

VARIATIONAL BEST-OF- N ALIGNMENT

Afra Amini Tim Vieira Elliott Ash Ryan Cotterell

ETH Zürich

{afra.amini, ryan.cotterell}@inf.ethz.ch

tim.f.vieira@gmail.com ashe@ethz.ch

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.

 <https://github.com/rycolab/vbon>

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 (Meta, 2023; OpenAI, 2023; Gemini, 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 the 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

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.

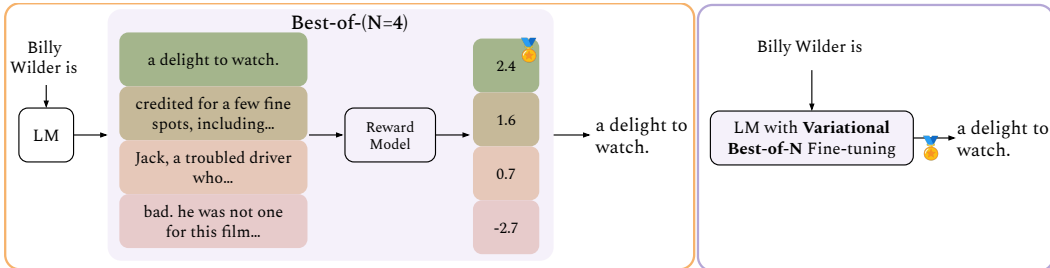


Figure 1: Best-of- N (on the left) is an effective alignment-via-inference method: it draws N samples from the language model, ranks them according to a reward model, and outputs the best sample. Variational Best-of- N (on the right) approximates this process via fine-tuning. The goal is to ensure that 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 .

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 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 (Meta, 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 vBoN on controlled generation and summarization tasks, comparing its performance to models fine-tuned using the KL-constrained RL objective. For controlled generation, our results indicate that models fine-tuned with the vBoN objective are more likely to fall on the Pareto frontier of the reward vs. KL curve compared to other fine-tuning-based alignment methods. This suggests that vBoN achieves a better balance between maximizing reward and maintaining

proximity to the reference model. On a summarization task, fine-tuning with vBoN yields higher reward values and greater win rates on average than models fine-tuned with the KL-constrained RL objective, further demonstrating its effectiveness.

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 $\pi_\theta(\mathbf{y} | \mathbf{x})$. 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,

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

where $D_{\text{KL}}(\cdot)$ is the KL divergence between two distributions, modulated by a hyperparameter β . This objective encourages the model to assign greater probability mass to high-reward outputs while simultaneously penalizing excessive divergence from the reference model. Levine (2018) shows that the policy that maximizes⁵ this objective (Eq. (1)) is

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

In simple terms, π_θ^* is the reference model reweighted by the exponentiated reward values and normalized by the partition function Z . However, direct sampling from π_θ^* is not feasible, as computing Z requires evaluating an infinite sum, making it intractable. 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 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 returns \mathbf{y}^* , where⁶

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

We present the probability distribution induced by BoN with π_{bon} . Notably, π_{bon} is *not* the optimal distribution under Eq. (1), the KL-constrained RL objective.⁷ Despite this, the BoN algorithm often

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.

⁵This formulation implicitly assumes that there exists a $\theta \in \Theta$ that achieves the unconstrained maximum.

⁶We assume that the argmax is unique, or ties are broken in a well-defined manner.

⁷Under simplifying assumptions is π_{bon} asymptotically (in string length) equivalent to π_θ^* (Yang et al., 2024).

performs well—even in comparison to RLHF-based methods. This naturally raises the question: under what optimization objective is π_{bon} the optimal distribution? To answer this question, we first compute the probability of strings under π_{bon} .

Proposition 1. *Suppose $r: \Sigma^* \rightarrow \mathbb{R}$ is a one-to-one mapping. Then, the probability of a string \mathbf{y} under π_{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 variational inference, we minimize the reverse KL divergence between π_{θ} and π_{bon} . Concretely,

$$\mathcal{J}^{\text{vBoN}}(\theta) = -D_{\text{KL}}(\pi_{\theta} \parallel \pi_{\text{bon}}) = \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \left[\log \pi_{\text{bon}}(\mathbf{y}) - \log \pi_{\theta}(\mathbf{y}) \right] \quad (5a)$$

$$= \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \left[\log \pi_{\text{bon}}(\mathbf{y}) \right] + H(\pi_{\theta}) \quad (5b)$$

$$= \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}} \left[\log \sum_{i=1}^N \binom{N}{i} F(r(\mathbf{y}))^{N-i} \pi_{\text{ref}}(\mathbf{y})^i \right] + H(\pi_{\theta}), \quad (5c)$$

where $H(\cdot)$ is the entropy of a distribution. Thus, Eq. (5) offers an answer to the question of what objective BoN optimizes. Inspecting the objective further, we see that 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 a 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 suggests that the vBoN objective may be less sensitive to outliers and the scale of rewards. This property is 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,

