

BEST-OF-N THROUGH THE SMOOTHING LENS: KL DIVERGENCE AND REGRET ANALYSIS

Anonymous authors

Paper under double-blind review

ABSTRACT

A simple yet effective method for inference-time alignment of generative models is Best-of- N (BoN), where N outcomes are sampled from a reference policy, evaluated using a proxy-reward model, and the highest-scoring one is selected. While prior work argues that BoN is almost optimal in reward vs KL tradeoffs, the effectiveness of BoN depends critically on the quality of the proxy-reward model used for selection. For this purpose, we study BoN through a smooth version known as Soft Best-of- N (SBoN) and develop a theoretical framework to address this gap. We analyze the scaling behaviour of BoN by providing bounds on the KL divergence between the SBoN policy and the reference policy, offering insights into how performance varies with the number of samples. We also study the regret gap, i.e., the gap between the expected true-reward under the optimal policy and the SBoN policy. Our theoretical and empirical findings show that smoothing helps SBoN mitigate reward overoptimization, especially when the quality of the proxy-reward is low.

1 INTRODUCTION

Large language models (LLMs) have transformed machine learning, achieving state-of-the-art results on a variety of tasks. Despite all advancements, LLMs can still generate undesirable outputs, such as toxic or factually incorrect responses. This has made alignment a central goal in modern LLM development (Achiam et al., 2023; Team et al., 2023).

Several post-hoc alignment methods have been proposed to address this challenge, including Reinforcement Learning from Human Feedback (RLHF) (Christiano et al., 2017; Ouyang et al., 2022), SLiC (Zhao et al., 2022), Direct Preference Optimization (Rafailov et al., 2023), controlled decoding (Mudgal et al., 2024) and Best-of- N (BoN) sampling (Stiennon et al., 2020). While these methods differ in their implementation which are ranging from training-time optimization to test-time selection, they can be viewed as approximating the solution to a KL-regularized reward maximization problem. The optimal solution to this problem is a tilted distribution over responses, which balances reward and proximity to the reference model (Yang et al., 2024).

In BoN as a test-time sampling strategy, given a prompt, N responses are sampled from the reference policy, and the one with the highest proxy-reward sample is selected. Empirically, BoN has been shown to achieve competitive or superior performance in the reward-versus-KL divergence trade-off when compared to RLHF and other alignment methods (Gao et al., 2023a; Mudgal et al., 2024) under *true-reward model*. Furthermore, under certain conditions, it asymptotically approximates the solution to the KL-regularized reward maximization objective (Yang et al., 2024). However, in practice, BoN relies on a learned proxy-reward model which is an approximation of the true-reward function, to guide this selection. As such, their effectiveness critically depends on both the *proxy-reward model* (estimation error of true-reward) and the *quality of the reference policy*.

Understanding how these two components, the quality of the proxy-reward model and the choice of reference policy, affect the alignment quality of test-time sampling algorithms is essential. There are different measures of alignment quality, including KL divergence¹ between aligned policy and reference policy and the *regret* defined as the gap between the expected true-reward under the optimal

¹Unless stated otherwise, all KL divergences are understood to be measured between the aligned policy and the reference policy.

(tilted) policy and the alignment policy. Note that minimizing the regret gap is critical to ensuring high-quality outputs and close performance to the optimal policy. Recent work by Gao et al. (2023a) and Hilton et al. (2022) has investigated the scaling laws governing reward model optimization in both reinforcement learning (RL) and BoN settings as a function of KL divergence between aligned policy and reference policy. They empirically demonstrate that, under proxy-reward models, the improvement in expected true-reward, relative to a reference policy, scales proportionally for both RL and BoN policies.

While recent work analyzes BoN under the idealized settings where there is no discrepancy between the proxy-reward and the true-reward (Yang et al., 2024; Beirami et al., 2024; Mroueh, 2024; Huang et al., 2025), our work relaxes this assumption to study the interplay between the reward discrepancy measured through regret and the KL-divergence. We present a theoretical study of **Soft Best-of- N (SBoN)**, a smoothed variant of BoN recently introduced by Mayrink Verdun et al. (2025); Jinnai et al. (2024). Unlike BoN, SBoN draws the final response *probabilistically* from the N candidates, yielding a policy that is tunable with a temperature parameter. Our analysis centres on two metrics:

- (a) the Kullback-Leibler divergence between SBoN policy (under the *true* reward or proxy-reward model) and the reference policy, and
- (b) the *regret*, i.e. the expected true-reward gap between optimal policy and SBoN policy.

We show how these results specialize to the BoN (as a limit of SBoN for the temperature goes to infinity) and quantify the estimation error incurred by using a *proxy* reward model instead of the true-reward. Finally, we characterize regimes in which SBoN attains *lower* regret bound than BoN when we use the proxy-reward model. Our main contributions are:

- We derive finite-sample upper bounds for KL divergence between the SBoN policy and reference policy, and upper and lower bounds for the regret gap of the SBoN policy, and we extend these bounds to BoN. These bounds reveal how the number of responses N , proxy-reward model quality and reference policy model affect performance.
- We quantify cases where SBoN performs better than BoN under overoptimization scenario where the proxy-reward model is used instead of the true-reward model.
- We provide experimental validation using various proxy-reward models to demonstrate SBoN’s advantages in the overoptimization scenario. Furthermore, we provide numerical experiments to evaluate our bounds.

2 RELATED WORKS

In this section, we discuss related works on BoN, the theoretical foundation of (Soft) BoN and overoptimization. More related works for the theoretical foundation of RLHF and smoothing of maximum are provided in the Appendix (App) B.

Best-of- N : Despite many recent advancements in alignment, a simple, popular, and well-performing method continues to be the BoN policy (Nakano et al., 2021; Stiennon et al., 2020; Beirami et al., 2024). In fact, Gao et al. (2023b); Mudgal et al. (2024); Eisenstein et al. show that BoN consistently achieves compelling win rate-KL tradeoff curves, often outperforming KL-regularized reinforcement learning and other more complex alignment strategies. LLaMA 2 (Touvron et al., 2023) leverages BoN outputs as teacher signals to further finetune the base model. Mudgal et al. (2024) extend BoN through Q-learning to block-wise BoN decoding. This empirical effectiveness has also inspired research into distilling BoN behaviour into standalone models (Amini et al., 2025; Sessa et al., 2024; Qiu et al., 2024). Hughes et al. (2024) utilize BoN as an effective method for jailbreaking, while BoN is also commonly used as a strong baseline for scaling inference-time compute (Brown et al., 2024; Snell et al., 2024). Given the broad success of BoN, we are motivated to theoretically investigate the BoN policies and the effect of the proxy-reward model (reward hacking) and the quality of the reference policy.

Theoretical Foundation of (Soft) BoN: KL divergence of BoN is studied in (Beirami et al., 2024; Mroueh, 2024) via information theoretical tools where the KL divergence of BoN sampling from the reference distribution is bounded by $\log(N) - (N - 1)/N$. Scaling laws governing reward as a function of KL divergence is empirically studied by Gao et al. (2023b) and theoretically formalized by Mroueh (2024). Furthermore, the asymptotic case and the equivalence of BoN to the KL-constrained reinforcement learning solution are studied by Yang et al. (2024) under the assump-

tion of access to optimal reward. Gui et al. (2024) further characterized the win rate–KL gap in the asymptotic regime where a model assigns extremely low likelihoods to successful completions. Furthermore, Sun et al. (2024) accelerated BoN using speculative rejection sampling. The regret of BoN under some assumptions is studied in (Huang et al., 2025). The convergence rate of the SBoN policy to the optimal tilted policy has been analyzed by Mayrink Verdun et al. (2025). Additionally, Geuter et al. (2025) investigate a variant of SBoN that incorporates speculative samples from a small auxiliary model, providing both theoretical and empirical insights. However, the regret gap and KL divergence of SBoN under overoptimization scenario remain largely unexplored in the existing literature.

3 PROBLEM FORMULATION

Notations: Upper-case letters denote random variables (e.g., Z), lower-case letters denote the realizations of random variables (e.g., z), and calligraphic letters denote sets (e.g., \mathcal{Z}). All logarithms are in the natural base. The set of probability distributions (measures) over a space \mathcal{X} with finite variance is denoted by $\mathcal{P}(\mathcal{X})$. Δ_N is N -simplex distribution set. The KL divergence between two probability distributions on \mathbb{R}^d with densities $p(x)$ and $q(x)$, such that $q(x) > 0$ when $p(x) > 0$, is $\text{KL}(p||q) := \int_{\mathbb{R}^d} p(x) \log(p(x)/q(x))dx$ (with $0 \cdot \log 0 := 0$). The total-variation distance is defined as $\mathbb{TV}(p, q) = \frac{1}{2} \int_{\mathcal{X}} |p(x) - q(x)|$. Furthermore, we define chi-square divergence as $\chi^2(p(x)||q(x)) = \int_{\mathcal{X}} \frac{p^2(x)}{q(x)} - 1$.

Preliminaries: Let the finite set² of prompts be \mathcal{X} and the discrete finite set of responses be \mathcal{Y} . Prompts are drawn from a distribution ρ over \mathcal{X} . A (stochastic) policy $\pi \in \Pi$ assigns, for every prompt $x \in \mathcal{X}$, a conditional distribution $\pi(\cdot | x)$ over \mathcal{Y} ; drawing $y \sim \pi(\cdot | x)$ yields a response. We treat the supervised-fine-tuned (SFT) model as a *reference policy*, denoted $\pi_{\text{ref}}(\cdot | x)$.

