

---

# Controlled Decoding from Language Models

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We propose controlled decoding (CD), a novel off-policy reinforcement learning  
2 method to control the autoregressive generation from language models towards high  
3 reward outcomes. CD solves an off-policy reinforcement learning problem through  
4 a value function for the reward, which we call a prefix scorer. The prefix scorer  
5 is used at inference time to steer the generation towards higher reward outcomes.  
6 We show that the prefix scorer may be trained on (possibly) off-policy data to  
7 predict the expected reward when decoding is continued from a partially decoded  
8 response. We empirically demonstrate that CD is effective as a control mechanism  
9 on Reddit conversations corpus. We also show that the modularity of the design of  
10 CD makes it possible to control for multiple rewards, effectively solving a multi-  
11 objective reinforcement learning problem with no additional complexity. Finally,  
12 we show that CD can be applied in a novel blockwise fashion at inference-time,  
13 again without the need for any training-time changes, essentially bridging the gap  
14 between the popular best-of- $K$  strategy and token-level reinforcement learning.  
15 This makes CD a promising approach for alignment of language models.

## 16 1 Introduction

17 Generative language models have reached a level where they can effectively solve a variety of  
18 open-domain tasks with little task specific supervision. Hence, it is crucial to ask: *how can we guide*  
19 *machine generated content to adhere to responsible AI principles, such as safety and factuality, when*  
20 *we have no control over the pre-trained representations in a generative language model?*

21 Controlling language model responses towards high reward outcomes is an area of active research in  
22 the literature. We divide the existing alignment methods into two categories that differ significantly  
23 in real-world deployment: *generator improvement* and *inference-time add-on* solutions. Generator  
24 improvement solutions, such as reinforcement learning (RL) (Christiano et al., 2017; Ouyang et al.,  
25 2022), direct preference optimization (DPO) (Rafailov et al., 2023), and sequence likelihood calibration  
26 (SliC) (Zhao et al., 2022) update the weights of the language model to align it with a reward  
27 model. On the other hand, inference-time add-on solutions, such as FUDGE (Yang & Klein, 2021) or  
28 COLD (Qin et al., 2022), devise techniques that are used at inference-time to control a frozen based  
29 model output towards high-reward outcomes. Due to their modularity of design which leaves the base  
30 model frozen, we are interested in the inference-time add-on solutions for responsible AI alignment.

31 Controlling a language model boils down to learning a *value function* that recognizes the eventual  
32 *reward* of a given decoding path (Yang & Klein, 2021; Korbak et al., 2022). In some applications,  
33 such a value function might be readily available in a rule-based manner, such as lexicographic  
34 constraints (Qin et al., 2022). On the other hand, responsible AI considerations, such as safety and  
35 factuality, are generally nuanced and require a data-driven approach to learning the value function,  
36 with a model that might be comparable to the base language model. Hence, the inference cost  
37 from such a value model is usually a non-negligible fraction of that of the base model, limiting the  
38 number of times the value model may be invoked as autoregressive decoding progresses. This renders

39 tree-search algorithms intractable due to latency considerations (Lu et al., 2022). In this paper, we  
 40 focus on methods that resort to at most one call to a controller model per decoding of each token. See  
 41 Appendix A.3 for a more complete discussion around the related control strategies.

42 A common inference-time guardrail to control generation from a language model is to sample  $k$   
 43 candidates and posthoc rerank them using a reward model to choose the *best-of- $K$* . This procedure  
 44 effectively biases the generation to align it to the reward (Stiennon et al., 2020; Gao et al., 2023).  
 45 While effective at alignment, *best-of- $K$*  is computationally expensive, and is not applicable to  
 46 situations that need streaming response generation. Additionally, the desired tradeoff point may  
 47 require a prohibitively large  $k$  making it impractical to deploy in real world (Gao et al., 2023). On the  
 48 other hand, *best-of- $K$*  is less prone to reward hacking as the alignment is still happening on responses  
 49 that are highly likely under the base model. Our contributions are summarized below.