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

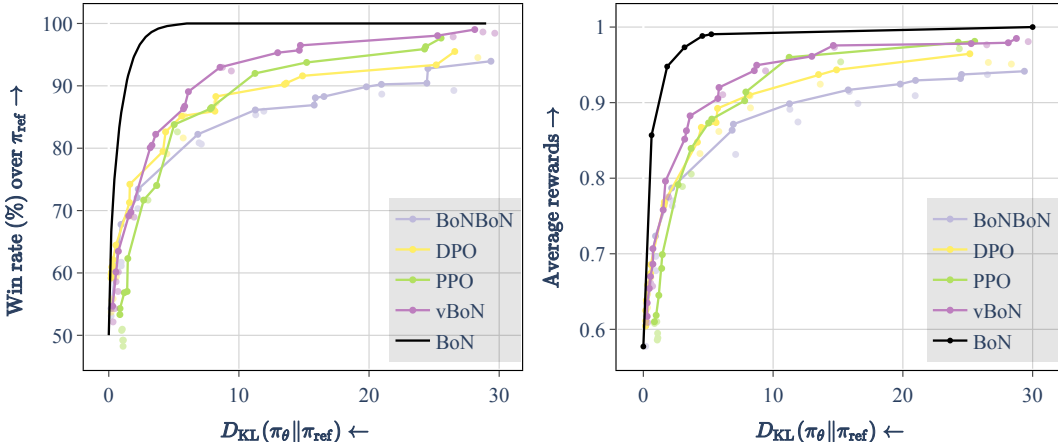
We approximate $F(r(\mathbf{y}))$ using M i.i.d. samples from π_{ref} , termed $\mathbf{y}'^{(1)}, \dots, \mathbf{y}'^{(M)} \stackrel{\text{i.i.d.}}{\sim} \pi_{\text{ref}}$, using which we compute $\hat{F}(r(\mathbf{y})) \stackrel{\text{def}}{=} \frac{1}{M} \sum_{m=1}^M \mathbb{1}\{r(\mathbf{y}'^{(m)}) < r(\mathbf{y})\}$. We then take the log of this Monte Carlo estimator as a biased, but consistent estimator of $\log F(\cdot)$ in Eq. (5).⁸ In §5.1, we empirically assess the number of samples needed for $\log \hat{F}$ to accurately approximate $\log F$.

⁸Using 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 THE BoN 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}^{vBoN} . Through this lower bound, we hope to achieve a deeper insight into 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 $\mathcal{J}^{vBoN}(\theta) \geq L(\theta)$, where

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

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 how the reward function is used. The KL-constrained RL objective aims to maximize the expected reward values, whereas vBoN maximizes the cumulative probability that strings sampled from the aligned model, π_θ , achieve higher rewards compared to those sampled from π_{ref} .

5 SENTIMENT CONTROL

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

The reference model, π_{ref} , is GPT-IMDB⁹, a GPT-2 (Radford et al., 2019) model fine-tuned on IMDB corpus (Maas et al., 2011). We use a binary sentiment classifier,¹⁰ denoted as p , with two classes $\{\text{POS}, \text{NEG}\}$ as the reward model, and define $r(\mathbf{y}) \stackrel{\text{def}}{=} p(\text{POS} \mid \mathbf{y})$. Following Rafailov et al. (2023), we sample 5000 movie reviews from the training set of IMDB dataset and for each sample, we randomly choose a prefix length from $\{2, \dots, 8\}$ and take that prefix as the prompt. We further generate 512 prompts in the same way from the test set of IMDB that we use to evaluate our models.

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

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

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

We fine-tune models by varying the degree of regularization. For BoN approaches, that is achieved by varying N , and for DPO and PPO, we vary β .¹³ Conveniently, N in vBoN is a hyperparameter, meaning that we do *not* need to generate more samples from π_{ref} when we increase N . However, with BoN and BoNBoN methods, we need to increase the number of samples from the reference model as we increase N .

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

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

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

⁹Specifically, we use <https://huggingface.co/lvwerra/gpt2-imdb>.