3.1 REWARD FUNCTION

We consider a (calibrated) true-reward function $r^*(y, x)$ and a (calibrated) proxy-reward function $\hat{r}(y, x)$, both mapping $\mathcal{Y} \times \mathcal{X} \rightarrow [0, 1]$ ³. The true-reward function is estimated via proxy-reward function $\hat{r}(y, x)$ using *some preference datasets*⁴. As shown by Balashankar et al. (2025), a calibrated reward function satisfies:

- Boundedness: for all x, y , we have $\hat{r}(y, x), r^*(y, x) \in [0, 1]$.
- Uniformity under the reference model: for each prompt $x \in \mathcal{X}$, if $Y \sim \pi_{\text{ref}}(\cdot | x)$ then $r(y, x) \sim \text{Unif}(0, 1)$ for $r \in \{\hat{r}, r^*\}$.

In practice, the proxy-reward model can be fit to a human-labeled preference dataset or to data annotated with true-rewards. Following (Huang et al., 2025), we assume for simplicity that $\hat{r}(y, x)$ is given. We define $\mathcal{Y}_{r^*}^*(x) = \arg \max_y r^*(y, x)$ as the set of maximizers for true-reward. Similarly, we can define $\hat{\mathcal{Y}}_{\hat{r}}(x)$ as the set of maximizers for proxy-reward model. More discussion regarding reward is provided in App. C.

Assumption 3.1 (Achievable maximum reward). *We assume that for $r \in \{\hat{r}, r^*\}$, we have $r(\hat{y}(x), x) = 1$ for all $\hat{y}(x) \in \arg \max_y r(y, x)$ and given $x \in \mathcal{X}$.*

In many settings, the reward function attains its maximum at specific responses. In particular, since a large language model (LLM) generates outputs using a finite vocabulary and a bounded number of tokens, the space of possible generations is finite, and thus the assumption holds trivially.

²For measure-theoretic simplicity and notational convenience, we assume finiteness for the set of prompts. Our results also hold for non-finite set.

³For the remainder of this paper, we assume all reward functions are calibrated (see App. C). For brevity, we refer to them simply as *reward functions*.

⁴In some cases, the reward model is not derived from human preference data. Instead, it is either deterministic (e.g., code execution scores) or provided by an automated classifier (e.g., for toxicity or sentiment).

3.2 SBoN ALGORITHM

Fix a prompt $x \in \mathcal{X}$ and draw N i.i.d. candidates $Y_{1:N} \sim \pi_{\text{ref}}(\cdot | x)$. Let $Z \in \{1, \dots, N\}$ denote the index of the selected response with distribution P_Z ; write $P_Z(i) = \Pr(Z = i)$. We seek a distribution over indices that maximizes the proxy-reward:

$$\max_{P_Z \in \Delta_N} \mathbb{E}_Z [\hat{r}(Y_Z, x)].$$

Without regularization, the optimizer is the deterministic Best-of- N (BoN) rule $P_Z = \delta_{i^*}$ with $i^* \in \arg \max_i \hat{r}(Y_i, x)$. Because \hat{r} is a *proxy* for the true-reward, this deterministic choice can overoptimize the proxy-reward and get response with lower true-reward. To smooth this, we add an entropy penalty with temperature $\beta > 0$:

$$\max_{P_Z \in \Delta_N} \mathbb{E}_Z [\hat{r}(Y_Z, x)] + \frac{1}{\beta} H(P_Z).$$

The unique solution is the softmax distribution, $P_Z(i) = \frac{\exp(\beta \hat{r}(Y_i, x))}{\sum_{j=1}^N \exp(\beta \hat{r}(Y_j, x))}$.

We then sample Z from this distribution and return Y_Z . We refer to this sampling rule as **Soft-BoN**, as introduced by Mayrink Verdun et al. (2025). *Note that the value of β parameter for SBoN is equal to β in tilted optimal policy. It is shown by Mayrink Verdun et al. (2025), the SBoN policy can be interpreted as a finite-sample approximation of the tilted optimal policy.*

We denote the final policy from SBoN via $\pi_{\hat{r}}^{\text{SBoN}}(y|x)$. Note that for $\beta \rightarrow \infty$ and $\beta \rightarrow -\infty$, we recover BoN and worst-of- N (WoN) (Balashankar et al., 2025), respectively. Furthermore, for $\beta \rightarrow 0$, we recover uniform sampling among the N response samples, which is equivalent to sampling from the reference model $\pi_{\text{ref}}(y|x)$. In (Mayrink Verdun et al., 2025, Lemma 1), the closed form solution of SBoN policy is derived,

$$\pi_{\hat{r}}^{\text{SBoN}}(y|x) = \frac{\pi_{\text{ref}}(y|x) \exp(\beta \hat{r}(y, x))}{Z_{N,\beta}}, \quad (1)$$

where $Z_{N,\beta} = \mathbb{E} \left[\left(\frac{1}{N} (\exp(\beta \hat{r}(y, x)) + \sum_{i=1}^{N-1} \exp(\beta \hat{r}(Y_i, x))) \right)^{-1} \right]^{-1}$. Similarly, we can define $\pi_{r^*}^{\text{SBoN}}(y|x)$ based on a true-reward model. For simplicity, we define BoN policies under true-reward and proxy-reward models as $\pi_{r^*}^{\text{BoN}}(y|x)$ and $\pi_{\hat{r}}^{\text{BoN}}(y|x)$, respectively. In this work, we focus on $\beta \geq 0$. Another motivation for SBoN based on the Gumbel-Max trick is provided in App. E.

3.3 OPTIMAL (TILTED) POLICY

For a given temperature $\beta > 0$, we seek a policy that remains close to π_{ref} while maximizing expected true-reward, leading to the KL-regularized objective

$$\max_{\pi \in \Pi} \mathbb{E}_{Y \sim \pi(\cdot|x)} [r^*(y, x)] - \frac{1}{\beta} \text{KL}(\pi(\cdot | x) \| \pi_{\text{ref}}(\cdot | x)). \quad (2)$$

The unique solution is the optimal *tilted* policy (Korbak et al., 2022b;a; Yang et al., 2024)

$$\pi_{\beta, r^*}(y|x) = \frac{\pi_{\text{ref}}(y|x) \exp(\beta r^*(y, x))}{Z_{r^*, Y}(x, \beta)}, \quad (3)$$

where $Z_{r^*, Y}(x, \beta) = \sum_{y \in \mathcal{Y}} \pi_{\text{ref}}(y|x) \exp(\beta r^*(y, x))$, is the normalizing (*partition*) function.

Note that, in practice, we do not have access to the closed form of reference policy $\pi_{\text{ref}}(y|x)$ and $r^*(y, x)$. We can only first estimate the true-reward function via proxy-reward function $\hat{r}(y, x)$ and then sample from $\pi_{\text{ref}}(y|x)$ and compute $\hat{r}(y, x)$ for each individual sample. Finally, we can apply inference time algorithms, e.g., BoN or SBoN (Mayrink Verdun et al., 2025), where N samples are generated from $\pi_{\text{ref}}(y|x)$ and we choose the sample with the highest proxy-reward function (BoN) or sampled from a distribution (SBoN) using the proxy-reward function. When only a proxy-reward function $\hat{r}(y, x)$ is available, we obtain the analogous partition function $Z_{\hat{r}, Y}(x, \beta)$ and policy $\pi_{\beta, \hat{r}}(\cdot|x)$.

We can also define the optimal policy under the true-reward model as,

$$\pi_{r^*}^*(y|x) = \arg \max_{\pi} \mathbb{E}_{Y \sim \pi(\cdot|x)} [r^*(Y, x)].$$

Similarly, we can define $\pi_{\hat{r}}^*(y|x)$ as the optimal policy under the proxy-reward model⁵.

As the reward functions (true and proxy) are bounded due to calibration, we can interpret optimal policies as the limit of tilted optimal policies for $r \in \{\hat{r}, r^*\}$,

$$\pi_{\infty, r}(\cdot|x) := \lim_{\beta \rightarrow \infty} \pi_{\beta, r}(\cdot|x)$$

where $\pi_{\infty, r^*}(\cdot|x)$ and $\pi_{\infty, \hat{r}}(\cdot|x)$ place all their probability mass on the maximizers of $r^*(y, x)$ and $\hat{r}(y, x)$, respectively. The connections between SBoN, BoN, optimal and optimal tilted policies under true or proxy-reward models are shown in Figure 1.

3.4 TILTED ERROR

Let’s define the tilted error as the tilted average of square estimation error of true-reward function for a given prompt x with parameter β ,⁶ as follows,

$$\varepsilon_{\beta, r}(x) := \frac{1}{\beta} \log \left(\mathbb{E}_{Y \sim \pi_{\text{ref}}(y|x)} [e^{\beta(r^*(Y, x) - \hat{r}(Y, x))^2}] \right). \quad (4)$$

A similar definition of estimation error is introduced in (Yang & Wibisono, 2022). When $\beta = 0$, the definition reduces to the mean-squared error, which is also introduced in (Huang et al., 2025). Letting $\beta \rightarrow \infty$ recovers the square of the supremum (infinity) norm ($\|\cdot\|_{\infty}$) of the estimation error between $r^*(y, x)$ and $\hat{r}(y, x)$. Therefore, the following properties hold for $\varepsilon_{\beta, r}(x)$,

- The tilted error is bounded, i.e., $\varepsilon_{\beta, r}(x) \in [0, 1]$.
- The tilted average of the estimation error is monotonically increasing in β .
- $\varepsilon_{\infty, r}(x) := \lim_{\beta \rightarrow \infty} \varepsilon_{\beta, r}(x) = \|r^*(Y, x) - \hat{r}(Y, x)\|_{\infty}^2$.

We assume that overoptimization regime happens whenever we have $\varepsilon_{\beta, r}(x) > 0$. In more details, $\varepsilon_{\beta, r}(x) > 0$ indicates reward misspecification. As noted in Gao et al. (2023b), misspecification is a necessary precondition for overoptimization, but overoptimization itself refers to the specific regime where the optimization strength is sufficiently high that the degradation due to error outweighs the improvement in the proxy-reward.

