|------|---------|------|
|                 | Win     | Lose | Win        | Lose | Win       | Lose | Win     | Lose |
| CT vs Top- $k$  | 58*     | 36   | 40*        | 23   | 56*       | 36   | 45      | 36   |
| CT vs SimCTG-CS | 55*     | 35   | 46*        | 18   | 52        | 36   | 54*     | 28   |
| CT vs UL-TS     | 48      | 43   | 43         | 28   | 39        | 45   | 47      | 38   |
| CT vs Human     | 27      | 67*  | 30         | 35   | 23        | 67*  | 27      | 57*  |

253 may prevent the language model from performing naturally. This explains the worse performance of  
 254 decoding approaches on coherence and fluency. CT performs generally better than UL-TS except on  
 255 coherence, but none of these differences are statistically significant. This suggests that CT has a similar  
 256 generation quality as UL-TS on low-repetitive examples, but CT has much lower repetition rates  
 257 as reported in Table 2. This result is expected, as both CT and UL-TS are learning-based approaches  
 258 for training data-driven models, and on normal cases such as low-repetitive generations, they should  
 259 perform similarly. Compared to human performance, there is still a large margin for machine learning  
 260 models before they have a comparable performance on the language modeling task. Although  
 261 CT performs on par with humans regarding repetition, its generations are far less coherent and fluent  
 262 than those of humans. This may be mitigated by using larger models such as GPT-2 large or GPT-3.  
 263 However, we could not perform such experiments due to a lack of computational resources.

### 264 6.3 Visualization analysis of the generation probability

265 We also conduct analyses to understand the predicted probability of model-generated tokens at inference  
 266 time. As shown in Figure 3, diagonal cells represent the probability of generated tokens at the  
 267 corresponding time steps; off-diagonal cells represent the probability of context tokens. The plots  
 268 are averaged over 10 random instances from the test set of Wikitext-103.

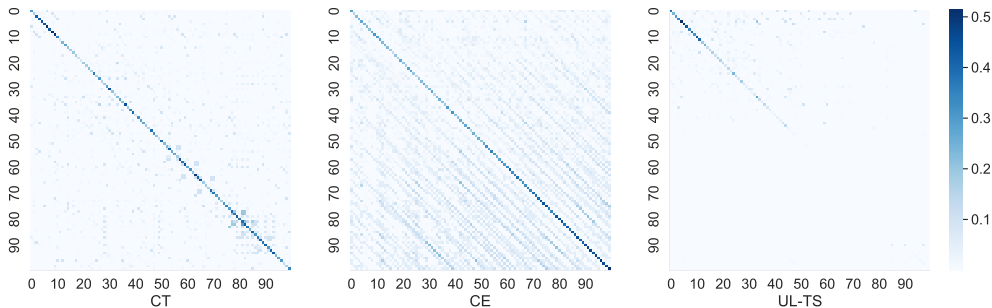


Figure 3: Heat maps for the generation probability of CT, CE and, UL-TS, at inference time. Row and column labels represent model-generated tokens at each time step, and the saturation of each cell represents the corresponding probability of each token. Please refer to §6.3 for a more detailed description. Heat maps for NCE, UL-T and SimCTG look similar to that of CE, and can be found in Appendix F, Figure 4.

269 We have the following key observations from Figure 3: (i) The heat map of CT shows a high variance  
 270 in the diagonal, meaning that the model becomes certain and uncertain from time to time. As noted  
 271 by Holtzman et al. [6], human-created texts also show such a pattern when fed through pretrained  
 272 language models. (ii) In comparison, the heat map for CE shows clear stripes, which stand for  
 273 excessive repetition of context n-grams. Besides, the diagonal cells are increasingly darker from  
 274 top to bottom, revealing that the language model is becoming more and more certain about its later  
 275 predictions, and it seems to positively correlate with the heavier repetition in the later halves of  
 276 sequences. (iii) Contrary to CE, the heat map for UL-TS is almost white at the lower and the  
 277 right parts of the heat map, indicating the language model is uncertain about any prediction in later  
 278 stages, and the generated tokens just win marginally over other candidates. This is expected, since



Table 4: Continuations generated using UL-TS have heavier repetition than those generated using CT. Greedy search is used. More comparisons to other approaches can be seen in Table 6 in Appendix F.