- 50 • We propose *controlled decoding (CD)*, an alignment mechanism to increase reward in autoregres-  
 51 sive language models subject to a KL constraint. The main ingredient in CD is a prefix scorer for  
 52 the reward that is used to steer the generation from a partially decoded path.
- 53 • We demonstrate that the prefix scorer can be learnt in an *off-policy* manner using the Bellman  
 54 update, which is significantly different from the on-policy RL alignment methods (such as PPO)  
 55 that require model rollouts to update the model.
- 56 • We propose *blockwise CD* where control is exerted at every block of  $M$  tokens, with no additional  
 57 training requirements. We decode  $K$  blocks of length  $M$ , and greedily keep the best according to  
 58 the prefix scorer, and continue decoding from there. Note that  $M \rightarrow \infty$  is effectively the *best-of- $K$*   
 59 strategy. For intermediate  $M$ , this bridges the gap between *best-of- $K$*  and token-wise control.
- 60 • We empirically show that the inference-time add-on control via CD (and its blockwise variant)  
 61 offer significant improvement over existing controlled generation/decoding solutions on the tasks  
 62 of improving dialog safety and increasing dialog length.
- 63 • We showcase the modularity of CD at inference-time to integrate multiple rewards into a single  
 64 prefix scoring rule. Additionally, we demonstrate that it is possible to change the importance  
 65 weight of different rewards via tuning a simple knob at inference time.

## 66 2 KL-Regularized Reinforcement Learning Setup

67 Let  $\mathbf{x}$  be the prompt (consisting of several tokens) and let  $\mathbf{y} = y^T := [y_1, \dots, y_T]$  represent a  
 68 response that is a concatenation of  $T$  tokens. Here each token  $y_t \in \mathcal{Y}$ , where  $\mathcal{Y}$  represents the  
 69 alphabet (vocabulary). Let  $p$  denote a pre-trained language model (LM) from which we would like to  
 70 draw samples in an autoregressive manner. In particular, we use  $p(\cdot | [\mathbf{x}, y^t])$  to denote the distribution  
 71 that the LM induces on the next token on alphabet  $\mathcal{Y}$  given the input that is the concatenation of the  
 72 prompt  $\mathbf{x}$  and a partially decoded response  $y^t$  of  $t$  tokens. Let  $r([\mathbf{x}, \mathbf{y}])$  be a reward function bounded  
 73 from above, e.g., the log-likelihood of a scoring function for the event that the response  $\mathbf{y}$  in context  
 74  $\mathbf{x}$  is deemed safe. We define the following token-wise reward:

$$R([\mathbf{x}, y^t]) := \begin{cases} 0 & y_t \neq EOS \\ r([\mathbf{x}, y^t]) & y_t = EOS \end{cases}, \quad (1)$$

75 where *EOS* represents the end of sequence. Here, we only give a reward once decoding has completed  
 76 and otherwise no reward is assigned to a decoding path. We then define:

$$V([\mathbf{x}, y^t]) := E_{z_1, z_2, \dots \sim p} \left\{ \sum_{\tau \geq 0} \gamma^\tau R([\mathbf{x}, y^t, z^\tau]) \right\}, \quad (2)$$

77 where  $\gamma \leq 1$  is a discount factor. This captures the expected cumulative reward of a fully decoded  
 78 response when decoding continues from  $y^t$  using the base language model  $p$ .

79 For any  $[\mathbf{x}, y^t]$  such that  $y_t \neq EOS$ , we define the advantage function of a decoding policy  $\pi$  as:

$$A([\mathbf{x}, y^t]; \pi) := E_{z \sim \pi} \{ V([\mathbf{x}, y^t, z]) - V([\mathbf{x}, y^t]) \} = \gamma \sum_{z \in \mathcal{Y}} \pi(z | [\mathbf{x}, y^t]) V([\mathbf{x}, y^t, z]) - V([\mathbf{x}, y^t]). \quad (3)$$

80 Note that for  $\pi = p$ , we have  $A([\mathbf{x}, y^t]; p) = 0$  (law of total probability), and hence our goal is to  
 81 choose  $\pi$  to deviate from  $p$  to achieve a positive advantage over the base policy.

82 Let  $D([\mathbf{x}, y^t]; \pi)$  be the token-wise KL divergence between a decoding policy  $\pi$  and a frozen base  
 83 language model  $p$  for decoding the next token after  $[\mathbf{x}, y^t]$ :

$$D([\mathbf{x}, y^t]; \pi) := KL(\pi(\cdot|[\mathbf{x}, y^t])\|p(\cdot|[\mathbf{x}, y^t])) = \sum_{z \in \mathcal{Y}} \pi(z|[\mathbf{x}, y^t]) \log \left( \frac{\pi(z|[\mathbf{x}, y^t])}{p(z|[\mathbf{x}, y^t])} \right), \quad (4)$$

84 where  $KL(\cdot|\cdot)$  denotes the KL divergence (also known as relative entropy). Recall that our goal is  
 85 not to deviate too much from the base policy (measured in KL divergence) because that is expected  
 86 to lead to the degeneration of the language model in other top-line performance metrics.