We define tilted error using (proxy and true) reward models rather than raw reward models, because our focus is on how rankings change under the proxy. For example, if the proxy-reward is a strictly increasing transform of the true-reward, the ranking is preserved; the Best-of-N (BoN) policy remains optimal and no overoptimization occurs. This behavior cannot be captured when working with the raw (uncalibrated) reward models. Note that in (Huang et al., 2025), raw (uncalibrated) reward models are utilized for error definition.

3.5 COVERAGE

For a given reward function $r(x, y)$, we define the tilted policy (softmax policy):

$$\pi_{\beta, r}(y|x) \propto \pi_{\text{ref}}(y|x) \exp(\beta r(y, x)).$$

⁵This policy, $\pi_{\hat{r}}^*(y|x)$ maximizes \hat{r}_c , it may be suboptimal or harmful under r^* due to Goodhart’s Law Gao et al. (2023b).

⁶To simplify notation, we adopt the same smoothing parameter β for the error definition as is used in the KL and tilted cases.

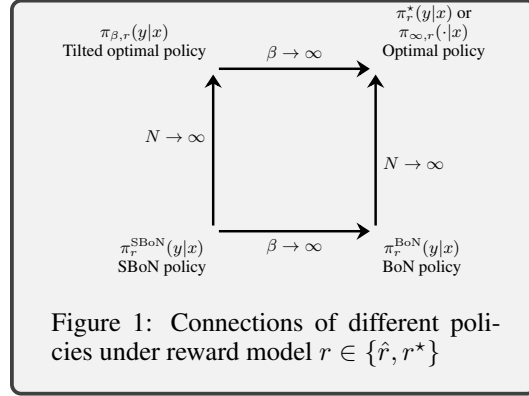


Figure 1: Connections of different policies under reward model $r \in \{\hat{r}, r^*\}$

Then, we introduce the coverage of tilted policy with respect to the reference policy as,

$$C_{\beta,r}(x) := \sum_{y \in \mathcal{Y}} \frac{\pi_{\beta,r}^2(y|x)}{\pi_{\text{ref}}(y|x)}. \quad (5)$$

We also define,

$$C_{\infty,\hat{r}}(x) := \lim_{\beta \rightarrow \infty} C_{\beta,r}(x).$$

This measure $C_{\beta,r}(x)$ can also be interpreted as a coverage constant, which is standard in KL-regularized policy learning. Furthermore, we can define the coverage of the tilted policy with respect to the reference policy as χ^2 -divergence between $\pi_{\beta,r}(y|x)$ and $\pi_{\text{ref}}(y|x)$, i.e., $\chi^2(\pi_{\beta,r}(y|x) \parallel \pi_{\text{ref}}(y|x))$. It ensures that the reference policy places sufficient probability mass on high-reward responses, thereby guaranteeing that the support of the optimal policy lies within the support of the reference. This prevents cases where optimal outputs are entirely excluded by the reference. Similar notions of coverage have been explored in [Huang et al. \(2025\)](#).

3.6 REGRET

For a given policy $\pi(Y|x)$, we define expected true-reward with respect to the policy (a.k.a. value function⁷) as

$$J_{r^*}(\pi(\cdot|x)) := \mathbb{E}_{Y \sim \pi(\cdot|x)}[r^*(Y, x)]. \quad (6)$$

For two policies, $\pi_1(\cdot|x)$ and $\pi_2(\cdot|x)$, we define the gap between these two policies as follows,

$$\Delta_{J_{r^*}}(\pi_1(\cdot|x), \pi_2(\cdot|x)) := J_{r^*}(\pi_1(\cdot|x)) - J_{r^*}(\pi_2(\cdot|x)). \quad (7)$$

We provide an upper bound on the gap of the SBoN solution, which is the gap between $\pi_{r^*}^*(\cdot|x)$ as the optimal policy and $\pi_{\hat{r}}^{\text{SBoN}}(\cdot|x)$,

$$\Delta_{J_{r^*}}(\pi_{r^*}^*(\cdot|x), \pi_{\hat{r}}^{\text{SBoN}}(\cdot|x)) = J_{r^*}(\pi_{r^*}^*(\cdot|x)) - J_{r^*}(\pi_{\hat{r}}^{\text{SBoN}}(\cdot|x)). \quad (8)$$

Regarding regret of the BoN, we consider $\pi_{\hat{r}}^{\text{BoN}}(\cdot|x)$ instead of $\pi_{\hat{r}}^{\text{SBoN}}(\cdot|x)$ in equation 8.

4 KL DIVERGENCE ANALYSIS

The KL divergence between the aligned policy and the reference policy, $\text{KL}(\pi_{r^*}^{\text{BoN}} \parallel \pi_{\text{ref}})$, is studied by [Beirami et al. \(2024\)](#); [Mroueh \(2024\)](#) from a theoretical perspective. In particular, [Beirami et al. \(2024\)](#) derives an upper bound on KL divergence for BoN policies under the assumptions of a bijective true-reward mapping and a finite output space:

$$\text{KL}(\pi_{r^*}^{\text{BoN}}(\cdot|x) \parallel \pi_{\text{ref}}(\cdot|x)) \leq \log(N) - 1 + \frac{1}{N}, \quad (9)$$

[Mroueh \(2024\)](#) relaxes the bijectivity assumption and derives similar bounds using information-theoretic tools. Under some assumptions, the bound in equation 9 is tight. Furthermore, using Pinsker's inequality, in a similar approach to ([Mroueh, 2024](#)), we have,

$$\mathbb{E}_{Y \sim \pi_{r^*}^{\text{SBoN}}(\cdot|x)}[r^*(Y, x)] \leq 0.5 + \sqrt{\frac{1}{2} \text{KL}(\pi_{r^*}^{\text{SBoN}}(\cdot|x) \parallel \pi_{\text{ref}}(\cdot|x))}. \quad (10)$$

Note that equation 10 implies that improvement of expected true-reward relative to the reference policy can not exceed the square root of the KL divergence. However, the analysis of KL divergence for the SBoN policy under the true-reward model is overlooked. Therefore, we first establish an upper bound on the KL divergence between the SBoN policy under the true-reward model and the reference policy, shedding light on its behaviour as a function of the number of samples N and temperature parameter β . All proof details are deferred to App. G.

⁷We can also consider $\mathbb{E}_{X \sim \rho(\cdot)}[J_{r^*}(\pi(\cdot|X))]$. All of our results also hold for expected version of value function.

Lemma 4.1. *The following upper bound holds on KL divergence between SBoN and reference policies for a given prompt $x \in \mathcal{X}$,*

$$\text{KL}(\pi_{r^*}^{\text{SBoN}}(y|x) \parallel \pi_{\text{ref}}(y|x)) \leq \log \left(\frac{N}{1 + (N-1) \exp(-\beta)} \right). \quad (11)$$

Using Lemma 4.1, we can observe that for BoN, $\beta \rightarrow \infty$, we have,

$$\text{KL}(\pi_{r^*}^{\text{BoN}}(y|x) \parallel \pi_{\text{ref}}(y|x)) \leq \log(N). \quad (12)$$

Comparing equation 12 with results in (Beirami et al., 2024; Mroueh, 2024), our result is derived from the SBoN asymptotic regime. Note that our bound is looser than the bound on KL divergence in equation 9. In contrast, our bound is general and can be applied to different β in SBoN. For $\beta = 0$, where our policy is the reference policy, our bound is tight. It is also important to note that the upper bound in Lemma 4.1 increases with the temperature parameter β for fixed N .

Recent works by Gao et al. (2023a) and Hilton et al. (2022) empirically demonstrate that, under a true-reward model, the improvement in expected true-reward, relative to a reference policy, *scales approximately proportionally to* $\sqrt{\text{KL}(\pi_{r^*}^{\text{BoN}} \parallel \pi_{\text{ref}})}$ for both RL and BoN policies. It is also observed by Gao et al. (2023b) that models optimized using proxy-rewards can suffer from overoptimization where the learned policy diverges further from the reference, the alignment may degrade. Despite theoretical advances, the KL divergence analysis for SBoN and BoN under the proxy-reward model remains largely unexplored. Therefore, we are interested in investigating the cost we have for estimation error of true-reward via proxy-reward model. For this aim, we first propose the following useful Lemma to study the closeness of the SBoN policy under the true-reward model to the SBoN policy under the proxy-reward model in KL divergence measure.

Lemma 4.2. *The following upper bound holds on the KL divergence between the SBoN policies under true-reward and proxy-reward models respectively,*

$$\text{KL}(\pi_{r^*}^{\text{SBoN}}(\cdot|x) \parallel \pi_{\hat{r}}^{\text{SBoN}}(\cdot|x)) \leq \frac{N\beta\sqrt{\varepsilon_{\beta,r}(x)}}{1 + (N-1) \exp(-\beta)} \left(\frac{N \exp(2\beta)}{(N-1)^2} + 1 \right). \quad (13)$$

Note that for $\beta = 0$, the upper bound in Lemma 4.2 is tight. The result in Lemma 4.2 quantifies the estimation error introduced by substituting a proxy-reward model for the true-reward model.

Next, we compare BoN and SBoN under overoptimization from KL-divergence perspective.