¹⁰Specifically, we use <https://huggingface.co/lvwerra/distilbert-imdb>.

¹¹One could argue that DPO has a slight advantage over other methods in this setup since it has seen 6 unique generations per prompt during training, while the others only have seen one (or 2 with BoNBoN). Nevertheless, we observe that vBoN is more effective than DPO.

¹²Following the authors’ recommendation, we set α so that both terms contribute equally to the final loss.

¹³See App. F for more details regarding the regularization hyperparameters.

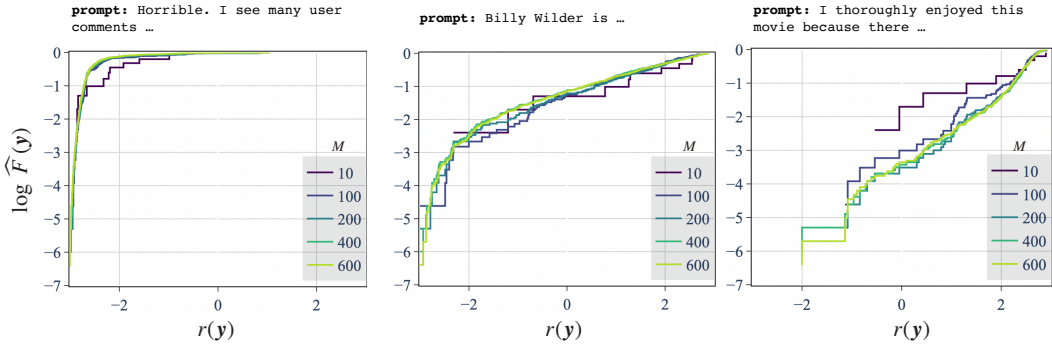


Figure 3: Estimates of $\log F(\cdot)$ with increasing the number of Monte Carlo samples. We test an adversarial prompt (left plot), a neutral prompt (middle plot), and a prompt with a positive sentiment (right plot). Overall, we hardly see any difference between the estimates after taking 200 samples. 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.

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

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

where \widehat{F}_{M_1} and \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-fine-tuning method, and along with other alignment-via-fine-tuning 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-fine-tuning 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.¹⁴ Since the F values are precomputed, we empirically observe that the 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. Our experiments 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 improved by taking into account the similarity between different prompts to learn an approximation to $\log F$ function, which we plan as future work.

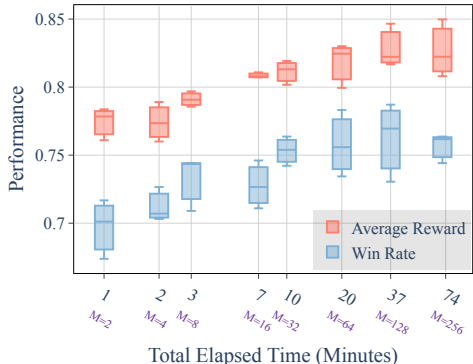


Figure 4: The average reward and win rate of the aligned models improve as we increase the sample size M used for approximating the vBoN loss function.

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).¹⁵ We use SFT to refer to this model in the plots. We use two separate reward models for training and evaluation: a pythia-2.8B¹⁶ reward model for fine-tuning and a larger pythia-6.9B¹⁷ model exclusively for evaluation.

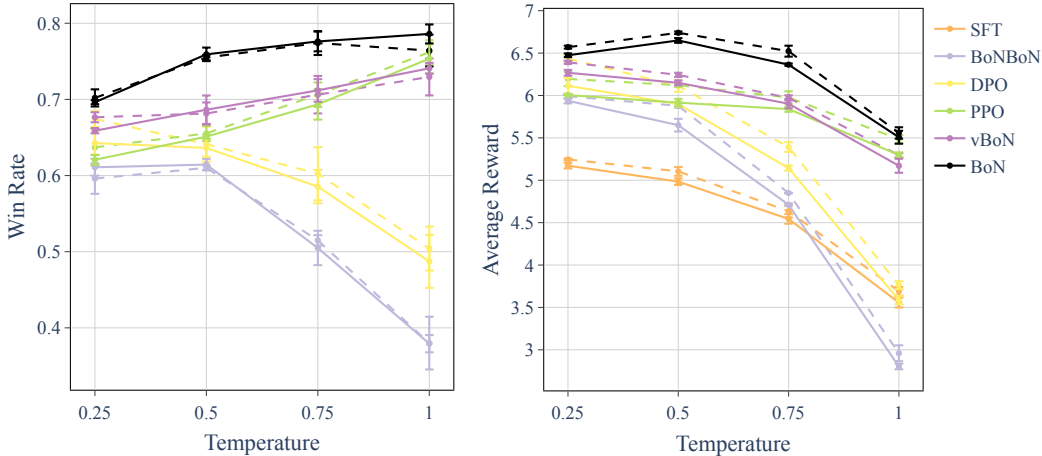
Dataset. To evaluate the generalization ability of the aligned models on out-of-distribution data, we fine-tune the models using only posts from the relationship and relationship_advice subreddits

