

VARIATIONAL BEST-OF- N ALIGNMENT

Anonymous authors

Paper under double-blind review

ABSTRACT

Best-of- N (Bo N) is a popular and effective algorithm for aligning language models to human preferences. The algorithm works as follows: at inference time, N samples are drawn from the language model, and the sample with the highest reward, as judged by a reward model, is returned as the output. Despite its effectiveness, Bo N is computationally expensive; it reduces sampling throughput by a factor of N . To make Bo N more efficient at inference time, one strategy is to fine-tune the language model to mimic what Bo N does during inference. To achieve this, we derive the distribution induced by the Bo N algorithm. We then propose to fine-tune the language model to minimize backward KL divergence to the Bo N distribution. Our approach is analogous to mean-field variational inference and, thus, we term it variational Bo N (vBo N). To the extent this fine-tuning is successful and we end up with a good approximation, we have reduced the inference cost by a factor of N . Our experiments on controlled generation and summarization tasks show that Bo N is the most effective alignment method, and our variational approximation to Bo N achieves the closest performance to Bo N and surpasses models fine-tuned using the standard KL-constrained RL objective. In the controlled generation task, vBo N appears more frequently on the Pareto frontier of reward and KL divergence compared to other alignment methods. In the summarization task, vBo N achieves high reward values across various sampling temperatures.

1 INTRODUCTION

Language models are pre-trained on large corpora to model a distribution over natural language text.¹ Beyond their initial pre-training, they are often additionally fine-tuned on domain-specific data through a process called **supervised fine-tuning (SFT)**. The goal of SFT is to enable the model to better perform various downstream tasks of interest. While the fine-tuned model, called the **reference model** in our paper, is indeed typically much better at performing the downstream task of interest, e.g., dialogue generation or summarization, it may still generate undesirable content, e.g., harmful or offensive text. To mitigate this issue, **aligning** the reference model to human preferences has become a fundamental step in the development of modern large language models (Touvron et al., 2023; OpenAI et al., 2023; Gemini et al., 2024).

The degree to which text is aligned with human preferences is typically operationalized using a real-valued reward function. Rather than constructing a reward function by hand, it is typically estimated from a dataset of human preferences.² And, after estimation, we expect the reward function to return higher values for text that is more likely to be preferred by humans, and lower values for text that is more likely to be dispreferred. Then, given an estimated reward function, an alignment algorithm further alters the reference models in a manner such that it places the highest probability on that text that is high reward under the reward model *and* high probability under the reference model.

Alignment algorithms can be taxonomized into two groups: (i) alignment via fine-tuning, where we change the language model’s parameters to achieve alignment (Christiano et al., 2017; Rafailov et al., 2023), and (ii) alignment via inference (Nakano et al., 2022; Mudgal et al., 2024). A common alignment-via-fine-tuning method is **reinforcement learning from human feedback (RLHF)**; (Christiano et al., 2017; Stiennon et al., 2020; Ouyang et al., 2022). RLHF typically consists of further

¹Many language models are also used to model text in non-natural languages, e.g., programming languages.

²In some cases, the reward model is not estimated from human preference data. It is either known, e.g., code-based execution scores, or given by a classifier, e.g., toxicity or sentiment classifiers.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

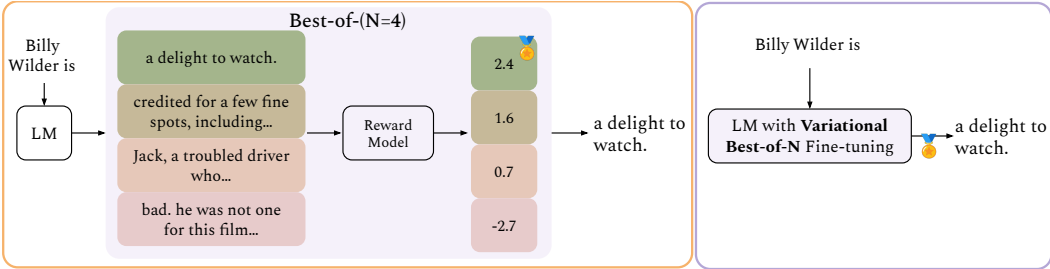


Figure 1: Best-of- N (on the left) is an effective alignment-via-inference method: it draws N samples from the language model, ranks them according to a reward model, and outputs the best sample. Variational Best-of- N (on the right) approximates this process via fine-tuning. The goal is to ensure that sampling a single string from the fine-tuned model produces a result equivalent to applying Best-of- N . This approach allows us to achieve similar performance while increasing the throughput by a factor of N .

fine-tuning the language model under a **KL-constrained RL objective**, which is made up of two terms: a term that encourages the model to maximize the reward, and a term that discourages high KL divergence between the language model and the reference model. This objective is often maximized with an RL algorithm, e.g., proximal policy optimization (PPO; Schulman et al., 2017). A common alignment-via-inference method is the Best-of- N (BoN; Stiennon et al., 2020) algorithm. As such, it does *not* require any fine-tuning of the language model. The algorithm is straightforward: One draws N samples from the reference model and returns the text that achieves the highest reward among those N samples. The BoN algorithm has also been effectively applied in controlled decoding (Yang & Klein, 2021; Mudgal et al., 2024) and to generate a dataset for supervised fine-tuning (Touvron et al., 2023).

Despite its simplicity, BoN has proven incredibly practical in generating high-reward text that still has a high probability under the reference model. Theoretically, Yang et al. (2024) prove that under some simplifying assumptions, the BoN distribution is asymptotically equivalent to the optimal distribution under the KL-constrained RL objective. Empirically, it has been repeatedly shown (Gao et al., 2023; Rafailov et al., 2023; Mudgal et al., 2024) that BoN often appears on the frontier of reward and KL curves, surpassing the performance of models fine-tuned with RLHF. However, the main factor preventing BoN from replacing fine-tuning methods for alignment is its significant computational overhead during inference. Even when sampling is done in parallel, BoN decreases the text generation throughput by a factor of N . This drawback limits its practicality for generating text from large language models.

To speed up BoN, we devise a scheme to convert it into an alignment-via-fine-tuning algorithm rather than an alignment-via-inference algorithm. To this end, we first formally derive the probability distribution induced by the BoN algorithm. Then we approximate this distribution by minimizing the reverse KL divergence between the language model and the BoN distribution. This leads to an optimization objective that we refer to as the vBoN objective. By analyzing a lower bound of this objective, we find that it behaves similarly to the KL-regularization objective in the limit, i.e., $N \rightarrow 1$ or $N \rightarrow \infty$. Importantly, the vBoN objective has a unique and useful property: it is insensitive to applying any monotonically increasing function to the reward values. This distinctive feature, along with the empirical success of the BoN algorithm, suggests that the vBoN objective is a promising and interesting objective to explore. Finally, we fine-tune the language model using PPO to optimize the vBoN objective. Our scheme, depicted in Fig. 1, allows us to achieve performance close to that of the BoN algorithm while increasing the inference throughput by a factor of N .

We experiment with our method on controlled generation and summarization tasks. We compare vBoN against models fine-tuned with the KL-constrained RL objective. In the controlled generation task, our results suggest that models fine-tuned with the vBoN objective are most likely to appear on the Pareto frontier of reward vs. KL curves compared to other alignment-via-finetuning methods, suggesting a better trade-off between attaining high rewards and not diverging too far from the reference model. Moreover, in the summarization task, we observe fine-tuning with vBoN leads to higher reward values and win rates on average compared to models fine-tuned with KL-constrained RL objective.

2 BACKGROUND: REINFORCEMENT LEARNING FROM HUMAN FEEDBACK

Let Σ be an **alphabet**, a finite, non-empty set of symbols.³ The elements of Σ may be characters, tokens, or words; the choice lies with the modeler. A **string** is a finite sequence of symbols drawn from Σ . A **language model** is a distribution over strings $\mathbf{y} \in \Sigma^*$, where Σ^* is the set of all strings over the alphabet Σ . In this paper, we consider language models, e.g., those based on neural networks, that are parameterized by a real vector $\theta \in \Theta$, denoted as π_θ . Furthermore, we restrict ourselves to language models that are differentiable functions of θ . In conditional generation tasks, e.g., summarization or dialogue generation, it is desirable to prompt the language model with a string $\mathbf{x} \in \Sigma^*$. Consequently, we consider prompted language models, i.e., those that give a conditional distribution over response strings \mathbf{y} , given a prompt string \mathbf{x} , as $\pi_\theta(\mathbf{y} \mid \mathbf{x})$. However, for notational convenience, we will drop the explicit conditioning on the prompt \mathbf{x} and simply write $\pi_\theta(\mathbf{y})$.

Algorithms for RLHF fine-tune the language model to increase the expected reward of the strings sampled from it while not diverging too far from the reference model. RLHF consists of three steps. First, the language model is fine-tuned on a task-specific dataset using the maximum-likelihood objective. Recall we term the language model after this step the reference model and show that with π_{ref} . Next, a **reward model** $r: \Sigma^* \rightarrow \mathbb{R}$ is trained to capture human preferences; the reward of a string is high if it is preferred by humans.⁴ Finally, the reference model is fine-tuned to maximize the KL-constrained RL objective,

$$\mathcal{J}^{\text{RL}}(\theta) = \mathbb{E}_{\mathbf{y} \sim \pi_\theta} [r(\mathbf{y})] - \beta D_{\text{KL}}(\pi_\theta \parallel \pi_{\text{ref}}), \quad (1)$$

where $D_{\text{KL}}(\cdot)$ is the KL divergence between two distributions, modulated by a hyperparameter β . This objective encourages the model to put more probability mass on strings that have high rewards under the reward model while penalizing it not to deviate too far from the reference model. Levine (2018) show that the optimal probability distribution that maximizes this objective is

$$\pi_\theta^*(\mathbf{y}) = \frac{1}{Z} \pi_{\text{ref}}(\mathbf{y}) \exp\left(\frac{1}{\beta} r(\mathbf{y})\right), \quad Z = \sum_{\mathbf{y} \in \Sigma^*} \pi_{\text{ref}}(\mathbf{y}) \exp\left(\frac{1}{\beta} r(\mathbf{y})\right). \quad (2)$$

π_θ^* is simply the reference model reweighted by the exponent of reward values and normalized by the partition function Z . Notably, we can not directly sample from π_θ^* because the partition function Z may be difficult to compute—it involves an infinite sum after all. However, a heuristic approach to sampling from π_θ^* would be to sample many strings from π_{ref} and only keep those that have high rewards. Indeed, this heuristic is the motivation behind the BoN algorithm.

3 DERIVING THE BEST-OF- N OBJECTIVE

Best-of- N algorithm is a simple alignment-via-inference algorithm. The algorithm works as follows. Let $Y_N = \{\mathbf{y}^{(n)}\}_{n=1}^N$ be the multi-set containing N i.i.d samples from π_{ref} . Then, BoN algorithm returns \mathbf{y}^* , where⁵

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}^{(n)} \in Y_N} r(\mathbf{y}^{(n)}). \quad (3)$$

We show the probability distribution induced from BoN sampling algorithm with π_{bon} . Importantly, π_{bon} is *not* the optimal distribution under Eq. (1), the KL-constrained RL objective.⁶ Nevertheless, the BoN algorithm often performs well—even compared to RLHF-based methods. This raises the question: under what optimization objective is π_{bon} the optimal distribution? To derive such an objective, we begin by computing the probability of strings under π_{bon} .

³Please refer to Tab. 3 for a summary of notations used throughout the paper.

⁴For example, in a summarization task, a preference dataset consists of a document, two candidate summaries for that document, and a label indicating which summary is preferred by humans. The reward model is trained on this dataset to maximize the likelihood of correctly predicting human preference.

⁵We assume that the argmax is unique, or ties are broken in a well-defined manner.

⁶Note that only under some simplifying assumptions, π_{bon} is asymptotically (in sequence length) equal to π_θ^* (Yang et al., 2024).

Proposition 1. Suppose $r: \Sigma^* \rightarrow \mathbb{R}$ is a one-to-one mapping. Then, the probability that a string $\mathbf{y} \sim \pi_{\text{bon}}$ is given by

$$\pi_{\text{bon}}(\mathbf{y}) = \sum_{i=1}^N \binom{N}{i} F(r(\mathbf{y}))^{N-i} \pi_{\text{ref}}(\mathbf{y})^i, \quad F(r(\mathbf{y})) \stackrel{\text{def}}{=} \mathbb{P}_{\mathbf{y}' \sim \pi_{\text{ref}}} (r(\mathbf{y}') < r(\mathbf{y})). \quad (4)$$

Proof. See App. B. ■

F can be understood as the strict cumulative density function of reward values under π_{ref} . In other words, $F(r(\mathbf{y}))$ represents the probability that a random sample drawn from π_{ref} has a reward value less than $r(\mathbf{y})$. We now describe how to fine-tune the language model to approximate π_{bon} . Similar to mean-field variational inference, we minimize the reverse KL divergence between π_{θ} and π_{bon} . Concretely,

$$\begin{aligned} \mathcal{J}^{\text{vBoN}}(\theta) &= -D_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{bon}}) = \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \left[\log \pi_{\text{bon}}(\mathbf{y}) - \log \pi_{\theta}(\mathbf{y}) \right] \\ &= \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \left[\log \pi_{\text{bon}}(\mathbf{y}) \right] + H(\pi_{\theta}) \\ &= \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \left[\log \sum_{i=1}^N \binom{N}{i} F(r(\mathbf{y}))^{N-i} \pi_{\text{ref}}(\mathbf{y})^i \right] + H(\pi_{\theta}), \quad (5) \end{aligned}$$

where $H(\cdot)$ is the entropy of a distribution. Eq. (5) is an entropy-regularized objective, where we use the probability of the string under the BoN distribution as the reward and discourage the model from having low entropy.

Monotonically invariant. An important property of the variational BoN objective is that it is invariant to applying any strictly monotonically increasing function to rewards. This is because the vBoN objective relies on reward values solely through F , which, as defined in Eq. (4), only depends on the ranking between the reward values and not their exact magnitude. This implies that the vBoN objective is insensitive to the outliers and the scale of rewards. This property is especially important as RL algorithms are notoriously sensitive to the scale of reward values (Henderson et al., 2018; Schaul et al., 2021).