Remark 4.3 (No overoptimization). *We can observe that for a given β , if we assume $\varepsilon_{\beta,r}(x) = 0$, we have $\text{KL}(\pi_{r^*}^{\text{SBoN}}(\cdot|x) \parallel \pi_{\hat{r}}^{\text{SBoN}}(\cdot|x)) = 0$. Note that, as mentioned in (Gao et al., 2023b), the expected true-reward under the aligned policy, relative to the reference policy, is proportional to the square root of KL divergence. Then a larger KL divergence is desirable in this context, as proposed by (Gao et al., 2023b), the BoN policy is preferred under no overoptimization scenario.*

Remark 4.4 (Overoptimization). *When $\varepsilon_{\beta,r}(x) > 0$, we have two conflicting goals in both Lemma 4.1 and Lemma 4.2: one suggesting for fixed N that β needs to be smaller for better estimation of the true policy by the proxy-reward model one given in Lemma 4.2, and another one suggesting a larger β to induce a better KL trade-off based on Gao et al. (2023b). Hence, for a given N , there exists an optimal β to balance between the estimation error of Lemma 4.2 and the scaling law under the SBoN policy for the true-reward model, Lemma 4.1. In this scenario, SBoN can lead to better tradeoffs than BoN. A similar discussion can be done for fixed β and varying N . Further analysis using reverse I-projection (Csiszár & Matus, 2003) in large N regime is provided in App. K.*

5 REGRET ANALYSIS

In this section, we derive theoretical regret bounds, *upper and lower bounds*, for SBoN and BoN based on reward models. First, we provide a helpful Lemma regarding the expected coverage assumption that can help us interpret the results of regret for BoN and SBoN. All proof details are deferred to App. H.

Lemma 5.1. *Under Assumption 3.1, it holds that $C_{\infty, \hat{r}}(x) = \frac{1}{\sum_i \pi_{\text{ref}}(y_{i,r}^{\max}(x)|x)}$, where $y_{i,r}^{\max}(x) \in \arg \max_y r(y, x)$.*

5.1 UPPER BOUND ON REGRET

Now, we derive an upper bound on the regret of SBoN.

Theorem 5.2 (Upper Bound on Regret of SBoN). *Under Assumption 3.1, the following upper bound holds on the optimal regret gap of the SBoN policy for any $\beta > 0$,*

$$\begin{aligned} \Delta_{J_{r^*}}(\pi_{r^*}^*(\cdot | x), \pi_{\hat{r}}^{\text{SBoN}}(\cdot | x)) &\leq \sqrt{\varepsilon_{\beta, r}(x)}(\sqrt{C_{\infty, \hat{r}}(x)} + \sqrt{C_{\infty, r^*}(x)}) \\ &\quad + 2\sqrt{\frac{1}{2} \log\left(1 + \frac{C_{\infty, \hat{r}}(x) - 1}{N}\right)} + \frac{\log(C_{\infty, r^*}(x))}{\beta}. \end{aligned}$$

Regret Bound of BoN Through Smoothing Lens: We now derive an upper bound on the regret of BoN by taking the asymptotic limit of the regret bound on regret of SBoN in Theorem 5.2.

Proposition 5.3 (Upper Bound on Regret of BoN). *Under Assumption 3.1, the following upper bound holds on the optimal regret gap of the BoN policy for any $\beta > 0$,*

$$\begin{aligned} \Delta_{J_{r^*}}(\pi_{r^*}^*(\cdot | x), \pi_{\hat{r}}^{\text{BoN}}(\cdot | x)) &\leq \sqrt{\varepsilon_{\infty, r}(x)}(\sqrt{C_{\infty, \hat{r}}(x)} + \sqrt{C_{\infty, r^*}(x)}) \\ &\quad + 2\sqrt{\frac{1}{2} \log\left(1 + \frac{C_{\infty, \hat{r}}(x) - 1}{N}\right)}. \end{aligned}$$

Remark 5.4 (Comparison with (Huang et al., 2025)). *The regret bound for BoN policy grows with the L_∞ -norm of the reward-model estimation error. In contrast to the result in (Huang et al., 2025), our bound remains finite whenever the overoptimization error vanishes, i.e., when $\varepsilon_{\infty, \beta}(x) = 0$ or N grows. We also derive results based on calibrated reward, instead of raw (uncalibrated) reward models. A full comparison with Huang et al. (2024) is provided in App. J.*

Overoptimization (asymptotic regime): Assume that $\varepsilon_{\beta, r}(x) > 0$ for every $\beta > 0$. Letting $N \rightarrow \infty$ and invoking Theorem 5.2, we obtain

$$\begin{aligned} \Delta_{J_{r^*}}(\pi_{r^*}^*(\cdot | x), \pi_{\hat{r}}^{\text{SBoN}}(\cdot | x)) &\leq \frac{\log C_{\infty, r^*}(x)}{\beta} \\ &\quad + \sqrt{\varepsilon_{\beta, r}(x)}(\sqrt{C_{\infty, \hat{r}}(x)} + \sqrt{C_{\infty, r^*}(x)}). \end{aligned} \tag{14}$$

Similarly, for BoN we have,

$$\Delta_{J_{r^*}}(\pi_{r^*}^*(\cdot | x), \pi_{\hat{r}}^{\text{BoN}}(\cdot | x)) \leq \sqrt{\varepsilon_{\infty, r}(x)}(\sqrt{C_{\infty, \hat{r}}(x)} + \sqrt{C_{\infty, r^*}(x)}). \tag{15}$$

No overoptimization: Assume that the overoptimization vanishes, i.e. $\varepsilon_{\beta, r}(x) = 0$ for every $\beta \in [0, \infty)$. Then the optimality gaps of the SBoN and BoN policies satisfy

$$\Delta_{J_{r^*}}(\pi_{r^*}^*(\cdot | x), \pi_{\hat{r}}^{\text{SBoN}}(\cdot | x)) \leq 2\sqrt{\frac{1}{2} \log\left(1 + \frac{C_{\infty, \hat{r}}(x) - 1}{N}\right)} + \frac{\log C_{\infty, r^*}(x)}{\beta}, \tag{16}$$

$$\Delta_{J_{r^*}}(\pi_{r^*}^*(\cdot | x), \pi_{\hat{r}}^{\text{BoN}}(\cdot | x)) \leq 2\sqrt{\frac{1}{2} \log\left(1 + \frac{C_{\infty, \hat{r}}(x) - 1}{N}\right)}. \tag{17}$$

5.2 LOWER BOUND ON REGRET

In this section, we complement the regret upper bounds (Theorem 5.2 and Proposition 5.3) with lower bounds that hold for any finite N and fixed $\beta \geq 0$. For the lower bound, the following assumptions are needed.

Assumption 5.5 (Margin Assumption). *Let $\gamma(x) = 1 - \sup_{y \notin \mathcal{Y}_{r^*}^*(x)} r^*(y, x)$. We assume that $\gamma(x) \in (0, 1)$.*

Note that Assumption 5.5 for strictly positive lower bound is needed.

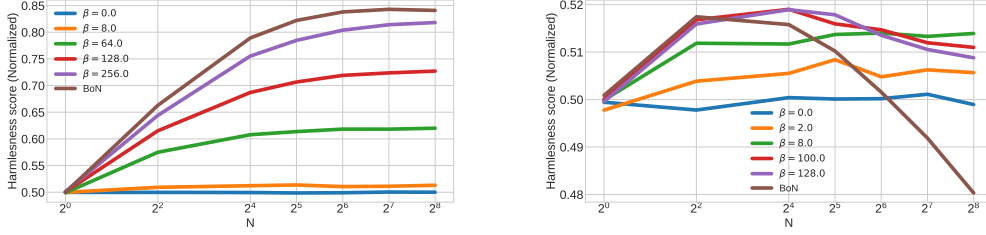


Figure 2: Soft Best-of-N experiment using a strong reward model (Left) and a weak one (Right). Number of samples versus Harmless score (higher is better). When the performance of the reward model is poor, BoN can lead to overoptimization, while the SBoN can help to mitigate it.

Theorem 5.6 (Lower Bound on Regret of SBoN). *Under the same Assumptions in Theorem 5.2 and Lemma H.3, the following lower bound holds on the regret of SBoN policy under the proxy-reward model,*

$$\begin{aligned} & \Delta_{J,r^*}(\pi_{r^*}^*(\cdot|x), \pi_{\hat{r}}^{\text{SBoN}}(\cdot|x)) \\ & \geq \gamma(x) \left((1 - C_{\infty,r^*}(x)^{-1})^N + \frac{N((1 - C_{\infty,r^*}(x)^{-1}) - (1 - C_{\infty,r^*}(x)^{-1})^N)}{N + (\exp(\beta) - 1)(N - 1)} \right) \\ & \quad - \min \left(\sqrt{1 - \exp(-U(N, \beta))}, \frac{N}{1 + (N - 1) \exp(-\beta)} \sqrt{\varepsilon_{\beta,r}(x)} \right) \end{aligned} \quad (18)$$

where $U(N, \beta) = \frac{N\beta\sqrt{\varepsilon_{\beta,r}(x)}}{1 + (N - 1) \exp(-\beta)} \left(\frac{N \exp(2\beta)}{(N - 1)^2} + 1 \right)$.

Using similar approach to Proposition 5.3, we can derive the lower bound on regret of the BoN policy through the lens of smoothing.

Proposition 5.7 (Lower Bound on Regret of BoN). *Under the same assumptions in Theorem 5.6, the following lower bound holds on regret of BoN policy,*

$$\Delta_{J,r^*}(\pi_{r^*}^*(\cdot|x), \pi_{\hat{r}}^{\text{BoN}}(\cdot|x)) \geq \gamma(x)(1 - C_{\infty,r^*}(x)^{-1})^N - \min(1, N\sqrt{\varepsilon_{\infty,r}(x)}). \quad (19)$$

Under no overoptimization ($\varepsilon_{\infty,r}(x) = 0$), due to $(1 - C_{\infty,r^*}(x)^{-1})^N \geq \exp(\frac{-N}{C_{\infty,r^*}(x)-1})$, the lower bound on regret of BoN is $O(\exp(-N))$. Furthermore, in App. I.2, we present a numerical example to illustrate the positivity of the lower bounds under certain conditions.

5.3 DISCUSSION

This section analyzes the performance of BoN and SBoN in the presence and absence of overoptimization.