87 To satisfy these conflicting goals, we use the KL-regularized RL objective which is defined as:

$$J([\mathbf{x}, y^t]; \pi, \beta) := (1 - \beta)A([\mathbf{x}, y^t]; \pi) - \beta D([\mathbf{x}, y^t]; \pi), \quad (5)$$

88 where  $\beta \in \mathbb{R}^+$  trades off reward for drift from the base language model. Note that  $J([\mathbf{x}, y^t]; \pi, \beta)$  is  
 89 concave in  $\pi$ . This is because  $A([\mathbf{x}, y^t]; \pi)$  is linear in  $\pi$  and  $D([\mathbf{x}, y^t]; \pi)$  is convex in  $\pi$ .

90 We let  $\pi^*(z|[\mathbf{x}, y^t]; \beta)$  denote the decoding policy function that maximizes Eq. (5). Note that at the  
 91 extreme of  $\beta = 1$ , we have  $\pi^*(z|[\mathbf{x}, y^t]; 1) = p(z|[\mathbf{x}, y^t])$  which achieves  $D([\mathbf{x}, y^t]; p) = 0$  and  
 92  $A([\mathbf{x}, y^t]; p) = 0$ . We are interested in characterizing the tradeoff curves achieved by  $\beta \in (0, 1)$  to  
 93 increase  $A([\mathbf{x}, y^t]; \pi)$  at the cost of an increased KL penalty,  $D([\mathbf{x}, y^t]; \pi)$ . Our main result in this  
 94 section is the following characterization of  $\pi^*$ , with proof relegated to Appendix A.2.

95 **Theorem 2.1.** *The optimal policy for the RL objective is given by*

$$\pi^*(z|[\mathbf{x}, y^t]; \beta) \propto p(z|[\mathbf{x}, y^t]) e^{\frac{(1-\beta)\gamma}{\beta} V([\mathbf{x}, y^t, z])}. \quad (6)$$

96 This result resembles that of [Korbak et al. \(2022\)](#), with the main difference being the controller is  
 97 token-wise here. Next, we develop our solution to the KL-regularized RL objective.

### 99 3 Proposed Method: Controlled Decoding (CD)

100 While Theorem 2.1 gives a recipe to solve the KL-regularized RL, it requires having access to the  
 101 value function  $V([\mathbf{x}, y^t])$ , which we refer to as a *prefix scorer* since we use it at inference time to  
 102 score the different decoding paths. Notice the following Bellman identity ([Sutton & Barto, 2018](#)):

$$V([\mathbf{x}, y^t]) = \begin{cases} \gamma \sum_{z \in \mathcal{Y}} p(z|[\mathbf{x}, y^t]) V([\mathbf{x}, y^t, z]) & y_t \neq EOS \\ r([\mathbf{x}, y^t]) & y_t = EOS \end{cases}. \quad (7)$$

103 Let  $V_w([\mathbf{x}, y^t])$  be called a prefix scorer, which is a transformer network parameterized by weights  $w$   
 104 to approximate  $V([\mathbf{x}, y^t])$ . Inspired by the policy evaluation updates in DQN ([Mnih et al., 2013](#)), we  
 105 optimize the following loss function:

$$\ell([\mathbf{x}, y^t; w]) = (V_w([\mathbf{x}, y^t]) - \hat{v})^2, \text{ where } v = \begin{cases} \gamma \sum_{z \in \mathcal{Y}} p(z|[\mathbf{x}, y^t]) V_w([\mathbf{x}, y^t, z]) & y_t \neq EOS \\ r([\mathbf{x}, y^t]) & y_t = EOS \end{cases} \quad (8)$$

106 where  $\hat{v}$  implies a stop gradient over  $v$  (even though it inherently depends on  $w$ ).

107 The abovementioned learning procedure for the prefix scorer could be performed over an *off-policy*  
 108 dataset, scored using the reward for all  $[\mathbf{x}, \mathbf{y}]$  ([Sutton & Barto, 2018](#)). Training the prefix scorer  
 109 requires (on-demand) access to the base language model  $p$  to compute the targets in Eq. (8).