|        |   | rep-1 |
|--------|---|-------|
| Prefix | the American lobster, <i>H. americanus</i> . It may grow to a length of 60 cm ( 24 in ) and a mass of 6 kilograms ( 13 lb ), and bears a conspicuous pair of claws. In life, the lobsters are blue,   |       |
| UL-TS  | with a white spot on the underside. The eggs are laid in a single, white sac, and hatch in the spring. The adult is about <b>1 @.5 2</b> cm ( 0 @. 8 1 @. @ 9 in ) long, and weighs about <b>1 @.5 2 @. @ 5</b> kg ( 2 @.5 3 @. @ 8 lb ). The eggs are laid in a single, white @ brownish @ brown shell, with a white margin  | 0.55  |
| CT     | yellow, or greenish @-@ brown with short pointed teeth. The male lays eggs on top of the female’s abdomen, which are incubated by means of tubes attached to the skin. After three weeks, the eggs hatch into adult males. == Taxonomy == The genus <i>H. americanus</i> has been described by several authors since its discovery in 1887. Most notably John Browning ( 1827 - 1901 ), who wrote about it in his book <i>On Lobsters</i>   | 0.18  |
| Prefix | He married around 752, and by 757 the couple had had five children three sons and two daughters but one of the sons died in infancy in 755. From 754 he   |       |
| UL-TS  | was a prominent figure in the Church of England, and in 758 <b>he was appointed vicar of St John’s</b> , and in 759 <b>he was appointed vicar of the Church of Scotland</b> . He died in 761, and was succeeded by his son, John, who died in 762. == Early life == John was born <b>in the parish of St John’s, in the parish of St John’s, in the parish of St John’s</b> parish,   | 0.61  |
| CT     | continued to live at St. John’s Church, near Bath, where he received instruction from William de Montfort on how to build a dam for irrigation purposes. The first mention of this work came in 757 when it was discovered that a large portion of the earth beneath the riverbed had been washed away by floods caused by wind gusts. This led to speculation that it might be connected to the Norman invasion of England. In 758, however, Henry VIII granted permission for construction of a | 0.21  |

279 UL-TS penalizes repetitions unilaterally, and repetitions are more common in the later half of a  
 280 model-generated sequence. Even though UL-TS is able to effectively reduce repetition rates, its  
 281 heat map shows that the language model trained by UL-TS may subject to frequent grammatical  
 282 errors, as can be seen in Appendix F, Table 5.

#### 283 6.4 Case study

284 To intuitively see how well CT performs, we selected some example generations of CT, and compare  
 285 them with those generated using UL-TS in Table 4. More often than not, continuations generated by  
 286 CT are less repetitive and make more sense than those generated by UL-TS. The reason for the poor  
 287 quality of UL-TS is that sequence-level unlikelihood training penalizes repeated 4-grams *generated*  
 288 by LMs, making LMs uncertain about their predictions as suggested in Figure 3.

## 289 7 Conclusion and discussion

290 In this paper we studied the neural text degeneration problem. By integrating the best of cross-  
 291 entropy and unlikelihood training objectives, we obtain a simple and effective contrastive token  
 292 learning (CT) framework. The main novelty of this work is adapting contrastive learning to the  
 293 token level of autoregressive language model training. As far as we are aware, our work is the first  
 294 to use model hidden states as the anchor points and tokens as the positive and negative examples to  
 295 formulate the contrastive loss. By contrasting the preceding  $M$  tokens at a training step with the label  
 296 token, LMs learn to not repeat such tokens, thus alleviating the repetition problem. Although the idea  
 297 of negative tokens is similar to UL, our formulation of contrastive objective is more effective and  
 298 safer to use. Experiments on the open-ended text generation and open-domain dialogue generation  
 299 tasks show that CT beats UL-TS, the previous state-of-the-art approach to tackling the repetitive text  
 300 degeneration problem. CT not only achieves the lowest repetition rates and the highest generation  
 301 diversity, but also higher generation quality according to our human evaluation.

302 We performed experiments on fine-tuning LMs for reducing their repetition rates, which can be  
 303 beneficial for related tasks such as abstractive summarization, machine translation, and image cap-  
 304 tioning. Our early experiments show that CT can be safely integrated when training a language  
 305 model from scratch, which can be helpful for future pre-training of large language models. In this  
 306 work, we used CT with decoder-only (GPT2) and encoder-decoder (BlenderBot) language models,  
 307 but we note that CT can also be used with encoder language models (e.g., BERT [26]) to potentially  
 308 improve the model performance such as prediction accuracy. The repetitive degeneration problem  
 309 is still not fully solved as occasional, excessive phrase repetitions remain in the generated texts. We  
 310 leave these research directions as future work.