Approximating $\log F(\cdot)$. Maximizing Eq. (5) requires us to compute $\log F(\cdot)$ for any $r(\mathbf{y})$. This, however, is computationally expensive, as we have to sum over the probabilities of all strings that have rewards less than $r(\mathbf{y})$. Fortunately, we can instead maximize a lower bound of Eq. (5) using a Monte Carlo estimator of $F(\cdot)$. Concretely, we can write $F(\cdot)$ as an expectation,

$$F(r(\mathbf{y})) = \mathbb{E}_{\mathbf{y}' \sim \pi_{\text{ref}}} [\mathbb{1}\{r(\mathbf{y}') < r(\mathbf{y})\}]. \quad (6)$$

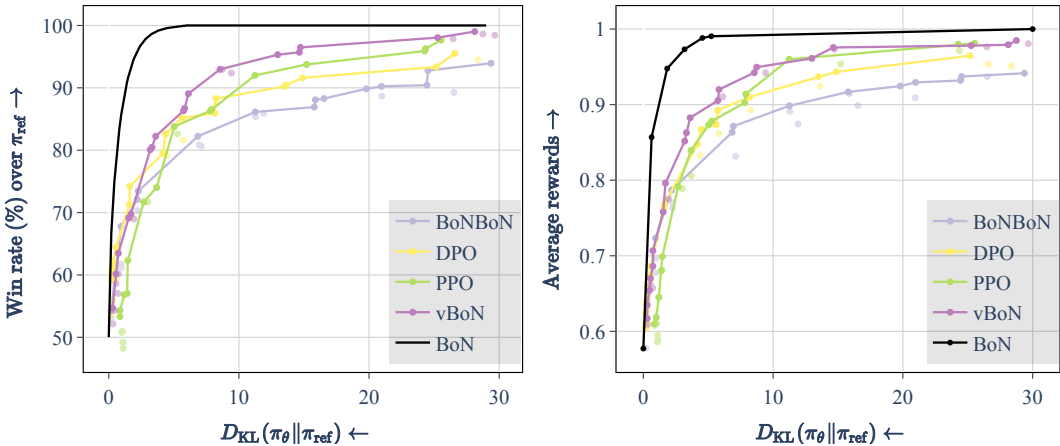
We approximate $F(r(\mathbf{y}))$ using M i.i.d. samples from π_{ref} , termed $\mathbf{y}'^{(1)}, \dots, \mathbf{y}'^{(M)}$ i.i.d. π_{ref} , and $\hat{F}(r(\mathbf{y})) \stackrel{\text{def}}{=} \frac{1}{M} \sum_{m=1}^M \mathbb{1}\{r(\mathbf{y}'^{(m)}) < r(\mathbf{y})\}$. We then take log of this Monte Carlo estimator as a biased, but consistent estimator of $\log F(\cdot)$ in Eq. (5).⁷ In §5.1 we empirically assess the number of samples we need so that $\log \hat{F}$ converges to $\log F$.

⁷Using Jensen’s inequality, we show biasedness. Concretely, note the following lower bound

$$\log F(r(\mathbf{y})) = \log \mathbb{E}_{\mathbf{y}'^{(1)}, \dots, \mathbf{y}'^{(M)}} \left[\frac{1}{M} \sum_{m=1}^M \mathbb{1}\{r(\mathbf{y}'^{(m)}) < r(\mathbf{y})\} \right] \quad (7a)$$

$$\geq \mathbb{E}_{\mathbf{y}'^{(1)}, \dots, \mathbf{y}'^{(M)}} \left[\log \left(\frac{1}{M} \sum_{m=1}^M \mathbb{1}\{r(\mathbf{y}'^{(m)}) < r(\mathbf{y})\} \right) \right], \quad (7b)$$

where Jensen’s inequality is applicable because log is concave. Consistency can be shown with an application of the delta method (§5.5.4; Casella & Berger, 2001).



(a) 4% of points on Pareto front belong to BoNBoN, 4% to PPO, 42% to DPO, and 50% to vBoN. (b) 7% of points on Pareto front belong to BoNBoN, 10% DPO, 33% PPO, and 50% vBoN.

Figure 2: Steering generated movie reviews towards positive sentiment. Points that are not on the Pareto front of each method have lower opacity. BoN is the most effective approach in achieving high win rates and high rewards while not diverging too far from the reference model. Our variational approximation to BoN gets closest to the performance of BoN compared to other fine-tuning methods, as reflected in the percentage of times it appears on the Pareto front.

4 COMPARING BO N AND RL OBJECTIVES

To explore the connection between the vBoN objective and the KL-regularized RL objective, we derive a lower bound for $\mathcal{J}^{\text{vBoN}}$. Through this lower bound, we can get more insights on how the reward function is used in the variational BoN objective, and why this objective discourages high KL divergence from the reference model.

To derive such a lower bound, we substitute the BoN distribution in Eq. (4) into the vBoN objective in Eq. (5). We then simplify this objective to arrive at the following theorem.

Theorem 2. We have $\mathcal{J}^{\text{vBoN}}(\theta) \geq L(\theta)$, where

$$L(\theta) \stackrel{\text{def}}{=} (N - 1) \mathbb{E}_{\mathbf{y} \sim \pi_\theta} \left[\log F(r(\mathbf{y})) \right] - D_{\text{KL}}(\pi_\theta \parallel \pi_{\text{ref}}). \tag{8}$$

Proof. See App. D. ■

Empirically, we observe that models that are fine-tuned to maximize $L(\theta)$ perform competitively to the ones that are fine-tuned to maximize the vBoN objective; see App. G for experimental results. Interestingly, if we compare Eq. (8) to the KL-constrained RL objective, Eq. (1), we see they have a very similar structure. We observe that N (in the vBoN objective) acts as a regularization parameter. As $N \rightarrow 1$, the optimal distribution gets closer to π_{ref} , which has the same effect as $\beta \rightarrow \infty$ in Eq. (1). Furthermore, as $N \rightarrow \infty$, the optimal distribution only generates the string with the maximum rewards, which is equivalent to $\beta \rightarrow 0$ in Eq. (1). Importantly, in both limits, the optimal distribution under the KL-regularized RL objective and the vBoN objective are equivalent.

The main difference between the KL-constrained RL objective Eq. (1) and the derived vBoN lower bound Eq. (8) is in the reward function. With the KL-constrained RL objective, we aim to maximize the expected reward values. In contrast, with vBoN, we maximize the cumulative probability that strings sampled from the aligned model, π_θ , achieve higher rewards compared to those sampled from π_{ref} .

270 5 SENTIMENT CONTROL

271
272 We now employ the variational BoN objective, Eq. (5), to fine-tune language models. We perform an
273 open-ended text generation task where the goal is to generate movie reviews with positive sentiment.

274 The reference model, π_{ref} , is GPT-IMDB⁸, a GPT-2 (Radford et al., 2019) model fine-tuned on IMDB
275 corpus (Maas et al., 2011). We use a binary sentiment classifier,⁹ denoted as p , with two classes
276 $\{\text{POS}, \text{NEG}\}$ as the reward model, and define $r(\mathbf{y}) \stackrel{\text{def}}{=} p(\text{POS} \mid \mathbf{y})$. Following Rafailov et al. (2023), we
277 sample 5000 movie reviews from the training set of IMDB dataset and for each sample, we randomly
278 choose a prefix length between 2 – 8 and take that prefix as the prompt. We further generate 512
279 prompts in the same way from the test set of IMDB that we use to evaluate our models.
280

281 We fine-tune the reference model with PPO using the vBoN objective Eq. (5). Then, we compare
282 the performance of the fine-tuned model (vBoN) to the exact BoN (BoN), i.e., applying BoN at
283 inference time.

284 We implement and compare the following existing methods for language model alignment:

- 285 • **BoN-SFT**: Perhaps the most straightforward way to approximate BoN distribution is to fine-tune
286 the model to maximize the likelihood of the samples taken with BoN algorithm. Unfortunately, we
287 find that SFT is incapable of achieving a good trade-off between achieving high rewards and low
288 KL divergence, see App. H (Fig. 7) for the experimental results.
- 289 • **PPO**: We use PPO to optimize the KL-constrained objective in Eq. (1). We use the default
290 hyperparameters in trlx library (Havrilla et al., 2023) for fine-tuning with PPO.
- 291 • **DPO**. Direct preference optimization (DPO; Rafailov et al., 2023) is a popular alternative to RLHF
292 that does not require training a reward model. Following DPO’s experimental setup, we generate 6
293 reviews per prompt and use the resulting 12 pairwise comparisons per prompt to construct DPO’s
294 contrastive loss.¹⁰
- 295 • **BoNBoN**: Concurrent work (Gui et al., 2024) explores another approach to approximate BoN
296 distribution. Assuming that the reference model distribution π_{ref} is continuous, Gui et al. (Theorem
297 3; 2024) prove that the expected difference between the relative likelihood, i.e., $\frac{\pi_{\text{bon}}(\cdot)}{\pi_{\text{ref}}(\cdot)}$, of the
298 Best-of- N response and the Worst-of- N response is $\frac{1}{2\beta} = \frac{1}{(N-1) \sum_{k=1}^{N-1} 1/k}$. They use this property
299 to construct a loss function similar to that of IPO (Azar et al., 2023). Furthermore, they add another
300 term to the loss function, which simply maximizes the likelihood of the Best-of- N response. The
301 final loss function is a convex combination of the IPO-like loss and the negative log-likelihood loss,
302 regulated by a hyperparameter α .¹¹
303
304

305 We fine-tune models by varying the degree of regularization. For BoN approaches, that is achieved
306 by varying N , and for DPO and PPO, we vary β .¹² Conveniently, N in vBoN is a hyperparameter,
307 meaning that we do *not* need to generate more samples from π_{ref} when we increase N . However,
308 with BoN and BoNBoN methods, we need to increase the number of samples from the reference
309 model as we increase N .

310 We generate movie reviews based on prompts from our test set using fine-tuned models and then
311 measure three metrics: (i) KL divergence between the fine-tuned model and the reference model; (ii)
312 win rate, defined as the percentage of times the fine-tuned model’s generations receive higher rewards
313 compared to the reference model’s generations; and (iii) average rewards obtained by the fine-tuned
314 model’s sampled strings.

315 For the BoN method, we report the empirical upper bound of $\log N - \frac{N-1}{N}$ for KL divergence
316 (Beirami et al., 2024; Mroueh, 2024) in our plots. Furthermore, the win rate of BoN over the
317 reference model can be computed analytically and is equal to $\frac{N}{N+1}$.

318 ⁸Specifically, we use <https://huggingface.co/lvwerra/gpt2-imdb>.

319 ⁹Specifically, we use <https://huggingface.co/lvwerra/distilbert-imdb>.

320 ¹⁰One could argue that DPO has a slight advantage over other methods in this setup since it has seen 6 unique
321 generations per prompt during training, while the others only have seen one (or 2 with BoNBoN). Nevertheless,
322 we observe that vBoN is more effective than DPO.

323 ¹¹Following the authors’ recommendation, we set α so that both terms contribute equally to the final loss.

¹²See App. F for more details regarding the regularization hyperparameters.

We visualize the win rate vs. KL curves in Fig. 2a, and Fig. 2b the average rewards of generations under π_θ vs. the KL divergence. As expected, BoN is the most effective approach; however, this comes at an extra inference cost that grows with N . We observe that among the fine-tuning methods, our variational approximation to BoN gets closest to the performance of BoN, as it appears more often on the Pareto front of the two curves compared to other methods. Notably, we observe that DPO performs better than PPO in terms of win rates but worse in terms of average rewards; this could be attributed to the contrastive nature of DPO’s loss function.

5.1 ERROR IN ESTIMATING $\log F(\cdot)$

We empirically quantify the error when estimating $\log F(\cdot)$ with a finite number of i.i.d samples from π_{ref} . To get a better intuition on the error of our estimators, in Fig. 3, we visualize the estimators for 3 different prompts: one adversarial prompt (left plot), where the prompt itself has a negative sentiment, one neutral prompt (middle plot), and one prompt with a positive sentiment (right plot). We vary the number of Monte Carlo samples from 10 to 600. We observe that for all the 3 prompts, the estimated CDF hardly changes after 200 samples. When using the adversarial prompt, the reward distribution is negatively peaked, and the estimated CDF does not change after taking only 100 samples.

We then quantify the change in the estimator by performing a two-sample Kolmogorov–Smirnov test (Hodges, 1958). This test measures the closeness of two empirical cumulative distribution functions. Concretely, the test statistic is

$$\sup_{\mathbf{y} \in \Sigma^*} \left| \widehat{F}_{M_1}(r(\mathbf{y})) - \widehat{F}_{M_2}(r(\mathbf{y})) \right|, \quad (9)$$

where $\widehat{F}_{M_1}, \widehat{F}_{M_2}$ are estimated CDFs from M_1 and M_2 samples respectively. The statistics show the magnitude of the difference between the two empirical distributions of samples. The null hypothesis is that the two distributions are identical.

In Tab. 1, for each sample size M , we compare the estimated CDF with M samples to the estimated CDF with 600 samples. If the two distributions are identical according to the test, we can reliably use the M sample to estimate the CDF. We report the number of prompts (out of 5000 prompts) for which we reject the null hypothesis, meaning that the distributions are not identical. Furthermore, for those prompts, we report the average test statistics and p -values. In general, for very few prompts, the null hypothesis is rejected. Moreover, with 250 samples, the estimated CDFs are identical to the estimated CDF with 600 samples for all prompts.

Table 1: Measuring the estimation error with increasing the sample size. After 250 samples, the estimated CDF is unchanged for all the prompts.

M	Rejection rate	Test statistics	p -value
5	6.14%	0.63	0.02
20	4.02%	0.33	0.03
100	1.14%	0.17	0.02
200	0.06%	0.12	0.02
250	0	-	-

5.2 EFFICIENCY ANALYSIS

We break down the efficiency analysis into 3 main parts: (i) the inference cost, (ii) the preference optimization cost, (iii) and the preprocessing cost.

Inference Cost. As discussed earlier, vBoN is an alignment-via-finetuning method, and along with other alignment-via-finetuning methods, it is N times more efficient at inference compared to BoN.