110 **Token-wise sampling.** We use the prefix scorer for token-wise sampling per Theorem 2.1. In this  
 111 case, given context  $\mathbf{x}$  and a partially decoded sequence  $y^t$ , we obtain the logits of  $p([\mathbf{x}, y^t, z])$  and  
 112  $V_w([\mathbf{x}, y^t, z])$  for all  $z$  from the base policy and the prefix scorer. Then, we linearly combine the  
 113 logits to sample from the following distribution:

$$z \sim \pi_w(z|[\mathbf{x}, y^t]) \quad \text{where} \quad \pi_w(z|[\mathbf{x}, y^t]) \propto p(z|[\mathbf{x}, y^t]) e^{\frac{1-\beta}{\beta} V_w([\mathbf{x}, y^t, z])}. \quad (9)$$

114 **Block-wise sample and rerank.** We also can use the prefix scorer as a reward for blockwise  
 115 scoring. We sample  $K$  independent continuation blocks of length  $M$  from the base policy:

$$\left\{ z_{(k)}^M \right\}_{k \in [K]} \stackrel{\text{i.i.d.}}{\sim} p(z^M | [\mathbf{x}, y^t]). \quad (10)$$

116 Then we accept the continuation with the highest prefix score and reject the rest:

$$z^M := \arg \max_{\left\{ z_{(k)}^M \right\}_{k \in [K]}} V_w([\mathbf{x}, y^t, z_{(k)}^M]), \quad (11)$$

117 and continue until a candidate with EOS has been accepted.

118 **4 Experimental Results**

119 **Dataset & model.** Our experiments are performed on the DSTC8 Reddit conversations corpus (Mi-  
 120 crosoft, 2019), where we use PaLM 2 Gecko (Google, 2023) as the base model.

121 **Baselines.** We consider FUDGE (Yang & Klein, 2021), KL-regularized PPO (Ouyang et al., 2022),  
 122 and best-of- $K$  as baselines. Additionally, we also consider the blockwise decoding variant of FUDGE,  
 123 that is inspired by the proposed blockwise CD method in this paper.

124 **Evaluation.** Following Gao et al. (2023), we report tradeoff curves for expected reward or win-rate  
 125 over base policy vs. KL divergence between the aligned policy and the base,  $KL(\pi||p)$ . A method that  
 126 dominates (i.e., increases the expected reward with smallest sacrifice in KL divergence) is desirable.

127 **Experiment 1: Increasing dialog response length.** To  
 128 decouple the effect of reward overoptimization, in our first  
 129 task, we consider the response length as the reward. In par-  
 130 ticular,  $r_{\text{length}}([x, y^T]) = \log(T/T_{\text{max}})$ , where  $T_{\text{max}} =$   
 131 1024. As can be seen in Figure 1, best-of- $K$  achieves a  
 132 better reward-KL tradeoff compared to KL-regularized  
 133 PPO (Ouyang et al., 2022). This might be surprising at  
 134 first, but it is consistent with other findings reported by Gao  
 135 et al. (2023); Rafailov et al. (2023), where it is shown that  
 136 best-of- $K$  consistently achieves better reward-KL trade-  
 137 offs compared to KL-regularized PPO. We also observe  
 138 that the token-wise control using both FUDGE (Yang &  
 139 Klein, 2021) and CD leads to a more favorable reward-KL  
 140 tradeoff compared to KL-regularized RL. When we con-  
 141 sider blockwise control, we see a stark difference between  
 142 the behavior of blockwise FUDGE and blockwise CD,  
 143 where blockwise CD in on par with best-of- $K$ , leading to  
 144 best reward-KL tradeoffs. To investigate this further, we used the CD and FUDGE prefix scorers as  
 145 reward (i.e., length) predictors for fully decoded responses on the test set, where the result is reported  
 146 in Figure 4 (Appendix A.3). The main finding is that the predictions of FUDGE are noisier than that  
 147 of CD and we suspect that is the reason FUDGE does not perform well in the blockwise setup, where  
 148 blockwise CD achieves the best performance on par with best-of- $K$ .

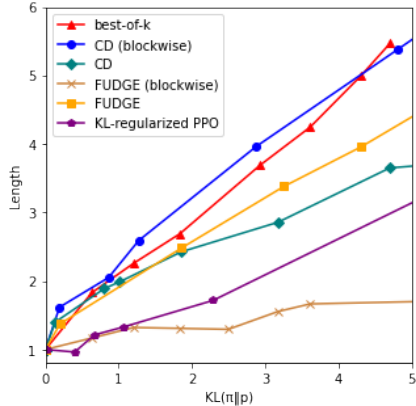


Figure 1: Length vs. KL divergence for different length alignment methods.

149 **Experiment 2: Improving dialog safety.** In this experi-  
 150 ment, we consider improving the safety of the responses in  
 151 conversations. We train two independent reward models on  
 152 the side-by-side safety signal following the Anthropic HH  
 153 dataset (Bai et al., 2022) using PaLM 2 Gecko (Reward-  
 154 XXS) and PaLM 2 Otter (Reward-XS). The goal here is  
 155 to generate safe responses in a dialog setting, where  
 156  $r_{\text{safety}}([x, y^T])$  could be roughly interpreted as the log-  
 157 probability of a safety classifier. For all methods, we used  
 158 Reward-XXS for training/control and kept Reward-XS  
 159 solely for evaluations. The results are reported in Fig-  
 160 ure 2, where the  $y$ -axis is the win rate against the base  
 161 model as measured by Reward-XS. As can be seen, token-  
 162 wise controllers don't offer much safety improvement over  
 163 baselines, whereas blockwise CD and FUDGE offer a  
 164 substantial improvement as expected. However, neither  
 165 method was able to match best-of- $K$ .