Remark 5.8 (Overoptimization). *Considering the upper bounds in equation 14 and equation 15, we define,*

$$g(\beta) = \beta(\sqrt{\varepsilon_{\infty,r}(x)} - \sqrt{\varepsilon_{\beta,r}(x)}), \quad \beta \geq 0.$$

Because $g(0) = g(\infty) = 0$ and $g(\beta) \geq 0$ for all β , there exists at least one maximizer $\beta^ \in (0, \infty)$ such that $g(\beta^*) = \max_{\beta \geq 0} g(\beta)$.*

If $\frac{\log C_{\infty,r^,\text{ref}}(x)}{\sqrt{C_{\infty,\hat{r},\text{ref}}(x)} + \sqrt{C_{\infty,r^*,\text{ref}}(x)}} \leq g(\beta^*)$, then the upper bound in equation 14 does not exceed equation 15, and hence the bound on the regret of the SBoN policy is tighter than the bound on the regret of the BoN policy under the proxy-reward model. An analogous comparison can be carried out for any fixed β and changing N . A numerical example for the existence of optimal β based on this approach is provided in App. I.2.*

Remark 5.9 (No overoptimization). *By Lemma 5.1, $C_{\infty,r^*,\text{ref}}(x) \geq 1$; consequently, the bound in equation 17 is tighter than the bound in equation 16. Similar discussion can be provided for lower bound.*

Remark 5.10 (Quality of reference policy). *The upper bounds in Proposition 5.3 (or Theorem 5.2) depend on two quantities,*

$$C_{\infty, r^*}(x) = \frac{1}{\sum_i \pi_{\text{ref}}(y_{i, r^*}^{\max}(x) | x)}, \quad \text{and} \quad C_{\infty, \hat{r}}(x) = \frac{1}{\sum_i \pi_{\text{ref}}(y_{i, \hat{r}}^{\max}(x) | x)},$$

which represent its quality under the true-reward model and proxy-reward model, respectively. A larger value for these quantities implies that the reference model rarely generates optimal responses, thereby degrading performance. We can observe that upper bounds increase by increasing $C_{\infty, r^}(x)$ and $C_{\infty, \hat{r}}(x)$.*

6 EMPIRICAL EVIDENCE

To support our theoretical analysis, we conducted experiments comparing Soft Best-of-N (SBoN) across different regularization strengths and reward model qualities. We used the Olmo-2 1B model (OLMo et al., 2024) as the generator and prompts from the Attaq dataset (Kour et al., 2023). For each prompt, we generated multiple responses and selected one using SBoN with varying temperature values β . We ran two experimental conditions: one using a strong proxy-reward model (ArmoRM 8B (Wang et al., 2024)) which is close to true-reward model, and another using a weaker proxy-reward model (Beaver 7B RM (Dai et al., 2023)). We use LLM-as-a-Judge Zheng et al. (2023) as our r^* . **To match our theoretical setting, we perform empirical calibration of each reward model by sampling 256 responses for every query and calculate the quantiles.** As shown in Figure 2, when the reward model is weak, performance degrades for large N due to reward hacking. However, the smoothing in SBoN helps mitigate this degradation. This observation is also aligned with our theoretical analysis and discussion in Section 4, where under overoptimization there exists a β for a given N which outperforms BoN. For more details, see App. I. We also studied the behavior of our upper bound on the KL divergence between the SBoN policy and the reference policy, Lemma 4.2, at App. I.2. More experiments with a medium weak reward model are provided in App. I.1.

7 CONCLUSION AND FUTURE WORK

In this work, we establish a theoretical foundation for alignment strategies based on Soft Best-of-N (SBoN) and Best-of-N (BoN) policies. Specifically, we derive upper bounds on the KL divergence between the aligned policy such as SBoN or BoN and the reference policy. We also studied the regret gap between the optimal policy and the aligned policy, e.g., BoN and SBoN policies. We further analyze how errors in reward estimation affect performance in both KL divergence and regret gap. Notably, both our theoretical analysis and empirical results demonstrate that, under a proxy-reward model where overoptimization happens, SBoN perform better than BoN under some conditions.

Existing literature Beirami et al. (2024); Mroueh (2024), suggests a tighter upper bound on the KL divergence for BoN policy. Our derived result, however, provides a looser bound for BoN. Consequently, an interesting direction for future research is the derivation of an upper bound that is asymptotically tight for the BoN policy. Furthermore, while our current analysis utilizes a shared β for both the tilted error and the SBoN policy, exploring decoupled temperature parameters for these components remains a promising avenue for future study.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Afra Amini, Tim Vieira, Elliott Ash, and Ryan Cotterell. Variational best-of-n alignment. In *International Conference on Learning Representations (ICLR)*, 2025.
- Gholamali Aminian, Amir R Asadi, Idan Shenfeld, and Youssef Mroueh. Theoretical analysis of kl-regularized rlhf with multiple reference models. *arXiv preprint arXiv:2502.01203*, 2025.
- Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- Ananth Balashankar, Ziteng Sun, Jonathan Berant, Jacob Eisenstein, Michael Collins, Adrian Hutter, Jong Lee, Chirag Nagpal, Flavien Prost, Aradhana Sinha, Ananda Theertha Suresh, and Ahmad Beirami. Infalign: Inference-aware language model alignment. *International Conference on Machine Learning (ICML)*, 2025.
- 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. *International Conference on Machine Learning (ICML)*, 2024.
- Vivek S Borkar. Q-learning for risk-sensitive control. *Mathematics of Operations Research*, 2002.
- 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. *arXiv preprint arXiv:2407.21787*, 2024.
- Clément L Canonne. A short note on an inequality between kl and tv. *arXiv preprint arXiv:2202.07198*, 2022.
- Pierre Cardaliaguet, François Delarue, Jean-Michel Lasry, and Pierre-Louis Lions. *The master equation and the convergence problem in mean field games:(ams-201)*. Princeton University Press, 2019.
- Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, et al. Open problems and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint arXiv:2307.15217*, 2023.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- Thomas Coste, Usman Anwar, Robert Kirk, and David Krueger. Reward model ensembles help mitigate overoptimization. In *The Twelfth International Conference on Learning Representations*.
- Imre Csiszár and Frantisek Matus. Information projections revisited. *IEEE Transactions on Information Theory*, 49(6):1474–1490, 2003.
- Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. Safe rlhf: Safe reinforcement learning from human feedback. *arXiv preprint arXiv:2310.12773*, 2023.
- Jacob Eisenstein, Chirag Nagpal, Alekh Agarwal, Ahmad Beirami, Alexander Nicholas D’Amour, Krishnamurthy Dj Dvijotham, Adam Fisch, Katherine A Heller, Stephen Robert Pfohl, Deepak Ramachandran, et al. Helping or herding? reward model ensembles mitigate but do not eliminate reward hacking. In *First Conference on Language Modeling*.
- L Gao, J Tow, B Abbasi, S Biderman, S Black, A DiPofi, C Foster, L Golding, J Hsu, A Le Noach, et al. A framework for few-shot language model evaluation. URL <https://zenodo.org/records/10256836>, 7, 2023a.

- Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In *International Conference on Machine Learning*, pp. 10835–10866. PMLR, 2023b.
- Jonathan Geuter, Youssef Mroueh, and David Alvarez-Melis. Guided speculative inference for efficient test-time alignment of llms. *arXiv preprint arXiv:2506.04118*, 2025.
- Lin Gui, Cristina Gârbacea, and Victor Veitch. Bonbon alignment for large language models and the sweetness of best-of-n sampling. *arXiv preprint arXiv:2406.00832*, 2024.
- Emil Julius Gumbel. Statistical theory of extreme values and some practical applications. *Nat. Bur. Standards Appl. Math. Ser. 33*, 1954.
- J. Hilton, P. Clark, et al. Measuring goodharts law: Towards an evaluation framework for open-ended generative models. <https://openai.com/index/measuring-goodharts-law>, 2022. Accessed: 2025-01-30.
- Ronald A Howard and James E Matheson. Risk-sensitive markov decision processes. *Management science*, 1972.
- Audrey Huang, Wenhao Zhan, Tengyang Xie, Jason D Lee, Wen Sun, Akshay Krishnamurthy, and Dylan J Foster. Correcting the mythos of kl-regularization: Direct alignment without overoptimization via chi-squared preference optimization. *arXiv preprint arXiv:2407.13399*, 2024.
- Audrey Huang, Adam Block, Qinghua Liu, Nan Jiang, Dylan J Foster, and Akshay Krishnamurthy. Is best-of-n the best of them? coverage, scaling, and optimality in inference-time alignment. *arXiv preprint arXiv:2503.21878*, 2025.
- John Hughes, Sara Price, Aengus Lynch, Rylan Schaeffer, Fazl Barez, Sanmi Koyejo, Henry Sleight, Erik Jones, Ethan Perez, and Mrinank Sharma. Best-of-n jailbreaking. *arXiv preprint arXiv:2412.03556*, 2024.
- Yuki Ichihara, Yuu Jinnai, Tetsuro Morimura, Kenshi Abe, Kaito Ariu, Mitsuki Sakamoto, and Eiji Uchibe. Evaluation of best-of-n sampling strategies for language model alignment. *Transactions on Machine Learning Research*.
- Hamish Ivison, Yizhong Wang, Jiacheng Liu, Zeqiu Wu, Valentina Pyatkin, Nathan Lambert, Noah A Smith, Yejin Choi, and Hanna Hajishirzi. Unpacking dpo and ppo: Disentangling best practices for learning from preference feedback. *Advances in neural information processing systems*, 37:36602–36633, 2024.
- Yuu Jinnai, Tetsuro Morimura, Kaito Ariu, and Kenshi Abe. Regularized best-of-n sampling to mitigate reward hacking for language model alignment. In *ICML 2024 Workshop on Models of Human Feedback for AI Alignment*, 2024.
- Hadi Khalaf, Claudio Mayrink Verdun, Alex Oesterling, Himabindu Lakkaraju, and Flavio du Pin Calmon. Inference-time reward hacking in large language models. *arXiv preprint arXiv:2506.19248*, 2025.
- Maxim Khanov, Jirayu Burapachee, and Yixuan Li. Args: Alignment as reward-guided search. *arXiv preprint arXiv:2402.01694*, 2024.
- Tomasz Korbak, Hady Elsahar, Germán Kruszewski, and Marc Dymetman. On reinforcement learning and distribution matching for fine-tuning language models with no catastrophic forgetting. *Advances in Neural Information Processing Systems*, 35:16203–16220, 2022a.
- Tomasz Korbak, Ethan Perez, and Christopher L Buckley. RL with kl penalties is better viewed as bayesian inference. *arXiv preprint arXiv:2205.11275*, 2022b.
- Barry W Kort and Dimitri P Bertsekas. A new penalty function method for constrained minimization. In *IEEE Conference on Decision and Control and 11th Symposium on Adaptive Processes*, 1972.
- George Kour, Marcel Zalmanovici, Naama Zwerdling, Esther Goldbraich, Ora Nova Fandina, Ateret Anaby-Tavor, Orna Raz, and Eitan Farchi. Unveiling safety vulnerabilities of large language models. *arXiv preprint arXiv:2311.04124*, 2023.