Optimization Cost. We compare vBoN’s preference optimization cost to its closest alignment-via-finetuning counterpart, PPO. In the optimization loop, the main difference between PPO and vBoN is that vBoN requires computing the strict CDF function, F , using M samples. Crucially, N in vBoN serves as a regularization hyperparameter, and increasing N does *not* incur additional computation costs. To implement vBoN efficiently, we precompute the F function before starting the optimization loop. This means the computational overhead is incurred only once, regardless of the number of optimization runs.¹³ Since the F values are precomputed, we empirically observe that the

¹³This is particularly advantageous since practitioners often perform the optimization multiple times to test various hyperparameter settings.

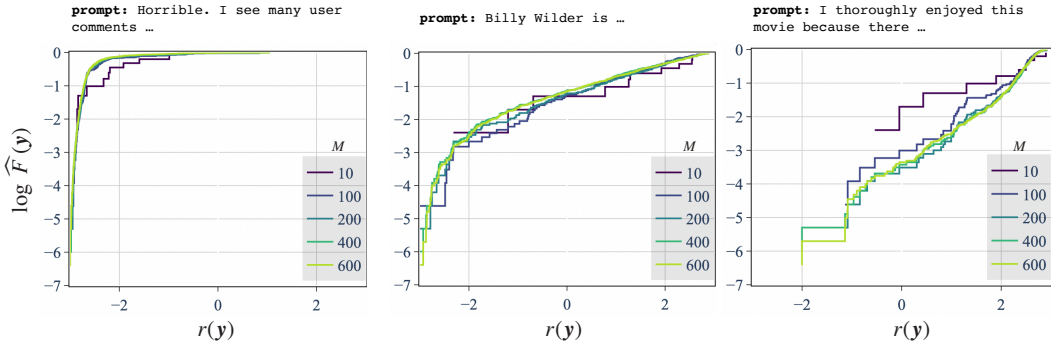


Figure 3: Estimates of $\log \hat{F}(\cdot)$ with increasing the number of Monte Carlo samples. We test an adversarial prompt (left plot), a neutral prompt (middle plot), and a prompt with a positive sentiment (right plot). Overall, we hardly see any difference between the estimates after taking 200 samples. For the adversarial prompt, the distribution of rewards is peaked, and we do not see any changes in our estimator after taking only 100 samples.

time needed to run the vBoN optimization loop is the same as running the PPO optimization loop, and the cost of evaluating F is negligible. Therefore, the main computational overhead in vBoN comes from precomputing $\log F(\cdot)$.

Preprocessing Cost. Estimating $\log F(\cdot)$ requires only forward passes through the LLM and reward model, without the need to compute and store gradients. This makes the process highly parallelizable. In our experiments, we utilize a memory-efficient library for LLM inference, (vLLM; Kwon et al., 2023), which allows us to perform these approximations efficiently.

We examine the impact of increasing the computational cost of vBoN by varying M , which directly affects the total elapsed time and downstream performance. For this analysis, we fix $N = 10$ and fine-tune the model using three random seeds. We report the average and standard deviation of reward values and win rates in Fig. 4 on a single A100-40GB GPU. Our results show that increasing M generally improves the aligned model’s rewards and win rates. Notably, even with $M = 32$ samples (taking only 10 minutes), the performance remains competitive with higher values of M . We hypothesize that the data efficiency of the simple Monte Carlo estimator can be greatly improved by taking into account the similarity between different prompts to learn an approximation to $\log F$ function, which we plan as future work.

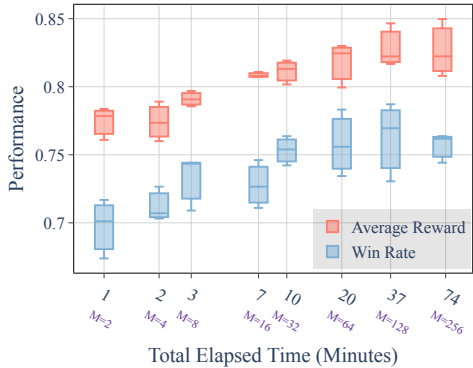


Figure 4: The average reward and win rate of the aligned models improve as we increase the sample size M used for approximating the vBoN loss function.

6 SUMMARIZATION

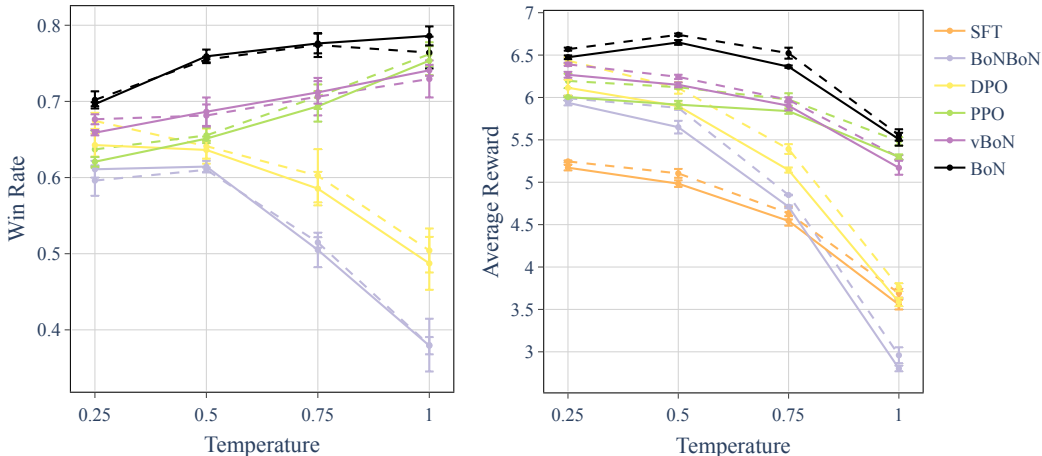
We further employ variational BoN in a summarization task, where the goal is to generate summaries that align with human preferences. The reference model, π_{ref} , is a pythia-2.8B model fine-tuned on human-written summaries of Reddit posts Stiennon et al. (2020).¹⁴ We use SFT to refer to this model in the plots. We use two separate reward models for training and evaluation: a pythia-2.8B¹⁵ reward model for fine-tuning and a larger pythia-6.9B¹⁶ model exclusively for evaluation.

¹⁴We use https://huggingface.co/cleanrl/EleutherAI_pythia-2.8b-deduped_sft_tldr.

¹⁵We use https://huggingface.co/cleanrl/EleutherAI_pythia-2.8b-deduped_reward_tldr.

¹⁶We use https://huggingface.co/cleanrl/EleutherAI_pythia-6.9b-deduped_reward_tldr.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
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
485



(a) Comparing the win rates of alignment methods against samples from the π_{ref} . vBoN achieves closer results to BoN compared to other alignment-via-finetuning methods. (b) Comparing the average rewards obtained from the evaluator reward model. BoN outperforms other alignment methods, and vBoN achieves closer results to BoN compared to other alignment-via-finetuning methods.

Figure 5: Performance of different alignment methods on the summarization task. Solid traces show the performance on in-distribution Reddit posts, while dashed lines demonstrate the out-of-distribution performance. Overall, BoN is the most effective approach in achieving high win rates and average rewards across all sampling temperatures. Our variational approximation to BoN (vBoN) gets closest to the performance of BoN, while being significantly cheaper at inference time.

Dataset. To evaluate the generalization ability of the aligned models on out-of-distribution data, we fine-tune the models using only posts from the relationship and relationship_advice subreddits of the Reddit TL;DR (Stiennon et al., 2020) dataset. We then assess the models’ performance on the two types of data by dividing the the test set into two equally-sized groups: in-distribution Reddit posts from the relationship and relationship_advice subreddits, and out-of-distribution posts from the rest of the subreddits. We visualize the performance of methods on in-distribution data with a solid trace and on out-of-distribution data with a dashed trace.

Experimental setup. We fine-tune the model with both the KL-constrained RL objective and vBoN objective for 10000 episodes. Similar to the previous experiment, we use 200 samples to estimate $\log F(\cdot)$ values. To create a smooth and continuous reward function, we further fit an exponential curve¹⁷ to the estimates. We set $N = 100$ for BoN and vBoN methods and the equivalent value of $\beta = 0.05$ for the KL-constrained RL objective. We closely follow Huang et al. (2024) for setting the hyperparameters of the PPO algorithm; please refer to App. F for more experimental details. After fine-tuning, we sample from the aligned models with different sampling temperatures $t \in [0.25, 0.5, 0.75, 1.]$, each with 3 different random seeds.

Win rates. In Fig. 5a we visualize the average and standard deviation of win rates compared against the samples from the SFT model. Notably, BoN achieves the highest win rates, which is consistent with findings from previous studies (Rafailov et al., 2023). We do not observe any significant differences between BoN performance on in-distribution (solid trace) and out-of-distribution data,¹⁸ which is expected as BoN is an alignment-via-inference method. Similarly, we mostly do not observe significant differences between in- and out-of-distribution performance of , indicating that these methods can generalize effectively in this experimental setup. DPO and BoNBoN only manage to perform competitively to other methods at lower temperatures (0.25, 0.5) and their performance drops

¹⁷We fit an exponential function of the form $f(x) = -a \exp(-bx)$ to the data using non-linear least squares method.

¹⁸The difference between the two data distributions becomes more apparent at temperature 1, potentially due to increased sample diversity in this setting.

Table 2: An example of summaries sampled at temperature 0.5 and their corresponding reward obtained from the evaluator reward model.

Content	Reward
SUBREDDIT: r/relationship_advice TITLE: Stuck in a rut and in need of advice/inspiration! POST: My boyfriend and I have been together for 3 years, and living together for 2. I'm quite the homebody, and when we first met, he was very outgoing and loved partying and socialising (although he was a student at the time). We're both working now, and most nights we find ourselves doing the same things: watching series (luckily we enjoy the same shows), playing Minecraft or playing various board games. We're tired after work, and can't bring ourselves to leave the house. The weekend is much the same – lots of sleep, or sitting around staring at one screen or another. We do party occasionally (we'll head to a pub once every few months) and there are a few mutual friends we enjoy spending time with, but I worry that we've become stuck in our boring ways. I really enjoy our lifestyle, and would be quite happy to never leave the house again, but I'm starting to feel guilty for turning him into a 50 year-old when he's only 24. Any ideas for shaking things up a little? Bear in mind that we live in a small town in South Africa, and neither of us has a car.	-
SFT: I'm stuck in a rut, and need to shake things up to see if it'll work out. Any advice?	3.08
PPO: In need of inspiration to break out of rut and live life fully! Any ideas welcome!	4.59
vBoN: Been happily living together for 2yr+, feeling bored after work regularly, looking for ideas to spice things up!	6.79
BoN: My boyfriend and I have been together for 3 years, and are both working full time. We spend most of our time in the house, and have become boring. What can we do to shake things up?	9.18

significantly at higher temperatures (0.75, 1). Importantly, while PPO and vBoN perform comparably at higher temperatures, vBoN significantly outperforms PPO at lower temperatures (0.25 and 0.5).

Average rewards. In Fig. 5b, we measure the average rewards across different temperatures. As the temperature increases, the average reward decreases consistently across all methods. This trend is also evident in the qualitative analysis in App. I, where we show sampled summaries at different temperatures. DPO and BoNBoN suffer more from increasing the temperature, as the average rewards get close to (or even worse than) the SFT average rewards. Generally, the average reward results align with the win-rate trends, and we observe that vBoN achieves significantly higher rewards compared to PPO at lower temperatures. In Tab. 2 we show an example of summaries generated from the fine-tuned models with their associated reward values.

7 CONCLUSION

Motivated by the effectiveness of the BoN algorithm, we formally derive a variational approximation to the distribution induced by BoN algorithm via fine-tuning language models. Our analysis highlights the similarities and distinctions between the variational BoN objective and the KL-constrained RL objectives. Our empirical findings reveal that models fine-tuned using the variational approximation to BoN not only attain high reward values but also maintain proximity to the reference models. Crucially, inference on the fine-tuned models with the vBoN objective remains as cost-effective as inference on the original reference model.

REFERENCES

- Mohammad Gheshlaghi Azar, Mark Rowland, Bilal Piot, Daniel Guo, Daniele Calandriello, Michal Valko, and Rémi Munos. A general theoretical paradigm to understand learning from human preferences. *Computing Research Repository*, arXiv:2310.12036, 2023. URL <https://arxiv.org/abs/2310.12036>.
- Ahmad Beirami, Alekh Agarwal, Jonathan Berant, Alexander D'Amour, Jacob Eisenstein, Chirag Nagpal, and Ananda Theertha Suresh. Theoretical guarantees on the best-of-n alignment policy.