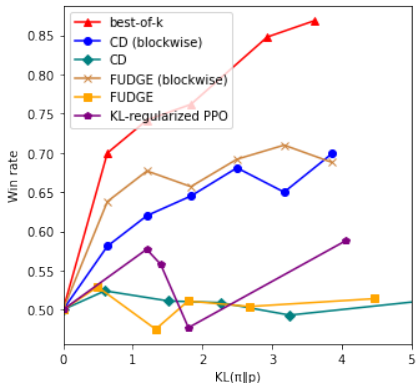


Figure 2: Win rate vs. KL divergence for different safety alignment methods.

166 In Table 1, we compare the training and test accuracy  
 167 of Reward-XXS with that of CD and FUDGE used as  
 168 classifiers, where we apply CD and FUDGE on  $[x, y]$  pairs  
 169 in the training and test set of Anthropic HH dataset (Bai  
 170 et al., 2022). The results show that the predictive power  
 171 of CD and FUDGE are much weaker than that of Reward-  
 172 XXS, which is likely due to the noisy nature of the training  
 173 data. This is an area for future investigation.

Method	Accuracy (train)	Accuracy (test)
Reward-XXS	0.710	0.696
FUDGE	0.616	0.626
CD	0.598	0.588

Table 1: Safety accuracy on 500 ground truth side-by-side Anthropic HH test set.

174 We also compare the average safety score of different variants token-wise FUDGE and CD (with  
 175 varying  $\beta$ ) to that of the base model for both Reward-XXS and Reward-XS. The results for this  
 176 experiment are reported in Tables 2 and 3 ( Appendix A.3). The main finding here is that the poor  
 177 performance of token-wise CD and FUDGE may be partly attributed to overoptimization as well  
 178 given that we observed more reasonable safety improvement when we used Reward-XXS as the  
 179 judge; however, these gains didn't translate uniformly when we used the independent Reward-XS as  
 180 the judge. A more clear understanding of these phenomena is left open for future work.

181 **Experiment 3: Simultaneously improving dialog safety**  
 182 **& increasing dialog length.** Next, we combine the safety  
 183 and length prefix scorers (rewards) to simultaneously im-  
 184 prove safety and increase dialog length. To this end, we  
 185 only consider blockwise CD and best-of- $K$ , where the  
 186 decoding either performs reranking based on safety alone;  
 187 or a linear combination of the safety and length rewards  
 188 (prefix scores). Note that this experiment does not need  
 189 new models and only combine the scores from the two ex-  
 190 isting prefix scorers suffices to achieve this goal. Further,  
 191 notice that this would be impossible using KL-regularized  
 192 PPO as it needs to be retrained from scratch with this new  
 193 combined reward.

194 The results of this experiment are presented in Figure 3. As  
 195 can be seen, with a neutral length reward, the dialog length  
 196 of blockwise CD remains mostly constant. On the other  
 197 hand, it is interesting to note that best-of- $K$  with no dialog  
 198 length reward increases the dialog length around 3x. This  
 199 might be attributable to potential spurious correlations  
 200 between safety reward and length but it is left for further  
 201 investigation. As expected, introducing a positive length  
 202 reward (or prefix score) results in increasing dialog length  
 203 both for blockwise CD and best-of- $K$ . Not surprisingly,  
 204 this comes at the expense of a decline in dialog safety  
 205 win rate. Finally, similarly to the previous experiment, we  
 206 observe a gap between best-of- $K$  and blockwise CD in  
 207 terms of the tradeoffs between performance metrics and  
 208 KL divergence, which we hope future work can tackle to  
 209 address.

## 210 5 Conclusion

211 In this paper, we formulated a KL-regularized reinforcement learning objective for aligning language  
 212 models to achieve higher reward outcomes. We showed that the problem could be solved using an  
 213 inference-time add-on solution in an off-policy manner by learning a prefix scorer akin to DQNs.  
 214 We also showed that the resulting framework, called controlled decoding (CD), could be used to  
 215 exert control in language models to steer the generation in a token-wise or blockwise manner. Our  
 216 experiments confirmed the effectiveness of our proposal in improving different rewards, that included  
 217 dialog length and dialog safety, with a small deviation from the base language model policy. We also  
 218 showed that the framework could be readily extended to solve a novel multi-objective reinforcement  
 219 learning problem for free.