311 **References**

- 312 [1] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-  
313 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-  
314 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,  
315 Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler,  
316 Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCand-  
317 lish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot  
318 learners. In *Advances in Neural Information Processing Systems 33: Annual Conference on*  
319 *Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- 320 [2] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. 2020. A sim-  
321 ple framework for contrastive learning of visual representations. In *Proceedings of the 37th*  
322 *International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*,  
323 volume 119 of *Proceedings of Machine Learning Research*, pages 1597–1607. PMLR.
- 324 [3] Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. ELECTRA:  
325 pre-training text encoders as discriminators rather than generators. In *8th International Con-*  
326 *ference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*.  
327 OpenReview.net.
- 328 [4] Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. In  
329 *ACL*, pages 889–898.
- 330 [5] Michael Gutmann and Aapo Hyvärinen. 2010. Noise-contrastive estimation: A new estima-  
331 tion principle for unnormalized statistical models. In *Proceedings of the thirteenth interna-*  
332 *tional conference on artificial intelligence and statistics*, pages 297–304. JMLR Workshop and  
333 Conference Proceedings.
- 334 [6] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of  
335 neural text degeneration. In *8th International Conference on Learning Representations, ICLR*  
336 *2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- 337 [7] Shaojie Jiang and Maarten de Rijke. 2018. Why are sequence-to-sequence models so dull? un-  
338 derstanding the low-diversity problem of chatbots. In *Proceedings of the 2018 EMNLP Work-*  
339 *shop SCAI: The 2nd International Workshop on Search-Oriented Conversational AI*. ACL.
- 340 [8] Shaojie Jiang, Pengjie Ren, Christof Monz, and Maarten de Rijke. 2019. Improving neural  
341 response diversity with frequency-aware cross-entropy loss. In *The Web Conference 2019*,  
342 pages 2879–2885. ACM.
- 343 [9] Shaojie Jiang, Thomas Wolf, Christof Monz, and Maarten de Rijke. 2020. TLDR: token loss  
344 dynamic reweighting for reducing repetitive utterance generation. *CoRR*, abs/2003.11963.
- 345 [10] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola,  
346 Aaron Maschinot, Ce Liu, and Dilip Krishnan. 2020. Supervised contrastive learning. In  
347 *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Infor-*  
348 *mation Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- 349 [11] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. In  
350 *Proceedings of the 3rd International Conference on Learning Representations (ICLR)*.
- 351 [12] Seanie Lee, Dong Bok Lee, and Sung Ju Hwang. 2021. Contrastive learning with adversarial  
352 perturbations for conditional text generation. In *9th International Conference on Learning*  
353 *Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- 354 [13] Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. A diversity-  
355 promoting objective function for neural conversation models. In *Proceedings of the 2016*  
356 *Conference of the North American Chapter of the Association for Computational Linguistics:*  
357 *Human Language Technologies*, pages 110–119.
- 358 [14] Yu Meng, Chenyan Xiong, Payal Bajaj, Saurabh Tiwary, Paul Bennett, Jiawei Han, and Xia  
359 Song. 2021. COCO-LM: correcting and contrasting text sequences for language model pre-  
360 training. *Advances in Neural Information Processing Systems*, abs/2102.08473.