- 540 *Computing Research Repository*, arXiv:2401.01879, 2024. URL <https://arxiv.org/abs/2401.01879>.
- 541
- 542
- 543 Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V. Le, Christopher Ré, and
- 544 Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling.
- 545 *Computing Research Repository*, arXiv:2407.21787, 2024. URL <https://arxiv.org/abs/2407.21787>.
- 546
- 547 George Casella and Roger L. Berger. *Statistical Inference*. Chapman and Hall/CRC, Pacific
- 548 Grove, CA, 2nd edition, 2001. ISBN 9781032593036. URL <https://www.routledge.com/Statistical-Inference/Casella-Berger/p/book/9781032593036>.
- 549
- 550 Eugene Charniak and Mark Johnson. Coarse-to-fine n-best parsing and MaxEnt discriminative
- 551 reranking. In Kevin Knight, Hwee Tou Ng, and Kemal Oflazer (eds.), *Proceedings of the 43rd*
- 552 *Annual Meeting of the Association for Computational Linguistics (ACL’05)*, pp. 173–180, Ann
- 553 Arbor, Michigan, June 2005. Association for Computational Linguistics. doi: 10.3115/1219840.
- 554 1219862. URL <https://aclanthology.org/P05-1022>.
- 555
- 556 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep rein-
- 557 forcement learning from human preferences. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach,
- 558 R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing*
- 559 *Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/d5e2c0adad503c91f91df240d0cd4e49-Paper.pdf.
- 560
- 561 Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan Zhang, Winnie Chow, Rui Pan, Shizhe Diao,
- 562 Jipeng Zhang, KaShun SHUM, and Tong Zhang. RAFT: Reward ranked finetuning for generative
- 563 foundation model alignment. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856.
- 564 URL <https://openreview.net/forum?id=m7p507zb1Y>.
- 565
- 566 Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization.
- 567 In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and
- 568 Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*,
- 569 volume 202 of *Proceedings of Machine Learning Research*, pp. 10835–10866. PMLR, 23–29 Jul
- 570 2023. URL <https://proceedings.mlr.press/v202/gao23h.html>.
- 571
- 572 Gemini, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan
- 573 Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson, Ioannis
- 574 Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillicrap,
- 575 Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul R. Barham, Tom Hennigan,
- 576 Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan Doherty, Eli Collins,
- 577 Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, Jack Krawczyk,
- 578 Cosmo Du, Ed Chi, Heng-Tze Cheng, Eric Ni, Purvi Shah, Patrick Kane, Betty Chan, Manaal
- 579 Faruqui, Aliaksei Severyn, Hanzhao Lin, YaGuang Li, Yong Cheng, Abe Ittycheriah, Mahdis
- 580 Mahdieh, Mia Chen, Pei Sun, Dustin Tran, Sumit Bagri, Balaji Lakshminarayanan, Jeremiah Liu,
- 581 Andras Orban, Fabian Gura, Hao Zhou, Xinying Song, Aurelien Boffy, Harish Ganapathy, Steven
- 582 Zheng, HyunJeong Choe, goston Weisz, Tao Zhu, Yifeng Lu, Siddharth Gopal, Jarrod Kahn,
- 583 Maciej Kula, Jeff Pitman, Rushin Shah, Emanuel Taropa, Majd Al Meray, Martin Baeuml, Zhifeng
- 584 Chen, Laurent El Shafey, Yujing Zhang, Olcan Sercinoglu, George Tucker, Enrique Piqueras,
- 585 Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs, Anas White, Anders
- 586 Andreassen, Tamara von Glehn, Lakshman Yagati, Mehran Kazemi, Lucas Gonzalez, Misha
- 587 Khalman, Jakub Sygnowski, Alexandre Frechette, Charlotte Smith, Laura Culp, Lev Proleev,
- 588 Yi Luan, Xi Chen, James Lottes, Nathan Schucher, Federico Lebron, Alban Rrustemi, Natalie
- 589 Clay, Phil Crone, Tomas Kocisky, Jeffrey Zhao, Bartek Perz, Dian Yu, Heidi Howard, Adam
- 590 Bloniarz, Jack W. Rae, Han Lu, Laurent Sifre, Marcello Maggioni, Fred Alcober, Dan Garrette,
- 591 Megan Barnes, Shantanu Thakoor, Jacob Austin, Gabriel Barth-Maron, William Wong, Rishabh
- 592 Joshi, Rahma Chaabouni, Deeni Fatiha, Arun Ahuja, Gaurav Singh Tomar, Evan Senter, Martin
- 593 Chadwick, Ilya Kornakov, Nithya Attaluri, Iñaki Iturrate, Ruibo Liu, Yunxuan Li, Sarah Cogan,
- Jeremy Chen, Chao Jia, Chenjie Gu, Qiao Zhang, Jordan Grimstad, Ale Jakse Hartman, Xavier
- Garcia, Thanumalayan Sankaranarayanan Pillai, Jacob Devlin, Michael Laskin, Diego de Las Casas,
- Dasha Valter, Connie Tao, Lorenzo Blanco, Adrià Puigdomènech Badia, David Reitter, Mianna
- Chen, Jenny Brennan, Clara Rivera, Sergey Brin, Shariq Iqbal, Gabriela Surita, Jane Labanowski,

594 Abhi Rao, Stephanie Winkler, Emilio Parisotto, Yiming Gu, Kate Olszewska, Ravi Addanki,
595 Antoine Miech, Annie Louis, Denis Teplyashin, Geoff Brown, Elliot Catt, Jan Balaguer, Jackie
596 Xiang, Pidong Wang, Zoe Ashwood, Anton Briukhov, Albert Webson, Sanjay Ganapathy, Smit
597 Sanghavi, Ajay Kannan, Ming-Wei Chang, Axel Stjerngren, Josip Djolonga, Yuting Sun, Ankur
598 Bapna, Matthew Aitchison, Pedram Pejman, Henryk Michalewski, Tianhe Yu, Cindy Wang, Juliette
599 Love, Junwhan Ahn, Dawn Bloxwich, Kehang Han, Peter Humphreys, Thibault Sellam, James
600 Bradbury, Varun Godbole, Sina Samangooei, Bogdan Damoc, Alex Kaskasoli, Sébastien M. R.
601 Arnold, Vijay Vasudevan, Shubham Agrawal, Jason Riesa, Dmitry Lepikhin, Richard Tanburn,
602 Srivatsan Srinivasan, Hyeontaek Lim, Sarah Hodgkinson, Pranav Shyam, Johan Ferret, Steven Hand,
603 Ankush Garg, Tom Le Paine, Jian Li, Yujia Li, Minh Giang, Alexander Neitz, Zaheer Abbas, Sarah
604 York, Machel Reid, Elizabeth Cole, Aakanksha Chowdhery, Dipanjan Das, Dominika Rogozińska,
605 Vitaliy Nikolaev, Pablo Sprechmann, Zachary Nado, Lukas Zilka, Flavien Prost, Luheng He,
606 Marianne Monteiro, Gaurav Mishra, Chris Welty, Josh Newlan, Dawei Jia, Miltiadis Allamanis,
607 Clara Huiyi Hu, Raoul de Liedekerke, Justin Gilmer, Carl Saroufim, Shruti Rijhwani, Shaobo Hou,
608 Disha Shrivastava, Anirudh Baddepudi, Alex Goldin, Adnan Ozturel, Albin Cassirer, Yunhan Xu,
609 Daniel Sohn, Devendra Sachan, Reinald Kim Amplayo, Craig Swanson, Dessie Petrova, Shashi
610 Narayan, Arthur Guez, Siddhartha Brahma, Jessica Landon, Miteyan Patel, Ruizhe Zhao, Kevin
611 Villela, Luyu Wang, Wenhao Jia, Matthew Rahtz, Mai Giménez, Legg Yeung, James Keeling,
612 Petko Georgiev, Diana Mincu, Boxi Wu, Salem Haykal, Rachel Saputro, Kiran Vodrahalli, James
613 Qin, Zeynep Cankara, Abhanshu Sharma, Nick Fernando, Will Hawkins, Behnam Neyshabur,
614 Solomon Kim, Adrian Hutter, Priyanka Agrawal, Alex Castro-Ros, George van den Driessche,
615 Tao Wang, Fan Yang, Shuo yiin Chang, Paul Komarek, Ross McIlroy, Mario Lučić, Guodong
616 Zhang, Wael Farhan, Michael Sharman, Paul Natsev, Paul Michel, Yamini Bansal, Siyuan Qiao,
617 Kris Cao, Siamak Shakeri, Christina Butterfield, Justin Chung, Paul Kishan Rubenstein, Shivani
618 Agrawal, Arthur Mensch, Kedar Soparkar, Karel Lenc, Timothy Chung, Aedan Pope, Loren
619 Maggiore, Jackie Kay, Priya Jhakra, Shibo Wang, Joshua Maynez, Mary Phuong, Taylor Tobin,
620 Andrea Tacchetti, Maja Trebacz, Kevin Robinson, Yash Katariya, Sebastian Riedel, Paige Bailey,
621 Kefan Xiao, Nimesh Ghelani, Lora Aroyo, Ambrose Slone, Neil Houlsby, Xuehan Xiong, Zhen
622 Yang, Elena Gribovskaya, Jonas Adler, Mateo Wirth, Lisa Lee, Music Li, Thais Kagohara, Jay
623 Pavagadhi, Sophie Bridgers, Anna Bortsova, Sanjay Ghemawat, Zafarali Ahmed, Tianqi Liu,
624 Richard Powell, Vijay Bolina, Mariko Inuma, Polina Zablotskaia, James Besley, Da-Woon Chung,
625 Timothy Dozat, Ramona Comanescu, Xiance Si, Jeremy Greer, Guolong Su, Martin Polacek,
626 Raphaël Lopez Kaufman, Simon Tokumine, Hexiang Hu, Elena Buchatskaya, Yingjie Miao,
627 Mohamed Elhawaty, Aditya Siddhant, Nenad Tomasev, Jinwei Xing, Christina Greer, Helen Miller,
628 Shereen Ashraf, Aurko Roy, Zizhao Zhang, Ada Ma, Angelos Filos, Milos Besta, Rory Blevins,
629 Ted Klimenko, Chih-Kuan Yeh, Soravit Changpinyo, Jiaqi Mu, Oscar Chang, Mantas Pajarskas,
630 Carrie Muir, Vered Cohen, Charline Le Lan, Krishna Haridasan, Amit Marathe, Steven Hansen,
631 Sholto Douglas, Rajkumar Samuel, Mingqiu Wang, Sophia Austin, Chang Lan, Jiepu Jiang, Justin
632 Chiu, Jaime Alonso Lorenzo, Lars Lowe Sjöstrand, Sébastien Cevey, Zach Gleicher, Thi Avrahami,
633 Anudhyan Boral, Hansa Srinivasan, Vittorio Selo, Rhys May, Konstantinos Aisopos, Léonard
634 Hussenot, Livio Baldini Soares, Kate Baumli, Michael B. Chang, Adrià Recasens, Ben Caine,
635 Alexander Pritzel, Filip Pavetic, Fabio Pardo, Anita Gergely, Justin Frye, Vinay Ramasesh, Dan
636 Horgan, Kartikeya Badola, Nora Kassner, Subhrajit Roy, Ethan Dyer, Víctor Campos Campos, Alex
637 Tomala, Yunhao Tang, Dalia El Badawy, Elspeth White, Basil Mustafa, Oran Lang, Abhishek Jindal,
638 Sharad Vikram, Zhitao Gong, Sergi Caelles, Ross Hemsley, Gregory Thornton, Fangxiaoyu Feng,
639 Wojciech Stokowiec, Ce Zheng, Phoebe Thacker, Çağlar Ünlü, Zhishuai Zhang, Mohammad Saleh,
640 James Svensson, Max Bileschi, Piyush Patil, Ankesh Anand, Roman Ring, Katerina Tsihlias, Arpi
641 Vezer, Marco Selvi, Toby Shevlane, Mikel Rodriguez, Tom Kwiatkowski, Samira Daruki, Keran
642 Rong, Allan Dafoe, Nicholas FitzGerald, Keren Gu-Lemberg, Mina Khan, Lisa Anne Hendricks,
643 Marie Pellat, Vladimir Feinberg, James Cobon-Kerr, Tara Sainath, Maribeth Rauh, Sayed Hadi
644 Hashemi, Richard Ives, Yana Hasson, Eric Noland, Yuan Cao, Nathan Byrd, Le Hou, Qingze
645 Wang, Thibault Sottiaux, Michela Paganini, Jean-Baptiste Lespiau, Alexandre Moufarek, Samer
646 Hassan, Kaushik Shivakumar, Joost van Amersfoort, Amol Mandhane, Pratik Joshi, Anirudh Goyal,
647 Matthew Tung, Andrew Brock, Hannah Sheahan, Vedant Misra, Cheng Li, Nemanja Rakićević,
Mostafa Dehghani, Fangyu Liu, Sid Mittal, Junhyuk Oh, Seb Noury, Eren Sezener, Fantine Huot,
Matthew Lamm, Nicola De Cao, Charlie Chen, Sidharth Mudgal, Romina Stella, Kevin Brooks,
Gautam Vasudevan, Chenxi Liu, Mainak Chain, Nivedita Melinkeri, Aaron Cohen, Venus Wang,
Kristie Seymore, Sergey Zubkov, Rahul Goel, Summer Yue, Sai Krishnakumaran, Brian Albert,
Nate Hurley, Motoki Sano, Anhad Mohanane, Jonah Joughin, Egor Filonov, Tomasz Kepa, Yomna