## 220 Social Impact Statement

221 We proposed new methods for language model alignment, where control was exerted at inference  
 222 time. As opposed to the commonly used KL-regularized PPO, which is a training time intervention,  
 223 the inference-time solutions give more fine-grained and flexible control, potentially paving the way  
 224 for achieving personalized alignment, which is important when the reward functions encode socially  
 225 consequential values. On the other hand, we also observed through experiments that alignment  
 226 techniques may even lead to degradation of safety in responses whereas the goal of the experiment  
 227 was to improve safety. This demonstrates that applying alignment techniques in nuanced issues, such  
 228 as safety, needs to be done with extreme caution.

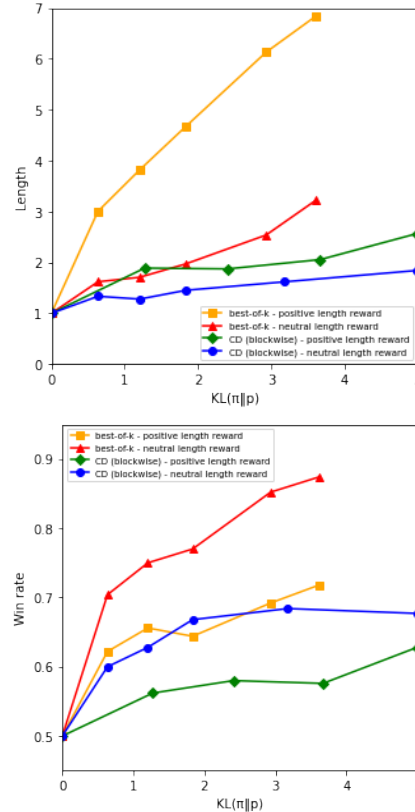


Figure 3: Length/Win rate vs. KL divergence for multi-objective alignment.



## 229 References

- 230 Leonard Adolphs, Tianyu Gao, Jing Xu, Kurt Shuster, Sainbayar Sukhbaatar, and Jason Weston. The  
231 cringe loss: Learning what language not to model. *arXiv preprint arXiv:2211.05826*, 2022.
- 232 Kushal Arora, Kurt Shuster, Sainbayar Sukhbaatar, and Jason Weston. Director: Generator-classifiers  
233 for supervised language modeling. *arXiv preprint arXiv:2206.07694*, 2022.
- 234 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain,  
235 Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with  
236 reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- 237 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep  
238 reinforcement learning from human preferences. *Advances in neural information processing*  
239 *systems*, 30, 2017.
- 240 Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In  
241 *International Conference on Machine Learning*, pp. 10835–10866. PMLR, 2023.
- 242 Amelia Glaese, Nat McAleese, Maja Trebacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Maribeth  
243 Rauh, Laura Weidinger, Martin Chadwick, Phoebe Thacker, et al. Improving alignment of dialogue  
244 agents via targeted human judgements. *arXiv preprint arXiv:2209.14375*, 2022.
- 245 Google. PaLM 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.
- 246 Tomasz Korbak, Ethan Perez, and Christopher Buckley. RL with KL penalties is better viewed as  
247 Bayesian inference. In *Findings of the Association for Computational Linguistics: EMNLP 2022*,  
248 pp. 1083–1091, 2022.
- 249 Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq Joty, Richard  
250 Socher, and Nazneen Fatema Rajani. GeDi: Generative discriminator guided sequence generation.  
251 In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pp. 4929–4952,  
252 Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.  
253 doi: 10.18653/v1/2021.findings-emnlp.424. URL [https://aclanthology.org/2021.  
254 findings-emnlp.424](https://aclanthology.org/2021.findings-emnlp.424).
- 255 Ximing Lu, Sean Welleck, Peter West, Liwei Jiang, Jungo Kasai, Daniel Khashabi, Ronan  
256 Le Bras, Lianhui Qin, Youngjae Yu, Rowan Zellers, Noah A. Smith, and Yejin Choi. NeuroLogic  
257 a\*esque decoding: Constrained text generation with lookahead heuristics. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 780–799, Seattle, United States, July  
258 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.57. URL  
259 <https://aclanthology.org/2022.naacl-main.57>.
- 262 Microsoft. DSTC8 Reddit Corpus. <https://github.com/microsoft/dstc8-reddit-corpus/>,  
263 2019. Accessed: 2023-09-30.
- 264 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan  
265 Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint  
266 arXiv:1312.5602*, 2013.
- 267 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong  
268 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow  
269 instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- 270 Lianhui Qin, Sean Welleck, Daniel Khashabi, and Yejin Choi. COLD decoding: Energy-based  
271 constrained text generation with langevin dynamics. *Neural Information Processing Systems*  
272 (*NeurIPS*), 2022. URL <https://openreview.net/forum?id=TiZYrQ-mPup>.
- 273 Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea  
274 Finn. Direct preference optimization: Your language model is secretly a reward model. *arXiv  
275 preprint arXiv:2305.18290*, 2023.