¹⁴This is particularly advantageous since practitioners often perform the optimization multiple times to test various hyperparameter settings.

¹⁵We use https://huggingface.co/cleanrl/ElleutherAI_pythia-2.8b-deduped__sft__tldr.

¹⁶We use https://huggingface.co/cleanrl/ElleutherAI_pythia-2.8b-deduped__reward__tldr.

¹⁷We use https://huggingface.co/cleanrl/ElleutherAI_pythia-6.9b-deduped__reward__tldr.



(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-fine-tuning 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-fine-tuning 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.

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 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 curve¹⁸ 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,¹⁹ 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 all alignment-via-fine-tuning methods, 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 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).

¹⁸We fit an exponential function of the form $f(x) = -a \exp(-bx)$ to the data using non-linear least squares.

¹⁹The 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

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.

ACKNOWLEDGEMENTS

We thank Ahmad Beirami for the fruitful discussion in the early stages of this project. We also thank Amrit Singh Bedi for identifying a typo in a previous version of the bound derivations. Finally, we thank the anonymous reviewers for their feedback. Afra Amini is supported by the ETH AI Center doctoral fellowship.

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. *Computing Research Repository*, arXiv:2401.01879, 2024. URL <https://arxiv.org/abs/2401.01879>.
- Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V. Le, Christopher Ré, and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling. *Computing Research Repository*, arXiv:2407.21787, 2024. URL <https://arxiv.org/abs/2407.21787>.
- George Casella and Roger L. Berger. *Statistical Inference*. Chapman and Hall/CRC, Pacific Grove, CA, 2nd edition, 2001. ISBN 9781032593036. URL <https://www.routledge.com/Statistical-Inference/Casella-Berger/p/book/9781032593036>.
- Eugene Charniak and Mark Johnson. Coarse-to-fine n -best parsing and MaxEnt discriminative reranking. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, 2005. doi: 10.3115/1219840.1219862. URL <https://aclanthology.org/P05-1022>.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems*, 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/d5e2c0adad503c91f91df240d0cd4e49-Paper.pdf.
- Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan Zhang, Winnie Chow, Rui Pan, Shizhe Diao, Jipeng Zhang, KaShun SHUM, and Tong Zhang. RAFT: Reward ranked finetuning for generative foundation model alignment. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=m7p507zb1Y>.
- Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In *Proceedings of the International Conference on Machine Learning*, Proceedings of Machine Learning Research, 2023. URL <https://proceedings.mlr.press/v202/gao23h.html>.
- Gemini. Gemini: A family of highly capable multimodal models. Technical report, Google, 2024. URL <https://arxiv.org/pdf/2312.11805>.
- Lin Gui, Cristina Gârbacea, and Victor Veitch. BoNBoN alignment for large language models and the sweetness of best-of- n sampling. *Computing Research Repository*, arXiv:2406.00832, 2024. URL <https://arxiv.org/pdf/2406.00832>.
- Alexander Havrilla, Maksym Zhuravynskyi, Duy Phung, Aman Tiwari, Jonathan Tow, Stella Biderman, Quentin Anthony, and Louis Castricato. trIX: A framework for large scale reinforcement learning from human feedback. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2023. doi: 10.18653/v1/2023.emnlp-main.530. URL <https://aclanthology.org/2023.emnlp-main.530>.
- Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger. Deep reinforcement learning that matters. In *Proceedings of the Conference on Artificial Intelligence and Innovative Applications of Artificial Intelligence Conference and AAAI Symposium on Educational Advances in Artificial Intelligence*, 2018. URL <https://dl.acm.org/doi/pdf/10.5555/3504035.3504427>.
- Joseph L. Hodges. The significance probability of the Smirnov two-sample test. *Arkiv för Matematik*, 3, 1958. URL <https://api.semanticscholar.org/CorpusID:121451525>.