648 Eldawy, Jiawern Lim, Rahul Rishi, Shirin Badieezadegan, Taylor Bos, Jerry Chang, Sanil Jain, Sri
649 Gayatri Sundara Padmanabhan, Subha Puttagunta, Kalpesh Krishna, Leslie Baker, Norbert Kalb,
650 Vamsi Bedapudi, Adam Kurzrok, Shuntong Lei, Anthony Yu, Oren Litvin, Xiang Zhou, Zhichun
651 Wu, Sam Sobell, Andrea Siciliano, Alan Papir, Robby Neale, Jonas Bragagnolo, Tej Toor, Tina
652 Chen, Valentin Anklin, Feiran Wang, Richie Feng, Milad Gholami, Kevin Ling, Lijuan Liu, Jules
653 Walter, Hamid Moghaddam, Arun Kishore, Jakob Adamek, Tyler Mercado, Jonathan Mallinson,
654 Siddhinita Wandekar, Stephen Cagle, Eran Ofek, Guillermo Garrido, Clemens Lombriser, Maksim
655 Mukha, Botu Sun, Hafeezul Rahman Mohammad, Josip Matak, Yadi Qian, Vikas Peswani, Pawel
656 Janus, Quan Yuan, Leif Schelin, Oana David, Ankur Garg, Yifan He, Oleksii Duzhyi, Anton
657 Älgmyr, Timothée Lottaz, Qi Li, Vikas Yadav, Luyao Xu, Alex Chinien, Rakesh Shivanna,
658 Aleksandr Chuklin, Josie Li, Carrie Spadine, Travis Wolfe, Kareem Mohamed, Subhabrata Das,
659 Zihang Dai, Kyle He, Daniel von Dincklage, Shyam Upadhyay, Akanksha Maurya, Luyan Chi,
660 Sebastian Krause, Khalid Salama, Pam G Rabinovitch, Pavan Kumar Reddy M, Aarush Selvan,
661 Mikhail Dektiarev, Golnaz Ghiasi, Erdem Guven, Himanshu Gupta, Boyi Liu, Deepak Sharma,
662 Idan Heimlich Shtacher, Shachi Paul, Oscar Akerlund, François-Xavier Aubet, Terry Huang, Chen
663 Zhu, Eric Zhu, Elico Teixeira, Matthew Fritze, Francesco Bertolini, Liana-Eleonora Marinescu,
664 Martin Bølle, Dominik Paulus, Khyatti Gupta, Tejasi Latkar, Max Chang, Jason Sanders, Roopa
665 Wilson, Xuwei Wu, Yi-Xuan Tan, Lam Nguyen Thiet, Tulsee Doshi, Sid Lall, Swaroop Mishra,
666 Wanming Chen, Thang Luong, Seth Benjamin, Jasmine Lee, Ewa Andrejczuk, Dominik Rabiej,
667 Vipul Ranjan, Krzysztof Styrz, Pengcheng Yin, Jon Simon, Malcolm Rose Harriott, Mudit Bansal,
668 Alexei Robsky, Geoff Bacon, David Greene, Daniil Mirylenka, Chen Zhou, Obaid Sarvana,
669 Abhimanyu Goyal, Samuel Andermatt, Patrick Siegler, Ben Horn, Assaf Israel, Francesco Pongetti,
670 Chih-Wei "Louis" Chen, Marco Selvatici, Pedro Silva, Kathie Wang, Jackson Tolins, Kelvin Guu,
671 Roey Yogev, Xiaochen Cai, Alessandro Agostini, Maulik Shah, Hung Nguyen, Noah Ó Donnaile,
672 Sébastien Pereira, Linda Friso, Adam Stambler, Adam Kurzrok, Chenkai Kuang, Yan Romanikhin,
673 Mark Geller, ZJ Yan, Kane Jang, Cheng-Chun Lee, Wojciech Fica, Eric Malmi, Qijun Tan, Dan
674 Banica, Daniel Balle, Ryan Pham, Yanping Huang, Diana Avram, Hongzhi Shi, Jasjot Singh, Chris
675 Hidey, Niharika Ahuja, Pranab Saxena, Dan Dooley, Srividya Pranavi Potharaju, Eileen O'Neill,
676 Anand Gokulchandran, Ryan Foley, Kai Zhao, Mike Dusenberry, Yuan Liu, Pulkit Mehta, Ragha
677 Kotikalapudi, Chalence Safranek-Shrader, Andrew Goodman, Joshua Kessinger, Eran Globen,
678 Prateek Kolhar, Chris Gorgolewski, Ali Ibrahim, Yang Song, Ali Eichenbaum, Thomas Brovelli,
679 Sahitya Potluri, Preethi Lahoti, Cip Baetu, Ali Ghorbani, Charles Chen, Andy Crawford, Shalini
680 Pal, Mukund Sridhar, Petru Gurita, Asier Mujika, Igor Petrovski, Pierre-Louis Cedoz, Chenmei Li,
681 Shiyuan Chen, Niccolò Dal Santo, Siddharth Goyal, Jitesh Punjabi, Karthik Kappaganthu, Chester
682 Kwak, Pallavi LV, Sarmishta Velury, Himadri Choudhury, Jamie Hall, Premal Shah, Ricardo
683 Figueira, Matt Thomas, Minjie Lu, Ting Zhou, Chintu Kumar, Thomas Jurdi, Sharat Chikkerur,
684 Yenai Ma, Adams Yu, Soo Kwak, Victor Åhdel, Sujeewan Rajayogam, Travis Choma, Fei Liu,
685 Aditya Barua, Colin Ji, Ji Ho Park, Vincent Hellendoorn, Alex Bailey, Taylan Bilal, Huanjie Zhou,
686 Mehrdad Khatir, Charles Sutton, Wojciech Rzadkowski, Fiona Macintosh, Konstantin Shagin, Paul
687 Medina, Chen Liang, Jinjing Zhou, Pararth Shah, Yingying Bi, Attila Dankovics, Shipra Banga,
688 Sabine Lehmann, Marissa Bredesen, Zifan Lin, John Eric Hoffmann, Jonathan Lai, Raynald Chung,
689 Kai Yang, Nihal Balani, Arthur Bražinskas, Andrei Sozanschi, Matthew Hayes, Héctor Fernández
690 Alcalde, Peter Makarov, Will Chen, Antonio Stella, Liselotte Sniijders, Michael Mandl, Ante
691 Kärroman, Paweł Nowak, Xinyi Wu, Alex Dyck, Krishnan Vaidyanathan, Raghavender R, Jessica
692 Mallet, Mitch Rudominer, Eric Johnston, Sushil Mittal, Akhil Udathu, Janara Christensen, Vishal
693 Verma, Zach Irving, Andreas Santucci, Gamaleldin Elsayed, Elnaz Davoodi, Marin Georgiev, Ian
694 Tenney, Nan Hua, Geoffrey Cideron, Edouard Leurent, Mahmoud Alnahlawi, Ionut Georgescu,
695 Nan Wei, Ivy Zheng, Dylan Scandinaro, Heinrich Jiang, Jasper Snoek, Mukund Sundararajan,
696 Xuezhi Wang, Zack Ontiveros, Itay Karo, Jeremy Cole, Vinu Rajashekhar, Lara Tumeh, Eyal Ben-
697 David, Rishub Jain, Jonathan Uesato, Romina Datta, Oskar Bunyan, Shimu Wu, John Zhang, Piotr
698 Stanczyk, Ye Zhang, David Steiner, Subhjit Naskar, Michael Azzam, Matthew Johnson, Adam
699 Paszke, Chung-Cheng Chiu, Jaume Sanchez Elias, Afroz Mohiuddin, Faizan Muhammad, Jin
700 Miao, Andrew Lee, Nino Vieillard, Jane Park, Jiageng Zhang, Jeff Stanway, Drew Garmon, Abhijit
701 Karmarkar, Zhe Dong, Jong Lee, Aviral Kumar, Luowei Zhou, Jonathan Evens, William Isaac,
Geoffrey Irving, Edward Loper, Michael Fink, Isha Arkatkar, Nanxin Chen, Izhak Shafraan, Ivan
Petrychenko, Zhe Chen, Johnson Jia, Anselm Levskaya, Zhenkai Zhu, Peter Grabowski, Yu Mao,
Alberto Magni, Kaisheng Yao, Javier Snaider, Norman Casagrande, Evan Palmer, Paul Suganthan,
Alfonso Castaño, Irene Giannoumis, Wooyeol Kim, Mikołaj Rybiński, Ashwin Sreevatsa, Jennifer
Prendki, David Soergel, Adrian Goedeckemeyer, Willi Gierke, Mohsen Jafari, Meenu Gaba, Jeremy

702 Wiesner, Diana Gage Wright, Yawen Wei, Harsha Vashisht, Yana Kulizhskaya, Jay Hoover, Maigo
703 Le, Lu Li, Chimezie Iwuanyanwu, Lu Liu, Kevin Ramirez, Andrey Khorlin, Albert Cui, Tian
704 LIN, Marcus Wu, Ricardo Aguilar, Keith Pallo, Abhishek Chakladar, Ginger Perng, Elena Allica
705 Abellan, Mingyang Zhang, Ishita Dasgupta, Nate Kushman, Ivo Penchev, Alena Repina, Xihui Wu,
706 Tom van der Weide, Priya Ponnappalli, Caroline Kaplan, Jiri Simsa, Shuangfeng Li, Olivier Dousse,
707 Fan Yang, Jeff Piper, Nathan Ie, Rama Pasumarthi, Nathan Lintz, Anitha Vijayakumar, Daniel
708 Andor, Pedro Valenzuela, Minnie Lui, Cosmin Paduraru, Daiyi Peng, Katherine Lee, Shuyuan
709 Zhang, Somer Greene, Duc Dung Nguyen, Paula Kurylowicz, Cassidy Hardin, Lucas Dixon, Lili
710 Janzer, Kiam Choo, Ziqiang Feng, Biao Zhang, Achintya Singhal, Dayou Du, Dan McKinnon,
711 Natasha Antropova, Tolga Bolukbasi, Orgad Keller, David Reid, Daniel Finchelstein, Maria Abi
712 Raad, Remi Crocker, Peter Hawkins, Robert Dadashi, Colin Gaffney, Ken Franko, Anna Bulanova,
713 Rémi Leblond, Shirley Chung, Harry Askham, Luis C. Cobo, Kelvin Xu, Felix Fischer, Jun Xu,
714 Christina Sorokin, Chris Alberti, Chu-Cheng Lin, Colin Evans, Alek Dimitriev, Hannah Forbes,
715 Dylan Banarse, Zora Tung, Mark Omernick, Colton Bishop, Rachel Sterneck, Rohan Jain, Jiawei
716 Xia, Ehsan Amid, Francesco Piccinno, Xingyu Wang, Praseem Banzal, Daniel J. Mankowitz, Alex
717 Polozov, Victoria Krakovna, Sasha Brown, MohammadHossein Bateni, Dennis Duan, Vlad Firoiu,
718 Meghana Thotakuri, Tom Natan, Matthieu Geist, Ser tan Girgin, Hui Li, Jiayu Ye, Ofir Roval,
719 Reiko Tojo, Michael Kwong, James Lee-Thorp, Christopher Yew, Danila Sinopalnikov, Sabela
720 Ramos, John Mellor, Abhishek Sharma, Kathy Wu, David Miller, Nicolas Sonnerat, Denis Vnukov,
721 Rory Greig, Jennifer Beattie, Emily Caveness, Libin Bai, Julian Eisenschlos, Alex Korchemniy,
722 Tomy Tsai, Mimi Jasarevic, Weize Kong, Phuong Dao, Zeyu Zheng, Frederick Liu, Fan Yang,
723 Rui Zhu, Tian Huey Teh, Jason Sanmiya, Evgeny Gladchenko, Nejc Trdin, Daniel Toyama, Evan
724 Rosen, Sasan Tavakkol, Linting Xue, Chen Elkind, Oliver Woodman, John Carpenter, George
725 Papamakarios, Rupert Kemp, Sushant Kaffle, Tanya Grunina, Rishika Sinha, Alice Talbert, Diane
726 Wu, Denese Owusu-Afriyie, Cosmo Du, Chloe Thornton, Jordi Pont-Tuset, Pradyumna Narayana,
727 Jing Li, Saaber Fatehi, John Wieting, Omar Ajmeri, Benigno Uria, Yeongil Ko, Laura Knight,
728 Amélie Héliou, Ning Niu, Shane Gu, Chenxi Pang, Yeqing Li, Nir Levine, Ariel Stolovich, Rebeca
729 Santamaria-Fernandez, Sonam Goenka, Wenny Yustalim, Robin Strudel, Ali Elqursh, Charlie
730 Deck, Hyo Lee, Zonglin Li, Kyle Levin, Raphael Hoffmann, Dan Holtmann-Rice, Olivier Bachem,
731 Sho Arora, Christy Koh, Soheil Hassas Yeganeh, Siim Pöder, Mukarram Tariq, Yanhua Sun,
732 Lucian Ionita, Mojtaba Seyedhosseini, Pouya Tafti, Zhiyu Liu, Anmol Gulati, Jasmine Liu, Xinyu
733 Ye, Bart Chrzaszcz, Lily Wang, Nikhil Sethi, Tianrun Li, Ben Brown, Shreya Singh, Wei Fan,
734 Aaron Parisi, Joe Stanton, Vinod Koverkathu, Christopher A. Choquette-Choo, Yunjie Li, TJ Lu,
735 Abe Ittycheriah, Prakash Shroff, Mani Varadarajan, Sanaz Bahargam, Rob Willoughby, David
736 Gaddy, Guillaume Desjardins, Marco Cornero, Brona Robenek, Bhavishya Mittal, Ben Albrecht,
737 Ashish Shenoy, Fedor Moiseev, Henrik Jacobsson, Alireza Ghaffarkhah, Morgane Rivièrè, Alanna
738 Walton, Clément Crepy, Alicia Parrish, Zongwei Zhou, Clement Farabet, Carey Radebaugh,
739 Praveen Srinivasan, Claudia van der Salm, Andreas Fildjeland, Salvatore Scellato, Eri Latorre-
740 Chimoto, Hanna Klimczak-Plucińska, David Bridson, Dario de Cesare, Tom Hudson, Piermaria
741 Mendolicchio, Lexi Walker, Alex Morris, Matthew Mauger, Alexey Guseynov, Alison Reid, Seth
742 Odoom, Lucia Loher, Victor Cotruta, Madhavi Yenugula, Dominik Grewe, Anastasia Petrushkina,
743 Tom Duerig, Antonio Sanchez, Steve Yadlowsky, Amy Shen, Amir Globerson, Lynette Webb,
744 Sahil Dua, Dong Li, Surya Bhupatiraju, Dan Hurt, Haroon Qureshi, Ananth Agarwal, Tomer Shani,
745 Matan Eyal, Anuj Khare, Shreyas Rammohan Belle, Lei Wang, Chetan Tekur, Mihir Sanjay Kale,
746 Jinliang Wei, Ruoxin Sang, Brennan Saeta, Tyler Liechty, Yi Sun, Yao Zhao, Stephan Lee, Pandu
747 Nayak, Doug Fritz, Manish Reddy Vuyyuru, John Aslanides, Nidhi Vyas, Martin Wicke, Xiao Ma,
748 Evgenii Eltyshv, Nina Martin, Hardie Cate, James Manyika, Keyvan Amiri, Yelin Kim, Xi Xiong,
749 Kai Kang, Florian Luisier, Nilesh Tripuraneni, David Madras, Mandy Guo, Austin Waters, Oliver
750 Wang, Joshua Ainslie, Jason Baldridge, Han Zhang, Garima Pruthi, Jakob Bauer, Feng Yang, Riham
751 Mansour, Jason Gelman, Yang Xu, George Polovets, Ji Liu, Honglong Cai, Warren Chen, XiangHai
752 Sheng, Emily Xue, Sherjil Ozair, Christof Angermueller, Xiaowei Li, Anoop Sinha, Weiren Wang,
753 Julia Wiesinger, Emmanouil Koukoumidis, Yuan Tian, Anand Iyer, Madhu Gurusurthy, Mark
754 Goldenson, Parashar Shah, MK Blake, Hongkun Yu, Anthony Urbanowicz, Jennimaria Palomaki,
755 Chrisantha Fernando, Ken Durden, Harsh Mehta, Nikola Momchev, Elahe Rahimtoroghi, Maria
Georgaki, Amit Raul, Sebastian Ruder, Morgan Redshaw, Jinhyuk Lee, Denny Zhou, Komal Jalan,
Dinghua Li, Blake Hechtman, Parker Schuh, Milad Nasr, Kieran Milan, Vladimir Mikulik, Juliana
Franco, Tim Green, Nam Nguyen, Joe Kelley, Aroma Mahendru, Andrea Hu, Joshua Howland, Ben
Vargas, Jeffrey Hui, Kshitij Bansal, Vikram Rao, Rakesh Ghiya, Emma Wang, Ke Ye, Jean Michel
Sarr, Melanie Moranski Preston, Madeleine Elish, Steve Li, Aakash Kaku, Jigar Gupta, Ice Pasupat,