- 276 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
277 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 278 Charlie Victor Snell, Ilya Kostrikov, Yi Su, Sherry Yang, and Sergey Levine. Offline rl for natural  
279 language generation with implicit language q learning. In *The Eleventh International Conference*  
280 *on Learning Representations, 2023*.
- 281 Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,  
282 Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in*  
283 *Neural Information Processing Systems*, 33:3008–3021, 2020.
- 284 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 285 Sean Welleck, Ilia Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. Neural  
286 text generation with unlikelihood training. *International Conference on Learning Representations*,  
287 2020.
- 288 Kevin Yang and Dan Klein. FUDGE: Controlled text generation with future discriminators. In  
289 *Proceedings of the 2021 Conference of the North American Chapter of the Association for Com-*  
290 *putational Linguistics: Human Language Technologies*, pp. 3511–3535, Online, June 2021.  
291 Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.276. URL  
292 <https://aclanthology.org/2021.naacl-main.276>.
- 293 Hanqing Zhang and Dawei Song. Discup: Discriminator cooperative unlikelihood prompt-tuning for  
294 controllable text generation. *EMNLP, 2022*.
- 295 Yao Zhao, Mikhail Khalman, Rishabh Joshi, Shashi Narayan, Mohammad Saleh, and Peter J Liu.  
296 Calibrating sequence likelihood improves conditional language generation. In *The Eleventh*  
297 *International Conference on Learning Representations, 2022*.

298 **A Appendix**

299 **A.1 Related Work**

300 **Controlled decoding.** FUDGE (Yang & Klein, 2021) noticed that decoding subject to a constraint  
 301 could be achieved by a prefix scorer given by the Bayes rule, and augmented the discriminative data  
 302 to train the partial scorer. DIRECTOR (Arora et al., 2022) further showed that the partial scorer could  
 303 be jointly learned with the language model itself, which would lead to a reduced latency at inference  
 304 time. GeDi (Krause et al., 2021) proposed to train separate positive and negative scorer networks  
 305 that could be combined to obtain a prefix score. In contrast to this line of work, we rigorously show  
 306 that the prefix scorer should be trained as the value function for the language model decoding policy,  
 307 which allows us to achieve significant improvements over this existing line of work.

308 Our work is also conceptually related to rule-based control. Lu et al. (2022) use tree-search with  
 309 a heuristic to determine the quality of a given decoding path to steer decoding towards favorable  
 310 outcomes. Qin et al. (2022) explore gradient-based sampling using Langevin dynamics which  
 311 significantly outperforms gradient-free sampling.

312 **Reinforcement learning (RL).** Another line of very relevant work is reinforcement learning subject  
 313 to a KL penalty with the language model. Korbak et al. (2022) observed that reinforcement learning  
 314 with a KL penalty could be viewed in a Bayesian manner with a corresponding reward function.  
 315 However, their work fell short of making the full connection in an autoregressive decoding setting,  
 316 which is our contribution in this work through CD via a variant of the Bellman update akin to deep  
 317 Q-learning (DQN) (Mnih et al., 2013). Another closely related work to ours is that of Snell et al.  
 318 (2023) that designs a value-based offline algorithm, albeit with a different learning objective than  
 319 ours (and that of the KL-regularized PPO).

320 Other related RL work includes generator improvement solutions through on-policy RL. Sparrow  
 321 (Glaese et al., 2022) showed that a variant of proximal policy optimization (PPO) (Schulman  
 322 et al., 2017) with an additional LM regularizer is effective at a variety of safety objectives and  
 323 alignment with human preference (Ouyang et al., 2022).

324 **Supervised learning from negative examples.** Another line of related work is supervised generator  
 325 improvement interventions. These include unlikelihood training (Welleck et al., 2020; Zhang &  
 326 Song, 2022), contrastive losses (Adolphs et al., 2022), and direct preference optimization (Rafailov  
 327 et al., 2023). In contrast to our work, these methods are all training-time interventions but they  
 328 could similarly be used to improve the likelihood of drawing positive examples by suppressing the  
 329 likelihood of negative ones.