- Tian Li, Ahmad Beirami, Maziar Sanjabi, and Virginia Smith. On tilted losses in machine learning: Theory and applications. *Journal of Machine Learning Research*, 24(142):1–79, 2023.
- Guan-Hong Liu and Evangelos A Theodorou. Deep learning theory review: An optimal control and dynamical systems perspective. *arXiv preprint arXiv:1908.10920*, 2019.
- Claudio Mayrink Verdun, Alex Oesterling, Himabindu Lakkaraju, and Flavio P Calmon. Soft best-of-n sampling for model alignment. *arXiv preprint arXiv:2505.03156*, 2025.
- Youssef Mroueh. Information theoretic guarantees for policy alignment in large language models. *arXiv preprint arXiv:2406.05883*, 2024.
- Sidharth Mudgal, Jong Lee, Harish Ganapathy, Yaguang Li, Tao Wang, Yanping Huang, Zhifeng Chen, Heng-Tze Cheng, Michael Collins, Trevor Strohman, et al. Controlled decoding from language models. In *International Conference on Machine Learning*, pp. 36486–36503. PMLR, 2024.
- Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, et al. 2 olmo 2 furious. *arXiv preprint arXiv:2501.00656*, 2024.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- EY Pee and Johannes O Royset. On solving large-scale finite minimax problems using exponential smoothing. *Journal of Optimization Theory and Applications*, 2011.
- Yury Polyanskiy and Yihong Wu. Information theory: From coding to learning, 2022.
- Jiahao Qiu, Yifu Lu, Yifan Zeng, Jiacheng Guo, Jiayi Geng, Huazheng Wang, Kaixuan Huang, Yue Wu, and Mengdi Wang. Treebon: Enhancing inference-time alignment with speculative tree-search and best-of-n sampling. *arXiv preprint arXiv:2410.16033*, 2024.
- 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. *arXiv preprint arXiv:2305.18290*, 2023.
- Pier Giuseppe Sessa, Robert Dadashi, Léonard Hussenot, Johan Ferret, Nino Vieillard, Alexandre Ramé, Bobak Shariari, Sarah Perrin, Abe Friesen, Geoffrey Cideron, et al. Bond: Aligning llms with best-of-n distillation. *arXiv preprint arXiv:2407.14622*, 2024.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- Yuda Song, Gokul Swamy, Aarti Singh, Drew Bagnell, and Wen Sun. The importance of online data: Understanding preference fine-tuning via coverage. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- 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. *Advances in Neural Information Processing Systems*, 33:3008–3021, 2020.
- Benedikt Stroebel, Sayash Kapoor, and Arvind Narayanan. Inference scaling f-laws: The limits of llm resampling with imperfect verifiers. *arXiv preprint arXiv:2411.17501*, 2024.
- Hanshi Sun, Momin Haider, Ruiqi Zhang, Huitao Yang, Jiahao Qiu, Ming Yin, Mengdi Wang, Peter Bartlett, and Andrea Zanette. Fast best-of-n decoding via speculative rejection. *arXiv preprint arXiv:2410.20290*, 2024.

- Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Haoxiang Wang, Wei Xiong, Tengyang Xie, Han Zhao, and Tong Zhang. Interpretable preferences via multi-objective reward modeling and mixture-of-experts. In *EMNLP*, 2024.
- Xueqin Wang, Yunlu Jiang, Mian Huang, and Heping Zhang. Robust variable selection with exponential squared loss. *Journal of the American Statistical Association*, 2013.
- Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang. Iterative preference learning from human feedback: Bridging theory and practice for rlhf under kl-constraint. In *Forty-first International Conference on Machine Learning*, 2024.
- Joy Qiping Yang, Salman Salamatian, Ziteng Sun, Ananda Theertha Suresh, and Ahmad Beirami. Asymptotics of language model alignment. In *2024 IEEE International Symposium on Information Theory (ISIT)*, pp. 2027–2032. IEEE, 2024.
- Kaylee Yingxi Yang and Andre Wibisono. Convergence of the inexact langevin algorithm and score-based generative models in kl divergence. *arXiv preprint arXiv:2211.01512*, 2022.
- Chenlu Ye, Wei Xiong, Yuheng Zhang, Nan Jiang, and Tong Zhang. A theoretical analysis of nash learning from human feedback under general kl-regularized preference. *arXiv preprint arXiv:2402.07314*, 2024.
- Wenhao Zhan, Masatoshi Uehara, Nathan Kallus, Jason D Lee, and Wen Sun. Provable offline preference-based reinforcement learning. *arXiv preprint arXiv:2305.14816*, 2023.
- Heyang Zhao, Chenlu Ye, Quanquan Gu, and Tong Zhang. Sharp analysis for kl-regularized contextual bandits and rlhf. *arXiv preprint arXiv:2411.04625*, 2024.
- Yao Zhao, Mikhail Khalman, Rishabh Joshi, Shashi Narayan, Mohammad Saleh, and Peter J Liu. Calibrating sequence likelihood improves conditional language generation. In *The Eleventh International Conference on Learning Representations*, 2022.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.

Appendix

Table of Contents

A	Table of Notations	16
B	Other Related Works	16
C	Calibrated Reward	17
D	Summary of KL divergence Results	17
E	Gumbel-Max trick	17
F	Technical Tools	18
G	Proof and Details of Section 4	21
H	Proof and details of Section 5	23
	H.1 Upper bound proof and details	23
	H.2 Lower bound proofs and details:	27
I	Experiments	29
	I.1 More Experiments	29
	I.2 Numerical Examples	30
J	Comparison with Huang et al. (2025)	33
K	Information projection across reward functions	34

A TABLE OF NOTATIONS

All notations are summarized in Table 1.

Notation	Definition
$r^*(y, x)$	(Calibrated) true-reward model
$\hat{r}(y, x)$	(Calibrated) proxy-reward model
N	Number of responses
x	Prompt
β	Temperature (regularization coefficient)
$\pi_{\beta, r}(y x)$	Tilted policy under reward function r
$\pi_r^{\text{SBoN}}(y x)$	SBoN policy under reward function r
$\pi_r^{\text{BoN}}(y x)$	BoN policy under reward function r
π_{ref}	Reference policy
$\text{KL}(p q)$	KL divergence between p and q distributions
$\pi_r^*(y x)$	Optimal policy
$C_{\beta, r}(x)$	Coverage constant under reward model r
$\mathcal{Y}_{r^*}^*(x)$	Set of maximizers for $r^*(y, x)$
$\varepsilon_{\beta, r}(x)$	Tilted Error

Table 1: Summary of notations in the paper

B OTHER RELATED WORKS

Smoothing of Maximum: Approximating the maximum operator using a smoothed or softmax-based surrogate is a widely adopted technique in machine learning. This approach is particularly useful in settings where the hard maximum is non-differentiable or leads to unstable optimization. For instance, in robust regression, smooth approximations to the max operator are used in min-max formulations to achieve tractable optimization under distributional shifts (Wang et al., 2013; Li et al., 2023). In sequential decision-making, similar ideas appear in risk-sensitive control and Q-learning, where the softmax of Q-values leads to stochastic policies that balance exploration and exploitation (Howard & Matheson, 1972; Borkar, 2002). In convex and non-convex optimization, smoothing the maximum objective has been shown to improve convergence properties (Kort & Bertsekas, 1972; Pee & Royset, 2011; Liu & Theodorou, 2019). The Soft Best-of-N (SBoN) framework, (Mayrink Verdun et al., 2025; Khanov et al., 2024; Jinnai et al., 2024), leverages this principle by replacing the hard selection of the highest-reward sample with a softmax-weighted sampling distribution. Regarding the SBoN, the empirical version of SBoN is introduced by (Khanov et al., 2024) as ARGS-stochastic, where a token from a probability distribution among the top-k candidate tokens is chosen. Then, the regularized version of BoN, which can be represented as SBoN, is discussed by (Jinnai et al., 2024). Given the broad success of SBoN, we are motivated to theoretically investigate the SBoN policies and the effect of the proxy-reward model (reward hacking) and the quality of the reference policy.

Theoretical Foundation of RLHF: Several works have studied the theoretical underpinnings of reverse KL-regularized RLHF, particularly in terms of sample complexity (Zhao et al., 2024; Xiong et al., 2024; Song et al., 2024; Zhan et al., 2023; Ye et al., 2024; Aminian et al., 2025). Note that, as the sampling distributions in BoN and SBoN are different, we can not apply RLHF analysis to these sampling strategies. Therefore, it is needed to develop new foundations for BoN and SBoN.