- 756 Da-Cheng Juan, Milan Someswar, Tejvi M., Xinyun Chen, Aida Amini, Alex Fabrikant, Eric Chu,
757 Xuanyi Dong, Amruta Muthal, Senaka Buthpitiya, Sarthak Jauhari, Nan Hua, Urvashi Khandelwal,
758 Ayal Hitron, Jie Ren, Larissa Rinaldi, Shahar Drath, Avigail Dabush, Nan-Jiang Jiang, Harshal
759 Godhia, Uli Sachs, Anthony Chen, Yicheng Fan, Hagai Taitelbaum, Hila Noga, Zhuyun Dai, James
760 Wang, Chen Liang, Jenny Hamer, Chun-Sung Ferng, Chenel Elkind, Aviel Atias, Paulina Lee, Vít
761 Listík, Mathias Carlen, Jan van de Kerkhof, Marcin Pikus, Krunoslav Zaher, Paul Müller, Sasha
762 Zykova, Richard Stefanec, Vitaly Gatsko, Christoph Hirnschall, Ashwin Sethi, Xingyu Federico
763 Xu, Chetan Ahuja, Beth Tsai, Anca Stefanoiu, Bo Feng, Keshav Dhandhanian, Manish Katyal,
764 Akshay Gupta, Atharva Parulekar, Divya Pitta, Jing Zhao, Vivaan Bhatia, Yashodha Bhavnani,
765 Omar Alhadlaq, Xiaolin Li, Peter Danenberg, Dennis Tu, Alex Pine, Vera Filippova, Abhipso
766 Ghosh, Ben Limonchik, Bhargava Urala, Chaitanya Krishna Lanka, Derik Clive, Yi Sun, Edward
767 Li, Hao Wu, Kevin Hongtongsak, Ianna Li, Kalind Thakkar, Kuanysh Omarov, Kushal Majmundar,
768 Michael Alverson, Michael Kucharski, Mohak Patel, Mudit Jain, Maksim Zabelin, Paolo Pelagatti,
769 Rohan Kohli, Saurabh Kumar, Joseph Kim, Swetha Sankar, Vineet Shah, Lakshmi Ramachandruni,
770 Xiangkai Zeng, Ben Bariach, Laura Weidinger, Tu Vu, Amar Subramanya, Sissie Hsiao, Demis
771 Hassabis, Koray Kavukcuoglu, Adam Sadovsky, Quoc Le, Trevor Strohman, Yonghui Wu, Slav
772 Petrov, Jeffrey Dean, and Oriol Vinyals. Gemini: A family of highly capable multimodal models.
773 Technical report, Google, 2024. URL <https://arxiv.org/pdf/2312.11805>.
- 774 Lin Gui, Cristina Gârbacea, and Victor Veitch. BoNBoN alignment for large language models and
775 the sweetness of best-of-n sampling. *Computing Research Repository*, arXiv:2406.00832, 2024.
776 URL <https://arxiv.org/pdf/2406.00832>.
- 777 Alexander Havrilla, Maksym Zhuravinskiy, Duy Phung, Aman Tiwari, Jonathan Tow, Stella Biderman,
778 Quentin Anthony, and Louis Castricato. trlX: A framework for large scale reinforcement learning
779 from human feedback. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of
780 the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 8578–8595,
781 Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.
782 emnlp-main.530. URL <https://aclanthology.org/2023.emnlp-main.530>.
- 783 Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger.
784 Deep reinforcement learning that matters. In *Proceedings of the Thirty-Second AAAI Con-
785 ference on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelli-
786 gence Conference and Eighth AAAI Symposium on Educational Advances in Artificial Intel-
787 ligence*, AAAI’18/IAAI’18/EAAI’18. AAAI Press, 2018. ISBN 978-1-57735-800-8. URL
788 <https://dl.acm.org/doi/pdf/10.5555/3504035.3504427>.
- 789 Joseph L. Hodges. The significance probability of the smirnov two-sample test. *Arkiv för Matematik*,
790 3:469–486, 1958. URL <https://api.semanticscholar.org/CorpusID:121451525>.
- 791 Shengyi Huang, Michael Noukhovitch, Arian Hosseini, Kashif Rasul, Weixun Wang, and Lewis Tun-
792 stall. The N+ implementation details of RLHF with PPO: A case study on TL;DR summarization.
793 In *First Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=kH0ZTa8e3>.
- 794
- 795 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.
796 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
797 serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating
798 Systems Principles*, 2023.
- 799
- 800 Sergey Levine. Reinforcement learning and control as probabilistic inference: Tutorial and review.
801 *Computing Research Repository*, arXiv:1805.00909, 2018. URL <https://arxiv.org/pdf/1805.00909>.
- 802
- 803 Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher
804 Potts. Learning word vectors for sentiment analysis. In Dekang Lin, Yuji Matsumoto, and Rada
805 Mihalcea (eds.), *Proceedings of the 49th Annual Meeting of the Association for Computational
806 Linguistics: Human Language Technologies*, pp. 142–150, Portland, Oregon, USA, June 2011.
807 Association for Computational Linguistics. URL <https://aclanthology.org/P11-1015>.
- 808
- 809 Youssef Mroueh. Information theoretic guarantees for policy alignment in large language models.
809 *Computing Research Repository*, arXiv:2406.05883, 2024. URL <https://arxiv.org/abs/2406.05883>.

- 810 Sidharth Mudgal, Jong Lee, Harish Ganapathy, YaGuang Li, Tao Wang, Yanping Huang, Zhifeng
811 Chen, Heng-Tze Cheng, Michael Collins, Trevor Strohman, Jilin Chen, Alex Beutel, and Ahmad
812 Beirami. Controlled decoding from language models. In *Proceedings of The 41st International
813 Conference on Machine Learning*, Proceedings of Machine Learning Research. PMLR, 2024. URL
814 <https://arxiv.org/pdf/2310.17022>.
- 815
- 816 Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher
817 Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, Xu Jiang, Karl Cobbe, Tyna Eloundou,
818 Gretchen Krueger, Kevin Button, Matthew Knight, Benjamin Chess, and John Schulman. WebGPT:
819 Browser-assisted question-answering with human feedback. *Computing Research Repository*,
820 arXiv:2112.09332, 2022. URL <https://arxiv.org/pdf/2112.09332>.
- 821
- 822 OpenAI, :, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Floren-
823 cia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red
824 Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavar-
825 ian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner,
826 Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim
827 Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey,
828 Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully
829 Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won
830 Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah
831 Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien
832 Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman,
833 Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni,
834 Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene,
835 Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He,
836 Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele,
837 Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain,
838 Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn,
839 Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish
840 Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Hendrik
841 Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich,
842 Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy
843 Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie
844 Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini,
845 Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne,
846 Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David
847 Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie
848 Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély,
849 Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo
850 Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano,
851 Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng,
852 Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto,
853 Michael Pokorny, Michelle Pocrass, Vitchyr Pong, Tolly Powell, Alethea Power, Boris Power,
854 Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis
855 Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted
856 Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel
857 Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon
858 Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky,
859 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang,
860 Nikolas Tezak, Madeleine Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston
861 Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,
862 Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason
863 Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff,
Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu,
Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba,
Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang,
William Zhuk, and Barret Zoph. GPT-4 technical report. Technical report, OpenAI, 2023. URL
<https://cdn.openai.com/papers/gpt-4.pdf>.

- 864 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
865 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser
866 Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan
867 Leike, and Ryan Lowe. Training language models to follow instructions with human feed-
868 back. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Ad-
869 vances in Neural Information Processing Systems*, volume 35, pp. 27730–27744. Curran Asso-
870 ciates, Inc., 2022. URL [https://proceedings.neurips.cc/paper_files/paper/2022/file/
871 b1efde53be364a73914f58805a001731-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/b1efde53be364a73914f58805a001731-Paper-Conference.pdf).
- 872 Alizée Pace, Jonathan Mallinson, Eric Malmi, Sebastian Krause, and Aliaksei Severyn. West-of-n:
873 Synthetic preference generation for improved reward modeling. *Computing Research Repository*,
874 arXiv:2401.12086, 2024. URL <https://arxiv.org/abs/2401.12086>.
- 875 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language mod-
876 els are unsupervised multitask learners, 2019. URL [https://d4mucfpksywv.cloudfront.net/
877 better-language-models/language_models_are_unsupervised_multitask_learners.pdf](https://d4mucfpksywv.cloudfront.net/better-language-models/language_models_are_unsupervised_multitask_learners.pdf).
- 878 Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea
879 Finn. Direct preference optimization: Your language model is secretly a reward model. In
880 *Advances in Neural Information Processing Systems*, volume 36. Curran Associates, Inc., 2023.
881 URL <https://arxiv.org/pdf/2305.18290.pdf>.
- 882 Tom Schaul, Georg Ostrovski, Iurii Kemaev, and Diana Borsa. Return-based scaling: Yet another
883 normalisation trick for deep rl. *Computing Research Repository*, arXiv:2105.05347, 2021. URL
884 <https://arxiv.org/abs/2105.05347>.
- 885 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
886 optimization algorithms. *Computing Research Repository*, arXiv:1707.06347, 2017. URL <https://arxiv.org/abs/1707.06347>.
- 887 Pier Giuseppe Sessa, Robert Dadashi, Léonard Hussenot, Johan Ferret, Nino Vieillard, Alexandre
888 Ramé, Bobak Shariari, Sarah Perrin, Abe Friesen, Geoffrey Cideron, Sertan Girgin, Piotr Stanczyk,
889 Andrea Michi, Danila Sinopalnikov, Sabela Ramos, Amélie Héliou, Aliaksei Severyn, Matt
890 Hoffman, Nikola Momchev, and Olivier Bachem. Bond: Aligning llms with best-of-n distillation.
891 *Computing Research Repository*, arXiv:2401.12086, 2024. URL [https://arxiv.org/abs/2401.
892 12086](https://arxiv.org/abs/2401.12086).
- 893 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute opti-
894 mally can be more effective than scaling model parameters. *Computing Research Repository*,
895 arXiv:2408.03314, 2024. URL <https://arxiv.org/abs/2408.03314>.
- 896 Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec
897 Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feed-
898 back. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Ad-
899 vances in Neural Information Processing Systems*, volume 33, pp. 3008–3021. Curran Asso-
900 ciates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/paper/2020/file/
901 1f89885d556929e98d3ef9b86448f951-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1f89885d556929e98d3ef9b86448f951-Paper.pdf).
- 902 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
903 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cris-
904 tian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu,
905 Wenying Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn,
906 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel
907 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya
908 Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar
909 Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan
910 Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan,
911 Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov,
912 Yuchen Zhang, Angela Fan, Melanیه Kambadur, Sharan Narang, Aurelien Rodriguez, Robert
913 Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat
914 models. Technical report, Meta, 2023. URL [https://ai.meta.com/research/publications/
915 llama-2-open-foundation-and-fine-tuned-chat-models/](https://ai.meta.com/research/publications/llama-2-open-foundation-and-fine-tuned-chat-models/).
- 916
- 917

918 Joy Qiping Yang, Salman Salamatian, Ziteng Sun, Ananda Theertha Suresh, and Ahmad Beirami.
919 Asymptotics of language model alignment. *Computing Research Repository*, arXiv:2404.01730,
920 2024. URL <https://arxiv.org/pdf/2404.01730>.
921
922 Kevin Yang and Dan Klein. FUDGE: Controlled text generation with future discriminators. In
923 *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computa-*
924 *tional Linguistics: Human Language Technologies*, pp. 3511–3535, Online, June 2021. Association
925 for Computational Linguistics. URL <https://aclanthology.org/2021.naacl-main.276>.
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

Symbol	Type	Explanation
Σ	alphabet	Σ is a set of symbols
\mathbf{y}, \mathbf{y}'	$\in \Sigma^*$	strings in Σ^*
\mathbf{x}	$\in \Sigma^*$	prompt string in Σ^*
$\boldsymbol{\theta}$	$\in \Theta$	A real vector representing the parameters of a language model
$\pi_{\boldsymbol{\theta}}$	language model	A language model parameterized by $\boldsymbol{\theta}$
π_{ref}	language model	A supervised-fine-tuned language model
r	$\Sigma^* \rightarrow \mathbb{R}$	A reward model
β	\mathbb{R}	Regularization parameter for the KL divergence term
F	$\mathbb{R} \rightarrow \mathbb{R}$	A strict cumulative density function of reward values under π_{ref}
N	\mathbb{Z}^+	Number of samples used in BoN algorithm
M	\mathbb{Z}^+	Number of samples used in the MC estimator

Table 3: A summary of the notation used in the paper

A RELATED WORK

Best-of- N . BoN is a straightforward alignment-via-inference algorithm to optimize the output of the language model using a trained reward model (Charniak & Johnson, 2005; Stiennon et al., 2020). Despite its simplicity, BoN performs comparably or even better than other alignment methods, such as RLHF and direct preference optimization (Nakano et al., 2022; Gao et al., 2023; Rafailov et al., 2023). However, as noted by Stiennon et al. (2020), BoN is an inefficient algorithm due to the reduced throughput at inference time.

Applications. BoN has been applied successfully at various stages of the development of language models. Touvron et al. (2023); Dong et al. (2023) employ iterative supervised fine-tuning on the outputs of the BoN algorithm to clone its behavior in the model. Pace et al. (2024) leverage BoN to enhance reward modeling by training the reward model on both the best and worst responses. Additionally, Brown et al. (2024); Snell et al. (2024) explore the scaling laws for alignment-via-inference methods and demonstrate how to utilize the limited inference budget to achieve the alignment.