- Shengyi Huang, Michael Noukhovitch, Arian Hosseini, Kashif Rasul, Weixun Wang, and Lewis Tunstall. The N+ implementation details of RLHF with PPO: A case study on TL;DR summarization. In *Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=KH02ZTa8e3>.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with PagedAttention. In *Proceedings of the ACM SIGOPS Symposium on Operating Systems Principles*, 2023. URL <https://arxiv.org/abs/2309.06180>.
- Sergey Levine. Reinforcement learning and control as probabilistic inference: Tutorial and review. *Computing Research Repository*, arXiv:1805.00909, 2018. URL <https://arxiv.org/pdf/1805.00909>.
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, 2011. URL <https://aclanthology.org/P11-1015>.
- Meta. Llama 2: Open foundation and fine-tuned chat models. Technical report, Meta, 2023. URL <https://ai.meta.com/research/publications/llama-2-open-foundation-and-fine-tuned-chat-models/>.
- Youssef Mroueh. Information theoretic guarantees for policy alignment in large language models. *Computing Research Repository*, arXiv:2406.05883, 2024. URL <https://arxiv.org/abs/2406.05883>.
- Sidharth Mudgal, Jong Lee, Harish Ganapathy, YaGuang Li, Tao Wang, Yanping Huang, Zhifeng Chen, Heng-Tze Cheng, Michael Collins, Trevor Strohman, Jilin Chen, Alex Beutel, and Ahmad Beirami. Controlled decoding from language models. In *Proceedings of The International Conference on Machine Learning*, Proceedings of Machine Learning Research. PMLR, 2024. URL <https://arxiv.org/pdf/2310.17022>.
- Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, Xu Jiang, Karl Cobbe, Tyna Eloundou, Gretchen Krueger, Kevin Button, Matthew Knight, Benjamin Chess, and John Schulman. WebGPT: Browser-assisted question-answering with human feedback. *Computing Research Repository*, arXiv:2112.09332, 2022. URL <https://arxiv.org/pdf/2112.09332>.
- OpenAI. GPT-4 technical report. Technical report, OpenAI, 2023. URL <https://cdn.openai.com/papers/gpt-4.pdf>.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/b1efde53be364a73914f58805a001731-Paper-Conference.pdf.
- Alizée Pace, Jonathan Mallinson, Eric Malmi, Sebastian Krause, and Aliaksei Severyn. West-of-n: Synthetic preference generation for improved reward modeling. *Computing Research Repository*, arXiv:2401.12086, 2024. URL <https://arxiv.org/abs/2401.12086>.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners, 2019. URL https://d4mucfpksyw.cloudfront.net/better-language-models/language_models_are_unsupervised_multitask_learners.pdf.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. In *Advances in Neural Information Processing Systems*, 2023. URL <https://arxiv.org/pdf/2305.18290.pdf>.

- Tom Schaul, Georg Ostrovski, Iurii Kemaev, and Diana Borsa. Return-based scaling: Yet another normalisation trick for deep RL. *Computing Research Repository*, arXiv:2105.05347, 2021. URL <https://arxiv.org/abs/2105.05347>.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *Computing Research Repository*, arXiv:1707.06347, 2017. URL <https://arxiv.org/abs/1707.06347>.
- Pier Giuseppe Sessa, Robert Dadashi, Léonard Hussenot, Johan Ferret, Nino Vieillard, Alexandre Ramé, Bobak Shariari, Sarah Perrin, Abe Friesen, Geoffrey Cideron, Sertan Girgin, Piotr Stanczyk, Andrea Michi, Danila Sinopalnikov, Sabela Ramos, Amélie Héliou, Aliaksei Severyn, Matt Hoffman, Nikola Momchev, and Olivier Bachem. BOND: Aligning LLMs with best-of-N distillation. *Computing Research Repository*, arXiv:2401.12086, 2024. URL <https://arxiv.org/abs/2401.12086>.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *Computing Research Repository*, arXiv:2408.03314, 2024. URL <https://arxiv.org/abs/2408.03314>.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. In *Advances in Neural Information Processing Systems*, 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/1f89885d556929e98d3ef9b86448f951-Paper.pdf.
- Joy Qiping Yang, Salman Salamatian, Ziteng Sun, Ananda Theertha Suresh, and Ahmad Beirami. Asymptotics of language model alignment. *Computing Research Repository*, arXiv:2404.01730, 2024. URL <https://arxiv.org/pdf/2404.01730>.
- Kevin Yang and Dan Klein. FUDGE: Controlled text generation with future discriminators. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2021. URL <https://aclanthology.org/2021.naacl-main.276>.