Overoptimization. Alignment methods are widely known to suffer from overoptimization, also known as misspecification, reward hacking, or Goodhart Law, where optimizing against a proxy-reward model leads to degraded performance compared to the true-reward model (Amodei et al., 2016; Casper et al., 2023; Gao et al., 2023b). This issue is particularly pronounced in inference-time alignment methods such as BoN, where an increasing number of responses N makes the overoptimization problem worse (Huang et al., 2025; Stroebl et al., 2024; Gao et al., 2023b). Huang et al. (2025) theoretically demonstrate that the BoN policy suffers from overoptimization when N is large, given a fixed estimation error in the reward model, and propose a solution based on a χ^2 -regularized

framework. Other approaches to mitigating this issue include ensembling strategies (Coste et al.; Eisenstein et al.) and regularization techniques (Ichihara et al.). In a concurrent line of work, Khalaf et al. (2025) introduce the Best-of-Poisson method to reduce overoptimization in inference-time algorithms. The overoptimization in BoN and SBoN is also studied by Khalaf et al. (2025) and a principled hedging framework is proposed to mitigate the overoptimization. In contrast, we study overoptimization in inference-time alignment methods SBoN and BoN from the perspectives of regret gap and KL divergence analysis.

C CALIBRATED REWARD

Inspired by (Balashankar et al., 2025), in this section, we provide more details regarding calibrated reward. A standard metric for evaluation of models is the *win-rate* relative to a base policy π_{ref} (Stiennon et al., 2020; Gao et al., 2023b). For a prompt x and responses y, z , define the win random variable under raw (uncalibrated) reward r_{uc} as

$$w_r(y, z | x) = \mathbf{1}\{r_{\text{uc}}(y, x) > r_{\text{uc}}(z, x)\} + \frac{1}{2} \mathbf{1}\{r_{\text{uc}}(y, x) = r_{\text{uc}}(z, x)\}.$$

Definition C.1 (Calibrated reward). *The calibrated reward of y under policy π is its expected win-rate probability against $z \sim \pi(\cdot | x)$:*

$$r_{c,\pi}(y, x) := \mathbb{E}_{z \sim \pi(\cdot | x)}[w_r(y, z | x)].$$

In practice, we consider $\pi = \pi_{\text{ref}}$, therefore we denote calibrated reward via $r(y, x)$ under reference policy. In the following, we provide some reasons for choosing calibrated reward instead of raw (uncalibrated) reward in our work,

- **Matches win-rate evaluation.** For any policies π_1, π_2 ,

$$W_r(\pi_1 \succ \pi_2 | x) := \mathbb{E}_{y \sim \pi_1(\cdot | x)}[r_{c,\pi_2}(y, x)],$$

where $W_r(\pi_1 \succ \pi_2 | x)$ is standard win-rate. So maximizing $\mathbb{E}_{Y \sim \pi}[r_{c,\pi_{\text{ref}}}(y, x)]$ directly optimizes standard win rate vs. the base model.

- **Invariance to score scaling.** If m is strictly increasing and $r' = m \circ r$, then

$$r'(y, x) = r(y, x),$$

making the target robust to arbitrary monotone reparameterizations of the reward (e.g., affine rescaling, temperature).

- **Unified, probabilistic scale.** For $y \sim \pi_{\text{ref}}(\cdot | x)$,

$$r(y, x) \sim \text{Unif}[0, 1],$$

independent of both r and π_{ref} . This normalizes per-prompt reward scales and interprets scores as win probabilities.

D SUMMARY OF KL DIVERGENCE RESULTS

In Table 2, we summarize results on KL divergences between the aligned and reference policies, along with corresponding upper bounds for both SBoN and BoN policies. Furthermore, in Table 3, we summarize results on KL divergences between aligned policies under true and proxy-reward models, along with upper bounds for SBoN and tilted policies.

E GUMBEL-MAX TRICK

We also provide an interpretation for SBoN from the Gumbel-Max trick. An alternative way to sample Z from

$$\Pr(Z = i) \propto \exp(\beta \hat{r}(Y_i, x))$$

is via the Gumbel-Max trick. We can draw independent Gumbel-distributed random variables $G_i \sim \text{Gumbel}(0, 1)$, $i = 1, \dots, n$, and then set

$$Z = \arg \max_{i \in \{1, \dots, N\}} \left[\hat{r}(Y_i, x) + \frac{G_i}{\beta} \right].$$

Table 2: KL divergences between the aligned and reference policies, along with corresponding upper bounds for both SBoN and BoN policies.

KL divergence Term	Theorem / Lemma	Upper Bound
$\text{KL}(\pi_{r^*}^{\text{SBoN}}(\cdot x) \parallel \pi_{\text{ref}}(\cdot x))$	Lemma 4.1	$\log\left(\frac{N}{1+(N-1)\exp(-\beta)}\right)$
$\text{KL}(\pi_{r^*}^{\text{BoN}}(\cdot x) \parallel \pi_{\text{ref}}(\cdot x))$	Theorem 3.1 in (Beirami et al., 2024) and Theorem 1 in (Mroueh, 2024)	$\log(N) - 1 + 1/N$

Table 3: KL divergences between aligned policies under true and proxy-reward models, along with upper bounds for SBoN and tilted policies.

KL divergence Term	Theorem / Lemma	Upper Bound
$\text{KL}(\pi_{r^*}^{\text{SBoN}}(\cdot x) \parallel \pi_{\hat{r}}^{\text{SBoN}}(\cdot x))$	Lemma 4.2	$\frac{N\beta\sqrt{\varepsilon_{\beta,r}(x)}}{1+(N-1)\exp(-\beta)} \left(\frac{N\exp(2\beta)}{(N-1)^2} + 1 \right)$
$\text{KL}(\pi_{\beta,r^*}(\cdot x) \parallel \pi_{\beta,\hat{r}}(\cdot x))$	Lemma F.8	$2\beta\sqrt{\varepsilon_{\beta,r}(x)} \left(\sqrt{\frac{\mathbb{E}[\exp(2\beta\hat{r}(Y,x))]}{\mathbb{E}^2[\exp(\beta\hat{r}(Y,x))]} + \sqrt{\frac{\mathbb{E}[\exp(2\beta r^*(Y,x))]}{\mathbb{E}^2[\exp(\beta r^*(Y,x))]} \right)$

By properties of the Gumbel distribution, this yields exactly the same softmax sampling law, without needing to compute the normalizing factor $\sum_{j=1}^N \exp(\beta\hat{r}(Y_j, x))$ explicitly (Gumbel, 1954). When $\beta \rightarrow \infty$, the effect of the Gumbel noises vanishes and the sampling strategy reduces to BoN.

F TECHNICAL TOOLS

We denote the set maximizers of the proxy-reward function via $\hat{\mathcal{Y}}(x) = \{\hat{y}_j(x)\}_{j=1}^{m(x)}$.

We introduce the functional derivative, see Cardaliaguet et al. (2019).

Definition F.1. (Cardaliaguet et al., 2019) *A functional $U : \mathcal{P}(\mathbb{R}^n) \rightarrow \mathbb{R}$ admits a functional derivative if there is a map $\frac{\delta U}{\delta m} : \mathcal{P}(\mathbb{R}^n) \times \mathbb{R}^n \rightarrow \mathbb{R}$ which is continuous on $\mathcal{P}(\mathbb{R}^n)$ and, for all $m, m' \in \mathcal{P}(\mathbb{R}^n)$, it holds that*

$$U(m') - U(m) = \int_0^1 \int_{\mathbb{R}^n} \frac{\delta U}{\delta m}(m_\beta, a) (m' - m)(da) d\beta,$$

where $m_\beta = m + \beta(m' - m)$.

Definition F.2 (Sensitivity of a policy). *We also define the sensitivity of a policy $\pi_r(y|x)$, which is a function of reward function $r(y, x)$, with respect to the reward function as*

$$\frac{\partial \pi}{\partial r}(r) := \lim_{\Delta r \rightarrow 0} \frac{\pi_r(y|x) - \pi_{r+\Delta r}(y|x)}{\Delta r}. \quad (20)$$

Lemma F.3 (Kantorovich-Rubenstein duality of total variation distance, see (Polyanskiy & Wu, 2022)). *The Kantorovich-Rubenstein duality (variational representation) of the total variation distance is as follows:*

$$\text{TV}(m_1, m_2) = \frac{1}{L} \sup_{g \in \mathcal{G}_L} \{ \mathbb{E}_{Z \sim m_1}[g(Z)] - \mathbb{E}_{Z \sim m_2}[g(Z)] \}, \quad (21)$$

where $\mathcal{G}_L = \{g : \mathcal{Z} \rightarrow \mathbb{R}, \|g\|_\infty \leq L\}$.

Lemma F.4 (Lemma 5.4 in (Aminian et al., 2025)). *Consider the softmax policy, $\pi_r^\beta(y|x) \propto \pi_{\text{ref}}(y|x) \exp(\beta r(y, x))$. Then, the sensitivity of the policy with respect to the reward function is*

$$\frac{\partial \pi_r^\beta}{\partial r}(r) = \beta \pi_r^\beta(y|x) (1 - \pi_r^\beta(y|x)).$$

Lemma F.5 (Pinskens Inequality (Canonne, 2022)). *For m_1 and m_2 , we have,*

$$\text{TV}(m_1, m_2) \leq \sqrt{\frac{1}{2} \text{KL}(m_2 \| m_1)}. \quad (22)$$

The following Lemmata are useful for our technical proofs.