Best-of- N as an alignment-via-fine-tuning method. Two concurrent efforts to ours have also attempted to convert BoN to an alignment-via-fine-tuning method. First, Gui et al. (2024) approximate the BoN by maximizing the likelihood of the Best-of- N response and adjusting the relative likelihood of the Best-of- N and the Worst-of- N response. Second, Sessa et al. (2024) similar to ours uses reinforcement learning to minimize the distance between the language model and the BoN policy. Different from ours, and to reduce the fine-tuning time, the authors use a crude estimation of $\log F$ and approximate the distance to Best-of- N by iteratively distilling the Best-of-2 model as a moving anchor.

B PROOF OF PROP. 1

Proposition 1. *Suppose $r: \Sigma^* \rightarrow \mathbb{R}$ is a one-to-one mapping. Then, the probability that a string $\mathbf{y} \sim \pi_{\text{bon}}$ is given by*

$$\pi_{\text{bon}}(\mathbf{y}) = \sum_{i=1}^N \binom{N}{i} F(r(\mathbf{y}))^{N-i} \pi_{\text{ref}}(\mathbf{y})^i, \quad F(r(\mathbf{y})) \stackrel{\text{def}}{=} \mathbb{P}_{\mathbf{y}' \sim \pi_{\text{ref}}} (r(\mathbf{y}') < r(\mathbf{y})). \quad (4)$$

Proof. The proof follows Casella & Berger (2001, Theorem 5.4.3). To compute $\pi_{\text{bon}}(\mathbf{y})$, we first define two events: (i) the event that all N samples have rewards less than or equal to $r(\mathbf{y})$, and (ii) the

event that all N samples have rewards less than $r(\mathbf{y})$. The probability of those events is as follows:¹⁹

$$p_1(\mathbf{y}) \stackrel{\text{def}}{=} \mathbb{P}(\text{all } N \text{ samples have rewards } \leq r(\mathbf{y})) = \left(F(r(\mathbf{y})) + \pi_{\text{ref}}(\mathbf{y}) \right)^N \quad (10a)$$

$$p_2(\mathbf{y}) \stackrel{\text{def}}{=} \mathbb{P}(\text{all } N \text{ samples have rewards } < r(\mathbf{y})) = F(r(\mathbf{y}))^N. \quad (10b)$$

Note that for Eq. (13a) to hold, we need the assumption that the reward function is a one-to-one mapping.²⁰ Furthermore, given this assumption, $\pi_{\text{bon}}(\mathbf{y})$ is the probability that *at least* one of the sampled strings out of N samples have the reward exactly equal to $r(\mathbf{y})$ and the rest of the samples have rewards less than or equal to $r(\mathbf{y})$. Given how we defined p_1 and p_2 , we have $\pi_{\text{bon}}(\mathbf{y}) = p_1(\mathbf{y}) - p_2(\mathbf{y})$.

$$\pi_{\text{bon}}(\mathbf{y}) = \left(F(r(\mathbf{y})) + \pi_{\text{ref}}(\mathbf{y}) \right)^N - F(r(\mathbf{y}))^N = \sum_{i=1}^N \binom{N}{i} F(r(\mathbf{y}))^{N-i} \pi_{\text{ref}}(\mathbf{y})^i. \quad (11)$$

■

C STRATEGIES FOR NON-INJECTIVE REWARD FUNCTIONS

If the reward function is not injective, we need a tie-breaking strategy for the BoN algorithm. We formalize this as defining a total order \prec_r on Σ^* as follows: for any two strings \mathbf{y}_1 and \mathbf{y}_2 , if $r(\mathbf{y}_1) < r(\mathbf{y}_2)$ then we have $\mathbf{y}_1 \prec_r \mathbf{y}_2$. If $r(\mathbf{y}_1) = r(\mathbf{y}_2)$ then $\mathbf{y}_1 \prec_r \mathbf{y}_2$ only if $\mathbf{y}_1 \prec \mathbf{y}_2$, where \prec is some arbitrary but fixed total order, e.g., lexicographic order. Therefore, we define $F(\mathbf{y})$ as

$$F(\mathbf{y}) \stackrel{\text{def}}{=} \mathbb{P}(\mathbf{y}' \prec_r \mathbf{y}). \quad (12)$$

We then need to define the two events and their probabilities, p_1 and p_2 , given this total order on strings, as follows:

$$p_1(\mathbf{y}) \stackrel{\text{def}}{=} \mathbb{P}(\text{all } N \text{ samples are } \preceq_r \mathbf{y}) = \left(F(\mathbf{y}) + \pi_{\text{ref}}(\mathbf{y}) \right)^N \quad (13a)$$

$$p_2(\mathbf{y}) \stackrel{\text{def}}{=} \mathbb{P}(\text{all } N \text{ samples are } \prec_r \mathbf{y}) = F(\mathbf{y})^N \quad (13b)$$

The rest of the proof is the same as with the one-to-one reward functions.

D PROOF OF THM. 2

Theorem 2. *We have $\mathcal{J}^{\text{BoN}}(\boldsymbol{\theta}) \geq L(\boldsymbol{\theta})$, where*

$$L(\boldsymbol{\theta}) \stackrel{\text{def}}{=} (N-1) \mathbb{E}_{\mathbf{y} \sim \pi_{\boldsymbol{\theta}}} \left[\log F(r(\mathbf{y})) \right] - D_{\text{KL}}(\pi_{\boldsymbol{\theta}} \parallel \pi_{\text{ref}}). \quad (8)$$

¹⁹The PMF of BoN is also derived by [Beirami et al. \(Lemma 1; 2024\)](#). In their notation, $p_1 = \mathcal{F}$ and $p_2 = \mathcal{F}^{-1}$.

²⁰If the reward function is not a one-to-one mapping, we need to devise a tie-breaking strategy. See App. C for further discussion.

1080 *Proof.* First, we prove $\mathcal{J}^{\text{vBoN}}(\boldsymbol{\theta}) \geq L(\boldsymbol{\theta})$.
 1081
 1082
 1083

$$1084 \quad D_{\text{KL}}(\pi_{\boldsymbol{\theta}} \parallel \pi_{\text{bon}}) = \mathbb{E}_{\mathbf{y} \sim \pi_{\boldsymbol{\theta}}} \left[\log \pi_{\boldsymbol{\theta}}(\mathbf{y}) - \log \pi_{\text{bon}}(\mathbf{y}) \right] \quad (14a)$$

$$1085 \quad = \mathbb{E}_{\mathbf{y} \sim \pi_{\boldsymbol{\theta}}} \left[\log \pi_{\boldsymbol{\theta}}(\mathbf{y}) - \log \sum_{i=1}^N \binom{N}{i} F(r(\mathbf{y}))^{N-i} \pi_{\text{ref}}(\mathbf{y})^i \right] \quad (14b)$$

$$1086 \quad \leq \mathbb{E}_{\mathbf{y} \sim \pi_{\boldsymbol{\theta}}} \left[\log \pi_{\boldsymbol{\theta}}(\mathbf{y}) - \log \sum_{i=1}^{N-1} \binom{N}{i} F(r(\mathbf{y}))^{N-i} \pi_{\text{ref}}(\mathbf{y})^i \right] \quad (14c)$$

$$1087 \quad \leq \mathbb{E}_{\mathbf{y} \sim \pi_{\boldsymbol{\theta}}} \left[\log \pi_{\boldsymbol{\theta}}(\mathbf{y}) - \log N F(r(\mathbf{y}))^{N-1} \pi_{\text{ref}}(\mathbf{y})^1 \right] \quad (14d)$$

$$1088 \quad \leq \mathbb{E}_{\mathbf{y} \sim \pi_{\boldsymbol{\theta}}} \left[\log \pi_{\boldsymbol{\theta}}(\mathbf{y}) - \log F(r(\mathbf{y}))^{N-1} \pi_{\text{ref}}(\mathbf{y}) \right] \quad (14e)$$

$$1089 \quad = \mathbb{E}_{\mathbf{y} \sim \pi_{\boldsymbol{\theta}}} \left[\log \pi_{\boldsymbol{\theta}}(\mathbf{y}) - \log \pi_{\text{ref}}(\mathbf{y}) - (N-1) \log F(r(\mathbf{y})) \right] \quad (14f)$$

$$1090 \quad = D_{\text{KL}}(\pi_{\boldsymbol{\theta}} \parallel \pi_{\text{ref}}) - (N-1) \mathbb{E}_{\mathbf{y} \sim \pi_{\boldsymbol{\theta}}} \left[\log F(r(\mathbf{y})) \right] \stackrel{\text{def}}{=} -L(\boldsymbol{\theta}). \quad (14g)$$

1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

The inequality in Eq. (14c) stems from the fact that we drop positive terms in the summation and only keep the first term. Therefore, the lower bound for our objective is:

$$\mathcal{J}^{\text{vBoN}}(\boldsymbol{\theta}) = -D_{\text{KL}}(\pi_{\boldsymbol{\theta}} \parallel \pi_{\text{bon}}) \geq (N-1) \mathbb{E}_{\mathbf{y} \sim \pi_{\boldsymbol{\theta}}} \left[\log F(r(\mathbf{y})) \right] - D_{\text{KL}}(\pi_{\boldsymbol{\theta}} \parallel \pi_{\text{ref}}). \quad (15)$$

■

Another approach to deriving a lower bound is by using the Jensen's inequality. By doing so, we arrive at the following theorem.

Theorem 3. Let $\alpha = \frac{(N+2)(N-1)}{2}$, $\beta = \frac{N(N+1)}{2}$, and $\gamma = \frac{N(N-1)}{2}$. Then, we have $\mathcal{J}^{\text{vBoN}}(\boldsymbol{\theta}) \geq L_1(\boldsymbol{\theta})$, where we further define

$$L_1(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \gamma \mathbb{E}_{\mathbf{y} \sim \pi_{\boldsymbol{\theta}}} \left[\log F(r(\mathbf{y})) \right] - \alpha H(\pi_{\boldsymbol{\theta}}) - \beta D_{\text{KL}}(\pi_{\boldsymbol{\theta}} \parallel \pi_{\text{ref}}). \quad (16)$$

1134 *Proof.*

$$1135 \quad D_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{bon}}) = \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \left[\log \pi_{\theta}(\mathbf{y}) - \log \pi_{\text{bon}}(\mathbf{y}) \right] \quad (17a)$$

$$1136 \quad = \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \left[\log \pi_{\theta}(\mathbf{y}) - \log \sum_{i=1}^N \binom{N}{i} F(r(\mathbf{y}))^{N-i} \pi_{\text{ref}}(\mathbf{y})^i \right] \quad (17b)$$

$$1137 \quad \leq \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \left[\log \pi_{\theta}(\mathbf{y}) - \sum_{i=1}^N \log \binom{N}{i} F(r(\mathbf{y}))^{N-i} \pi_{\text{ref}}(\mathbf{y})^i \right] \quad (17c)$$

$$1138 \quad = \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \left[\log \pi_{\theta}(\mathbf{y}) - \sum_{i=1}^N \log \binom{N}{i} - \sum_{i=1}^N \log F(r(\mathbf{y}))^{N-i} - \sum_{i=1}^N \log \pi_{\text{ref}}(\mathbf{y})^i \right] \quad (17d)$$

$$1139 \quad = \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \left[\log \pi_{\theta}(\mathbf{y}) - \sum_{i=1}^N \log \binom{N}{i} - \log F(r(\mathbf{y})) \sum_{i=1}^N (N-i) - \log \pi_{\text{ref}}(\mathbf{y}) \sum_{i=1}^N i \right] \quad (17e)$$

$$1140 \quad \leq \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \left[\log \pi_{\theta}(\mathbf{y}) - \frac{N(N-1)}{2} \log F(r(\mathbf{y})) - \frac{N(N+1)}{2} \log \pi_{\text{ref}}(\mathbf{y}) \right] \quad (17f)$$

$$1141 \quad = \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \left[\log \pi_{\theta}(\mathbf{y}) - \frac{N(N+1)}{2} \log \pi_{\text{ref}}(\mathbf{y}) - \frac{N(N-1)}{2} \log F(r(\mathbf{y})) \right] \quad (17g)$$

$$1142 \quad = \frac{N(N+1)}{2} D_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{ref}}) + \mathbb{E}_{\pi_{\theta}} \left[\frac{-(N+2)(N-1)}{2} \log \pi_{\theta}(\mathbf{y}) - \frac{N(N-1)}{2} \log F(r(\mathbf{y})) \right] \quad (17h)$$

$$1143 \quad = \frac{N(N+1)}{2} D_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{ref}}) + \frac{(N+2)(N-1)}{2} \mathbb{H}(\pi_{\theta}) - \mathbb{E}_{\pi_{\theta}} \left[\frac{N(N-1)}{2} \log F(r(\mathbf{y})) \right] \quad (17i)$$

1144 In Eq. (17c), because $-\log(x)$ is convex for $x \geq 0$, we applied Jensen's inequality to obtain the
1145 upper bound. Abstracting away from the three multiplicative factors, naming them γ , α and β , we
1146 end up with the following function

$$1147 \quad \mathcal{J}^{\text{vBon}}(\theta) = -D_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{bon}}) \geq \gamma \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \log F(r(\mathbf{y})) - \alpha \mathbb{H}(\pi_{\theta}) - \beta D_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{ref}}), \quad (18)$$

1148 which is a bound for some settings of γ , α and β . ■

1149 Importantly L_1 is a looser bound compared to L . We formalize this in the following theorem.

1150 **Theorem 4.** For every $\theta \in \Theta$, we have $L(\theta) \geq L_1(\theta)$.