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. Meta (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 of a string \mathbf{y} under π_{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:²⁰

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

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

Note that for Eq. (10a) to hold, we need the assumption that the reward function is a one-to-one mapping.²¹ 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 $\mathcal{J}^{\text{vBoN}}(\boldsymbol{\theta}) \geq L(\boldsymbol{\theta})$, where*

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

²⁰The PMF of BoN is also derived by Beirami et al. (Lemma 1; 2024). In their notation, $p_1 = \mathcal{F}$ and $p_2 = \mathcal{F}^{-1}$.

²¹If the reward function is not a one-to-one mapping, we need to devise a tie-breaking strategy. See App. C for further discussion.

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

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

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

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

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

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

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

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

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 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

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

Proof.

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

■

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.

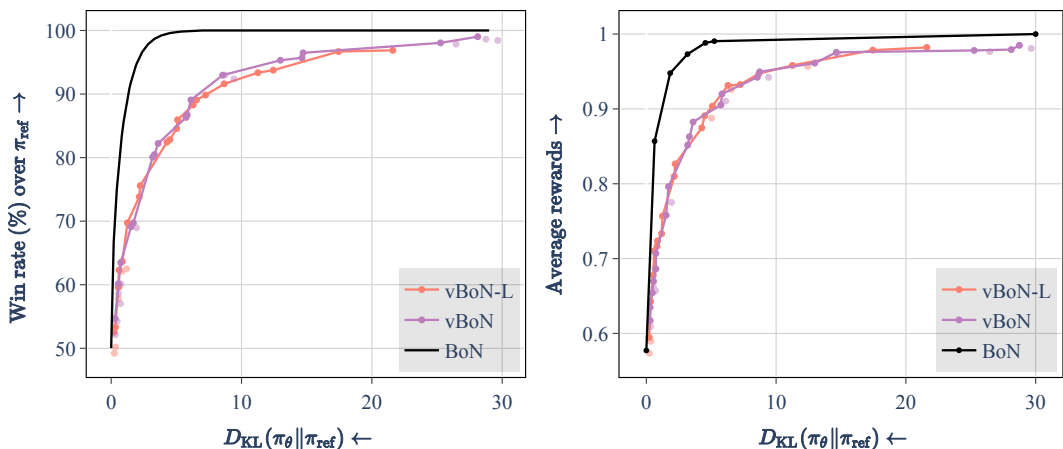
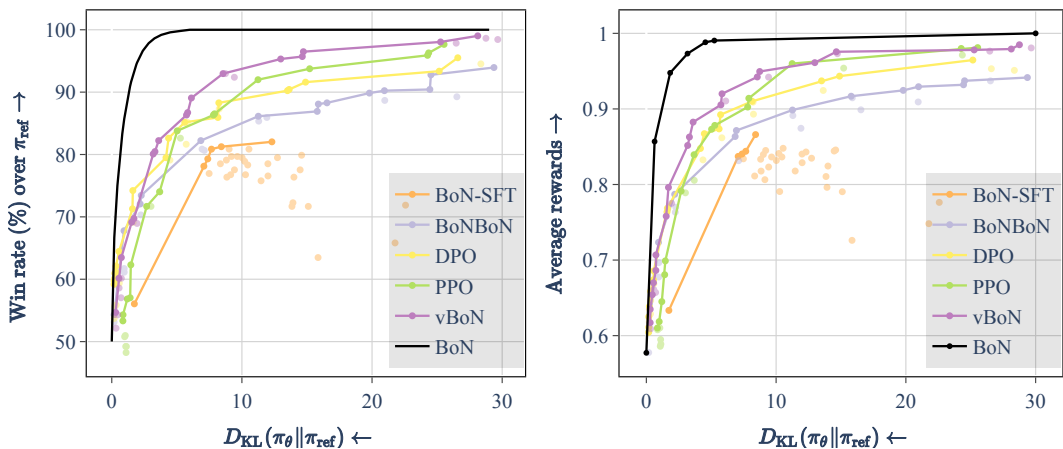


Figure 6: Comparing models trained with the vBoN objective and its lower bound (L). We observe that the performance of the two methods is very close to each other.



(a) 4% of points on Pareto front belong to BoNBoN, 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 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 the 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

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

Content	Reward
<p>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.</p>	-
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

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

Content	Reward
<p>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.</p>	-
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

Table 6: An example of summaries sampled at temperature 0.75 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: 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

Table 7: An example of summaries sampled at temperature 1. 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: 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