330 **A.2 Proof of Theorem 2.1**

331 *Proof of Theorem 2.1.* First notice that

$$J([\mathbf{x}, y^t]; \pi, \beta) = \sum_{z \in \mathcal{Y}} \pi(z | [\mathbf{x}, y^t]) \left( (1 - \beta)(\gamma V([\mathbf{x}, y^t, z]) - V([\mathbf{x}, y^t])) + \beta \log \left( \frac{p(z | [\mathbf{x}, y^t])}{\pi(z | [\mathbf{x}, y^t])} \right) \right) \quad (12)$$

$$= \beta \sum_{z \in \mathcal{Y}} \pi(z | [\mathbf{x}, y^t]) \log \left( \frac{p(z | [\mathbf{x}, y^t]) e^{\frac{1-\beta}{\beta}(\gamma V([\mathbf{x}, y^t, z]) - V([\mathbf{x}, y^t]))}}{\pi(z | [\mathbf{x}, y^t])} \right). \quad (13)$$

332 Now, let

$$q(z | [\mathbf{x}, y^t]; \beta) := \frac{p(z | [\mathbf{x}, y^t]) e^{\frac{(1-\beta)\gamma}{\beta} V([\mathbf{x}, y^t, z])}}{Z([\mathbf{x}, y^t]; \beta)}, \quad (14)$$

333 where

$$Z(\mathbf{x}, y^t; \beta) = \sum_{z \in \mathcal{Y}} p(z | \mathbf{x}, y^t) e^{\frac{(1-\beta)\gamma}{\beta} V(\mathbf{x}, y^t, z)}. \quad (15)$$

334 Thus,

$$J([\mathbf{x}, y^t]; \pi, \beta) = -\beta D(\pi(\cdot | [\mathbf{x}, y^t]) \| q(\cdot | [\mathbf{x}, y^t]; \beta)) + \beta \log Z([\mathbf{x}, y^t]; \beta), \quad (16)$$

335 which is maximized by

$$\pi(\cdot | [\mathbf{x}, y^t]) = q(\cdot | [\mathbf{x}, y^t]; \beta), \quad (17)$$

336 completing the proof.  $\square$



337 **A.3 Additional experimental results**

338 In this section, we provide some additional experimental results to better understand the prefix scorer  
 339 learnt via CD and FUDGE.

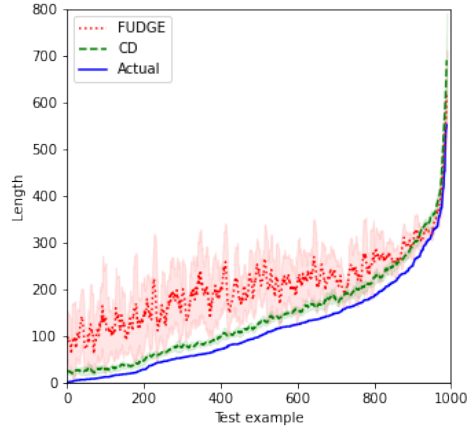


Figure 4: CD and FUDGE used to predict the length of a fully decoded response on Reddit corpus test set (Microsoft, 2019). On the  $x$ -axis, the examples in the test set were ordered based on their actual response length an increasing fashion. CD and FUDGE are applied to  $(x, y)$  pairs for all test set to predict the length. CD predictions are much better aligned with actual length, especially for pairwise comparison, whereas FUDGE predictions are noisy.

	FUDGE	CD
$\beta = 0$ (base)	1.0	1.0
$\beta = 0.10$	0.981	0.948
$\beta = 0.15$	0.959	0.961
$\beta = 0.20$	0.964	1.023
$\beta = 0.23$	0.926	0.990
$\beta = 2.00$	0.836	0.731

Table 2: Normalized average safety scores where the Reward-XS model (not used for alignment) is the judge, with 1000 generations from each model. The results are normalized to the average safety score of the base model. As can be seen, both prefix scorers generalize poorly which might be partly attributed to overoptimization.

	FUDGE	CD
$\beta = -0.50$	0.683	0.661
$\beta = -0.30$	0.792	0.745
$\beta = -0.20$	0.848	0.881
$\beta = 0$ (base)	1.0	1.0
$\beta = 0.100$	1.034	1.066
$\beta = 0.125$	1.009	1.007
$\beta = 0.150$	0.965	1.002
$\beta = 0.175$	0.984	1.021
$\beta = 0.200$	1.034	0.997
$\beta = 0.250$	1.011	1.04

Table 3: Normalized average safety scores where the Reward-XXS model (used for alignment) is the judge, with 1000 generations from each model. The results are normalized to the average safety score of the base model.