1151 *Proof.* We prove $-L_1(\theta) \geq -L(\theta)$, meaning that L is a tighter lower bound. According to Eq. (17f),
1152 we have:

$$1153 \quad -L_1(\theta) \geq \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \left[\log \pi_{\theta}(\mathbf{y}) - \sum_{i=1}^N \log F(r(\mathbf{y}))^{N-i} \pi_{\text{ref}}(\mathbf{y})^i \right] \quad (19a)$$

$$1154 \quad \geq \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \left[\log \pi_{\theta}(\mathbf{y}) - \sum_{i=1}^{N-1} \log F(r(\mathbf{y}))^{N-i} \pi_{\text{ref}}(\mathbf{y})^i \right] \quad (19b)$$

$$1155 \quad = \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \left[\log \pi_{\theta}(\mathbf{y}) - \log F(r(\mathbf{y}))^{N-1} \pi_{\text{ref}}(\mathbf{y}) \right] = -L(\theta). \quad (19c)$$

Hyperparameter	Value
Episodes	10000
Optimizer	AdamW ($\epsilon = 1e - 5, lr = 3e - 6$)
Scheduler	Linear
Batch Size	32
β (Both for vBoN and KL-constrained RL objective)	0.05
γ (Discount Factor)	1
λ (for GAE)	0.95
Number of PPO Update Iteration Per Epoch	4
PPO’s Policy Clipping Coefficient	0.2
Value Clipping Coefficient	0.2
Value Function Coefficient	0.2
Value Function Loss Clipping	True
Sampling Temperature	0.7

E vBoN PSEUDOCODE

Algorithm 1 The vBoN algorithm

```

1: procedure vBoN( $\pi_{\text{ref}}, r, N, E, B$ )  $\triangleright \mathcal{D}$ : the prompt dataset,  $E$ : number of epochs,  $B$  batch size
2:   Initialize  $\pi_{\theta}$  with  $\pi_{\text{ref}}$ 
3:   for  $E$  epochs :
4:     for each batch in  $\mathcal{D}$  :
5:        $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(B)} \sim \pi_{\theta}(\cdot)$   $\triangleright$  Sample 1 response for each prompt in the batch
6:       Compute  $r(\mathbf{y}^{(1)}), \dots, r(\mathbf{y}^{(B)})$ 
7:       Compute  $F(r(\mathbf{y}^{(1)})), \dots, F(r(\mathbf{y}^{(B)}))$ 
8:       Optimize  $\pi_{\theta}$  with Eq. (5) (or Eq. (8)) using PPO
9:   return  $\pi_{\theta}$ 

```

F EXPERIMENTAL DETAILS

Hyperparameter sweep in the sentiment experiment. To visualize the trade-off between the expected rewards and KL divergence, we vary the degree of the visualization using the following hyperparameters for each method:

- **BoN-SFT**: $N \in [10, 50, 90, 130, 170, 210, 250, 290, 330, 370, 410, 450, 490, 530, 570, 600]$ with 2 different seeds, resulting in 32 runs.
- **PPO**: $\beta \in [0.005, 0.01, 0.02, 0.03, 0.04, 0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 1., 2., 3., 4., 5.]$ with 2 different seeds, resulting in 32 runs.
- **DPO**: $\beta \in [0.01, 0.1, 0.2, 0.3, 0.4, 0.5, 1., 2., 3., 4., 5.]$ with 3 different seeds, resulting in 33 runs.
- **BoNBoN** and **vBoN**: $N \in [1, 2, 3, 4, 8, 16, 32, 64, 128, 256, 512]$ with 3 different seeds, resulting in 33 runs.
- **vBoN** with L bound: $\beta \in [0.005, 0.01, 0.02, 0.03, 0.04, 0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 1., 2., 3., 4., 5.]$ with 2 different seeds, resulting in 32 runs. Note that comparing Eq. (5) and Eq. (1), we have $N = \frac{1}{\beta} + 1$.

PPO Hyperparameters. In App. F, we include the hyperparameters used with the PPO algorithm for the summarization experiment.

G COMPARING THE vBoN OBJECTIVE AND L LOWER BOUND

We compare the performance of models fine-tuned with the vBoN objective and its lower bound (L) in Fig. 6. We observe that the performance of the models is very close to each other.

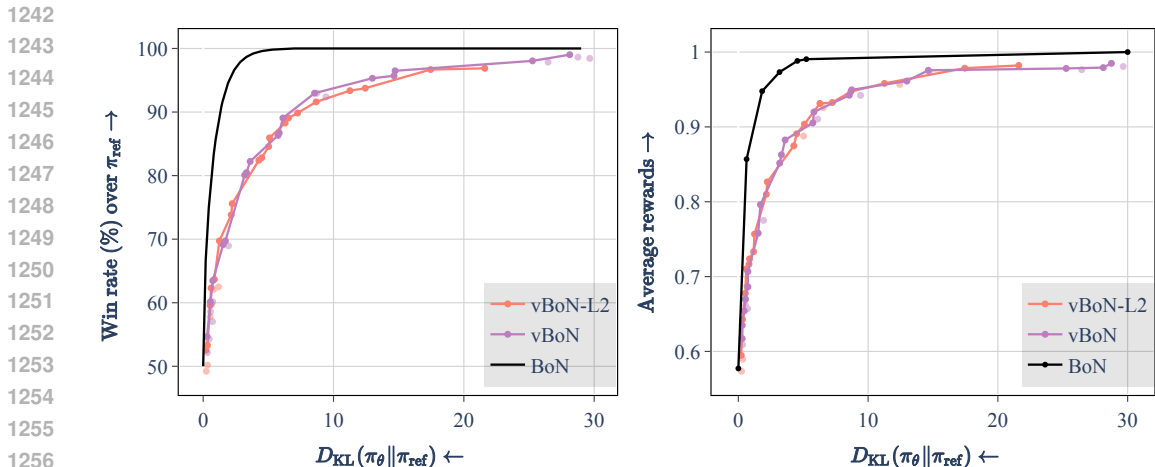
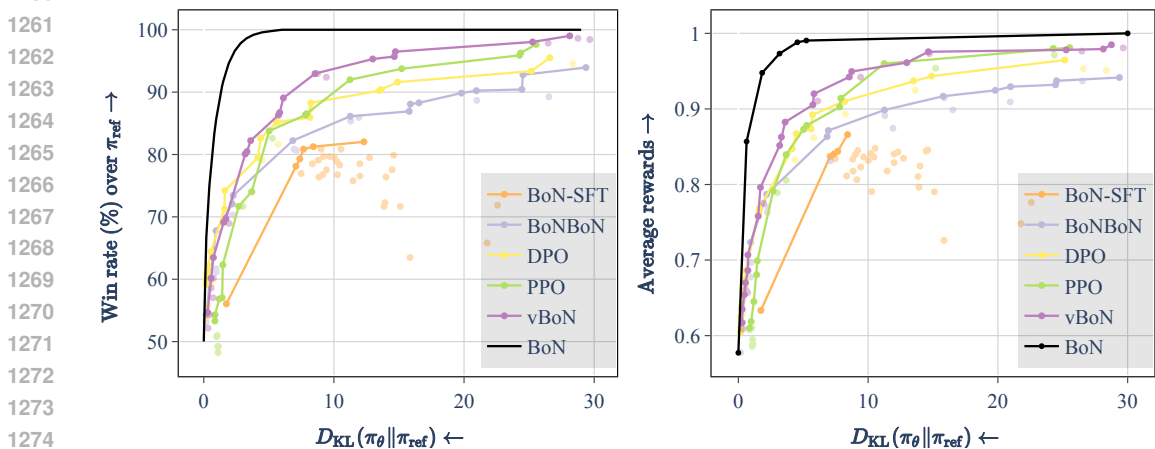


Figure 6: Comparing models trained with the vBoN objective and its lower bound (L). We observe that the performance of the two methods is very close to each other.



(a) 4% of points on Pareto front belong to BoNBoN, (b) 7% of points on Pareto front belong to BoNBoN, 4% to PPO, 42% to DPO, and 50% to vBoN. 10% DPO, 33% PPO, and 50% vBoN.

Figure 7: Steering generated movie reviews towards positive sentiment. Points that are not on the Pareto front have lower opacity.

H ADDITIONAL EXPERIMENTS WITH BoN-SFT

We further experiment with training with maximum likelihood objective on BoN generations when varying N . The results are depicted in Fig. 7. We observe that BoN diverges too much from the reference model compared to other fine-tuning methods for alignment.

I QUALITATIVE RESULTS

1296

1297

1298 Table 4: An example of summaries sampled at temperature 0.25 and their corresponding reward obtained from
1299 the evaluator reward model.

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

Content	Reward
SUBREDDIT: r/relationship_advice TITLE: Stuck in a rut and in need of advice/inspiration! POST: My boyfriend and I have been together for 3 years, and living together for 2. I'm quite the homebody, and when we first met, he was very outgoing and loved partying and socialising (although he was a student at the time). We're both working now, and most nights we find ourselves doing the same things: watching series (luckily we enjoy the same shows), playing Minecraft or playing various board games. We're tired after work, and can't bring ourselves to leave the house. The weekend is much the same – lots of sleep, or sitting around staring at one screen or another. We do party occasionally (we'll head to a pub once every few months) and there are a few mutual friends we enjoy spending time with, but I worry that we've become stuck in our boring ways. I really enjoy our lifestyle, and would be quite happy to never leave the house again, but I'm starting to feel guilty for turning him into a 50 year-old when he's only 24. Any ideas for shaking things up a little? Bear in mind that we live in a small town in South Africa, and neither of us has a car.	-
SFT: I'm stuck in a rut with my boyfriend, and I'd like to shake things up a little. Any ideas?	5.83
PPO: In need of inspiration for moving past boring routines and becoming active and open minded individuals. Any ideas welcome!	4.88
vBoN: Been dating for three years and living together for two. Feeling stuck in rut after working hard all week. Looking for ideas to shake things up!	7.39
BoN: Been together for 3 years, and living together for 2. We're both working now, and find ourselves doing the same things: watching series, playing games, or having sex. We're tired after work, and can't bring ourselves to leave	7.14

1326 Table 5: An example of summaries sampled at temperature 0.5 and their corresponding reward obtained from
1327 the evaluator reward model.

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

Content	Reward
SUBREDDIT: r/relationship_advice TITLE: Stuck in a rut and in need of advice/inspiration! POST: My boyfriend and I have been together for 3 years, and living together for 2. I'm quite the homebody, and when we first met, he was very outgoing and loved partying and socialising (although he was a student at the time). We're both working now, and most nights we find ourselves doing the same things: watching series (luckily we enjoy the same shows), playing Minecraft or playing various board games. We're tired after work, and can't bring ourselves to leave the house. The weekend is much the same – lots of sleep, or sitting around staring at one screen or another. We do party occasionally (we'll head to a pub once every few months) and there are a few mutual friends we enjoy spending time with, but I worry that we've become stuck in our boring ways. I really enjoy our lifestyle, and would be quite happy to never leave the house again, but I'm starting to feel guilty for turning him into a 50 year-old when he's only 24. Any ideas for shaking things up a little? Bear in mind that we live in a small town in South Africa, and neither of us has a car.	-
SFT: I'm stuck in a rut, and need to shake things up to see if it'll work out. Any advice?	3.08
PPO: In need of inspiration to break out of rut and live life fully! Any ideas welcome!	4.59
vBoN: Been happily living together for 2yr+, feeling bored after work regularly, looking for ideas to spice things up!	6.79
BoN: My boyfriend and I have been together for 3 years, and are both working full time. We spend most of our time in the house, and have become boring. What can we do to shake things up?	9.18

1350

1351

1352 Table 6: An example of summaries sampled at temperature 0.75 and their corresponding reward obtained from
1353 the evaluator reward model.

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

Content	Reward
SUBREDDIT: r/relationship_advice TITLE: Stuck in a rut and in need of advice/inspiration! POST: My boyfriend and I have been together for 3 years, and living together for 2. I'm quite the homebody, and when we first met, he was very outgoing and loved partying and socialising (although he was a student at the time). We're both working now, and most nights we find ourselves doing the same things: watching series (luckily we enjoy the same shows), playing Minecraft or playing various board games. We're tired after work, and can't bring ourselves to leave the house. The weekend is much the same – lots of sleep, or sitting around staring at one screen or another. We do party occasionally (we'll head to a pub once every few months) and there are a few mutual friends we enjoy spending time with, but I worry that we've become stuck in our boring ways. I really enjoy our lifestyle, and would be quite happy to never leave the house again, but I'm starting to feel guilty for turning him into a 50 year-old when he's only 24. Any ideas for shaking things up a little? Bear in mind that we live in a small town in South Africa, and neither of us has a car.	-
SFT: Boyfriend and I are stuck in a rut. We don't do anything fun, and I'm starting to feel guilty for not even trying to get away from it all.	4.41
PPO: In rutting boredom with boring life; any ideas for shaking things up a little? Any ideas for improving relationships?	4.76
vBoN: Been happily cohabiting for three years with boyfriend of two years. Feeling guilty for turning him into an old man! Any ideas for shaking things up?	6.02
BoN: Boyfriend and I are a little boring, and want to shake things up so we can go out and have fun but we're too stuck on the idea of staying in all day.	5.47

1379 Table 7: An example of summaries sampled at temperature 1. and their corresponding reward obtained from the
1380 evaluator reward model.

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

Content	Reward
SUBREDDIT: r/relationship_advice TITLE: Stuck in a rut and in need of advice/inspiration! POST: My boyfriend and I have been together for 3 years, and living together for 2. I'm quite the homebody, and when we first met, he was very outgoing and loved partying and socialising (although he was a student at the time). We're both working now, and most nights we find ourselves doing the same things: watching series (luckily we enjoy the same shows), playing Minecraft or playing various board games. We're tired after work, and can't bring ourselves to leave the house. The weekend is much the same – lots of sleep, or sitting around staring at one screen or another. We do party occasionally (we'll head to a pub once every few months) and there are a few mutual friends we enjoy spending time with, but I worry that we've become stuck in our boring ways. I really enjoy our lifestyle, and would be quite happy to never leave the house again, but I'm starting to feel guilty for turning him into a 50 year-old when he's only 24. Any ideas for shaking things up a little? Bear in mind that we live in a small town in South Africa, and neither of us has a car.	-
SFT: My boyfriend and I used to party all the time and are ending up reverting to the boring old lifestyle of a 24-year-old. How can we both break out of a rut and get some fun back in our relationship?	5.89
PPO: Taken classic boring suburban lifestyle and turned into boring 24yo. Looking for ideas for breaking out of rut and refreshing lifestyle. Any ideas welcome!	5.91
vBoN: boyfriend and I seem to have fallen into a rut of regularly doing mundane activities after working together non-stop for several months! ideas to spice things up?	6.57
BoN: in a relationship that's getting a bit stale, looking for some inspiration to make changes on a whim in hopes of rejuvenating it!	6.74