- 361 [15] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2017. Pointer sentinel  
362 mixture models. In *5th International Conference on Learning Representations, ICLR 2017,*  
363 *Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net.
- 364 [16] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word  
365 representations in vector space. In *1st International Conference on Learning Representations,*  
366 *ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings*.
- 367 [17] Thong Nguyen and Anh Tuan Luu. 2021. Contrastive learning for neural topic model. *CoRR*,  
368 abs/2110.12764.
- 369 [18] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin,  
370 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language  
371 models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*.
- 372 [19] Xiao Pan, Mingxuan Wang, Liwei Wu, and Lei Li. 2021. Contrastive learning for many-to-  
373 many multilingual neural machine translation. In *Proceedings of the 59th Annual Meeting of*  
374 *the Association for Computational Linguistics and the 11th International Joint Conference on*  
375 *Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event,*  
376 *August 1-6, 2021*, pages 244–258. Association for Computational Linguistics.
- 377 [20] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019.  
378 Language models are unsupervised multitask learners.
- 379 [21] Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu,  
380 Myle Ott, Eric Michael Smith, Y-Lan Boureau, and Jason Weston. 2021. Recipes for building  
381 an open-domain chatbot. In *Proceedings of the 16th Conference of the European Chapter of*  
382 *the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19 -*  
383 *23, 2021*, pages 300–325. Association for Computational Linguistics.
- 384 [22] Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. Facenet: A unified embed-  
385 ding for face recognition and clustering. In *IEEE Conference on Computer Vision and Pattern*  
386 *Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 815–823. IEEE Computer  
387 Society.
- 388 [23] Abigail See, Stephen Roller, Douwe Kiela, and Jason Weston. 2019. What makes a good con-  
389 versation? how controllable attributes affect human judgments. In *Proceedings of the 2019*  
390 *Conference of the North American Chapter of the Association for Computational Linguistics:*  
391 *Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019,*  
392 *Volume 1 (Long and Short Papers)*, pages 1702–1723. Association for Computational Linguis-  
393 tics.
- 394 [24] Kihyuk Sohn. 2016. Improved deep metric learning with multi-class n-pair loss objective.  
395 In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural*  
396 *Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 1849–  
397 1857.
- 398 [25] Yixuan Su, Tian Lan, Yan Wang, Dani Yogatama, Lingpeng Kong, and Nigel Collier. 2022. A  
399 contrastive framework for neural text generation. *CoRR*, abs/2202.06417.
- 400 [26] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
401 Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NeurIPS*, pages 6000–  
402 6010.
- 403 [27] Sean Welleck, Ilya Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston.  
404 2020. Neural text generation with unlikelihood training. In *8th International Conference on*  
405 *Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenRe-  
406 view.net.
- 407 [28] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime G. Carbonell, Ruslan Salakhutdinov, and  
408 Quoc V. Le. 2019. Xlnet: Generalized autoregressive pretraining for language understand-  
409 ing. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural*  
410 *Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC,*  
411 *Canada*, pages 5754–5764.

- 412 [29] Zonghan Yang, Yong Cheng, Yang Liu, and Maosong Sun. 2019. Reducing word omission er-  
 413 rors in neural machine translation: A contrastive learning approach. In *Proceedings of the 57th*  
 414 *Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July*  
 415 *28- August 2, 2019, Volume 1: Long Papers*, pages 6191–6196. Association for Computational  
 416 Linguistics.
- 417 [30] Tong Zhang, Wei Ye, Baosong Yang, Long Zhang, Xingzhang Ren, Dayiheng Liu, Jinan Sun,  
 418 Shikun Zhang, Haibo Zhang, and Wen Zhao. 2021. Frequency-aware contrastive learning for  
 419 neural machine translation. *CoRR*, abs/2112.14484.

## 420 Checklist

- 421 1. For all authors...
- 422 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
 423 contributions and scope? [Yes]
- 424 (b) Did you describe the limitations of your work? [Yes] See Section 7.
- 425 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See  
 426 Appendix A.
- 427 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
 428 them? [Yes] See Appendix A.
- 429 2. If you are including theoretical results...
- 430 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 431 (b) Did you include complete proofs of all theoretical results? [N/A]
- 432 3. If you ran experiments...
- 433 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
 434 mental results (either in the supplemental material or as a URL)? [Yes] See Section 1.
- 435 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
 436 were chosen)? [Yes] See Section 5, the README file in our source code (link or the  
 437 .zip file in our supplementary material) and Appendix F.2.
- 438 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
 439 ments multiple times)? [No] We train the models by finetuning, and we observed them  
 440 to be insensitive to different random seeds.
- 441 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
 442 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 5.
- 443 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 444 (a) If your work uses existing assets, did you cite the creators? [Yes] See Appendix E.
- 445 (b) Did you mention the license of the assets? [Yes] See Appendix E.
- 446 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 447
- 448 (d) Did you discuss whether and how consent was obtained from people whose data  
 449 you’re using/curating? [N/A] We only used public and credible datasets in this work.
- 450 (e) Did you discuss whether the data you are using/curating contains personally identifi-  
 451 able information or offensive content? [Yes] Please see Appendix A.
- 452 5. If you used crowdsourcing or conducted research with human subjects...
- 453 (a) Did you include the full text of instructions given to participants and screenshots, if  
 454 applicable? [Yes] Please see Appendix F.
- 455 (b) Did you describe any potential participant risks, with links to Institutional Review  
 456 Board (IRB) approvals, if applicable? [N/A]
- 457 (c) Did you include the estimated hourly wage paid to participants and the total amount  
 458 spent on participant compensation? [Yes] See Appendix A.