Lemma F.6. *The following upper bound holds,*

$$\log \left(\frac{Z_{r^*,Y}(x, \beta)}{Z_{\hat{r},Y}(x, \beta)} \right) \leq \beta \sqrt{\varepsilon_{\beta,r}(x)} \sqrt{C_{\beta,r^*}(x)}. \quad (23)$$

Proof.

$$\begin{aligned} \frac{Z_{\hat{r},Y}(x, \beta)}{Z_{r^*,Y}(x, \beta)} &= \frac{\sum_y \exp(\beta \hat{r}(y, x)) \pi_{\text{ref}}(y|x)}{\sum_y \exp(\beta r^*(y, x)) \pi_{\text{ref}}(y|x)} \\ &= \frac{\sum_y \exp(\beta(\hat{r}(y, x) - r^*(y, x))) \exp(\beta r^*(y, x)) \pi_{\text{ref}}(y|x)}{\sum_y \exp(\beta r^*(y, x)) \pi_{\text{ref}}(y|x)} \\ &= \sum_y \pi_{\beta,r^*}(y|x) \exp(\beta(\hat{r}(y, x) - r^*(y, x))) \end{aligned} \quad (24)$$

Due to convexity of $-\log(\cdot)$ and using Cauchy-Schwarz inequality, we have,

$$\begin{aligned} -\log \left(\frac{Z_{\hat{r},Y}(x, \beta)}{Z_{r^*,Y}(x, \beta)} \right) &\leq \beta \sum_y \pi_{\beta,r^*}(y|x) (r^*(y, x) - \hat{r}(y, x)) \\ &\leq \beta \sum_y \frac{\pi_{\beta,r^*}(y|x)}{\pi_{\text{ref}}(y|x)} (r^*(y, x) - \hat{r}(y, x)) \pi_{\text{ref}}(y|x) \\ &\leq \beta \sqrt{\sum_y (r^*(y, x) - \hat{r}(y, x))^2 \pi_{\text{ref}}(y|x)} \sqrt{\sum_y \frac{\pi_{\beta,r^*}^2(y|x)}{\pi_{\text{ref}}(y|x)}} \\ &= \sqrt{\beta} \sqrt{\sum_y \log(\exp(\beta(r^*(y, x) - \hat{r}(y, x))^2) \pi_{\text{ref}}(y|x))} \sqrt{C_{\beta,r^*}(x)} \\ &\leq \beta \sqrt{\frac{1}{\beta} \log \left(\sum_y \exp(\beta(r^*(y, x) - \hat{r}(y, x))^2) \pi_{\text{ref}}(y|x) \right)} \sqrt{C_{\beta,r^*}(x)} \\ &= \beta \sqrt{\varepsilon_{\beta,r}(x)} \sqrt{C_{\beta,r^*}(x)}, \end{aligned} \quad (25)$$

□

Lemma F.7. *The following holds,*

$$\text{KL}(\pi_{r^*}^*(\cdot | x) \| \pi_{\text{ref}}(\cdot | x)) \leq \log(C_{\infty,r^*}(x)). \quad (26)$$

Proof. Note that, we have,

$$\begin{aligned} \text{KL}(\pi_{r^*}^*(\cdot | x) \| \pi_{\text{ref}}(\cdot | x)) &\leq \log \left(\mathbb{E}_{Y \sim \pi_{r^*}^*} \left(\frac{\pi_{r^*}^*(\cdot | x)}{\pi_{\text{ref}}(\cdot | x)} \right) \right) \\ &\leq \log(C_{\infty,r^*}(x)). \end{aligned} \quad (27)$$

□

Lemma F.8. *The following upper bound holds,*

$$\text{KL}(\pi_{\beta,r^*}(y|x) \| \pi_{\beta,\hat{r}}(y|x)) \leq \beta \sqrt{\varepsilon_{\beta,r}(x)} (\sqrt{C_{\beta,r^*}(x)} + \sqrt{C_{\beta,\hat{r}}(x)}). \quad (28)$$

Proof.

$$\begin{aligned} \text{KL}(\pi_{\beta,r^*}(y|x) \| \pi_{\beta,\hat{r}}(y|x)) &= \sum_y \pi_{\beta,r^*}(y|x) \log \left(\frac{\pi_{\beta,r^*}(y|x)}{\pi_{\beta,\hat{r}}(y|x)} \right) \\ &= \beta \sum_y (r^*(y, x) - \hat{r}(y, x)) \pi_{\beta,r^*}(y|x) + \log(Z_{\hat{r},Y}(x, \beta) / Z_{r^*,Y}(x, \beta)) \\ &\leq \beta \sqrt{\varepsilon_{\beta,r}(x)} (\sqrt{C_{\beta,r^*}(x)} + \sqrt{C_{\beta,\hat{r}}(x)}), \end{aligned} \quad (29)$$

where the final inequality holds due to Lemma F.6 and applying Cauchy-Schwarz inequality. \square

Lemma F.9. Suppose that $f(Z) \in [0, Z_{\max}]$, $\mathcal{Z}_{\max} = \{z_{m,i}\}_{i=1}^m$ is the set of maximizers of $f(Z)$, i.e., $f(z) = Z_{\max}$ for $z \in \mathcal{Z}_{\max}$. Then we have,

$$\lim_{\beta \rightarrow \infty} \frac{\mathbb{E}[\exp(2\beta f(Z))]}{\mathbb{E}[\exp(\beta f(Z))]^2} = \frac{1}{\sum_{z \in \mathcal{Z}_{\max}} P(Z = z)}. \quad (30)$$

Proof.

$$\frac{\mathbb{E}[\exp(2\beta f(Z))]}{\mathbb{E}[\exp(\beta f(Z))]^2} = \frac{\mathbb{E}[\exp(2\beta(f(Z) - Z_{\max}))]}{\mathbb{E}[\exp(\beta(f(Z) - Z_{\max}))]^2} \quad (31)$$

$$\frac{\sum_j P(Z = z_j) \exp(2\beta(f(z_j) - Z_{\max}))}{(\sum_j P(Z = z_j) \exp(\beta(f(z_j) - Z_{\max})))^2} \quad (32)$$

Now, we have,

$$\lim_{\beta \rightarrow \infty} \frac{\mathbb{E}[\exp(2\beta f(Z))]}{\mathbb{E}[\exp(\beta f(Z))]^2} \quad (33)$$

$$= \lim_{\beta \rightarrow \infty} \frac{\sum_j P(Z = z_j) \exp(2\beta(f(z_j) - Z_{\max}))}{(\sum_j P(Z = z_j) \exp(\beta(f(z_j) - Z_{\max})))^2} \quad (34)$$

$$= \frac{\sum_{z \in \mathcal{Z}_{\max}} P(Z = z)}{(\sum_{z \in \mathcal{Z}_{\max}} P(Z = z))^2} \quad (35)$$

$$= \frac{1}{\sum_{z \in \mathcal{Z}_{\max}} P(Z = z)}, \quad (36)$$

where we used the fact that $\lim_{\beta \rightarrow \infty} \exp(\beta(z_j - Z_{\max})) = 0$ for $z_j < Z_{\max}$. \square

Lemma F.10 (Theorem 1 in (Mayrink Verdun et al., 2025)). For $\beta > 0$, and $N \geq 1$, we have,

$$\text{KL}(\pi_{\beta, r^*}(\cdot|x) \parallel \pi_{r^*}^{\text{SBoN}}(y|x)) \leq \log \left(1 + \frac{C_{\beta, r^*}(x) - 1}{N} \right). \quad (37)$$

Lemma F.11. For a given $x \in \mathcal{X}$, we have,

$$\left| \frac{\delta f(r)}{\delta r} \right| \leq \frac{N^2 \beta \exp(2\beta)}{(N-1)^2}, \quad (38)$$

where $f(r) = \log \left(\mathbb{E} \left[\frac{1}{\exp(\beta r) + \sum_{i=1}^{N-1} \exp(\beta R_i)} \right] \right)$, $r = r(y, x)$ and $R_i = r(Y_i, x)$.

Proof. Note that $\{R_i\}_{i=1}^{N-1}$ are i.i.d. . Therefore, we have,

$$\begin{aligned} \frac{\delta f(r(y, x))}{\delta r} &= \mathbb{E} \left[\frac{1}{\exp(\beta r) + \sum_{i=1}^{N-1} \exp(\beta R_i)} \right]^{-1} \frac{\delta \mathbb{E} \left[\frac{1}{\exp(\beta r) + \sum_{i=1}^{N-1} \exp(\beta R_i)} \right]}{\delta r} \\ &\leq \mathbb{E} \left[\frac{1}{\exp(\beta r) + \sum_{i=1}^{N-1} \exp(\beta R_i)} \right]^{-1} \\ &\quad \times \left(\sum_{k=1}^N \frac{\beta k \binom{N-1}{k-1} \exp(\beta r)}{(k \exp(\beta r) + N - 1 - k)^2} (1 - P(R = r))^{N-k} P^{k-1}(R = r) \right) \\ &\leq \mathbb{E} \left[\frac{1}{\exp(\beta r) + \sum_{i=1}^{N-1} \exp(\beta R_i)} \right]^{-1} \\ &\quad \times \frac{\beta \exp(\beta)}{(N-1)^2} \left(\sum_{k=1}^N k \binom{N-1}{k-1} (1 - P(R = r))^{N-k} P^{k-1}(R = r) \right) \\ &\leq \frac{N \beta \exp(2\beta)}{(N-1)^2} (1 + (N-1)P(R = r)) \\ &\leq \frac{N^2 \beta \exp(2\beta)}{(N-1)^2}. \end{aligned} \quad (39)$$

For first inequality provide more details in the following. Lets consider the conditional expectation where we condition on having exactly $(k - 1)$ of $(N - 1)$ R_i s equal to r . Therefore, we have,

$$\begin{aligned} & \frac{\delta \mathbb{E} \left[\frac{1}{\exp(\beta r) + \sum_{i=1}^{N-1} \exp(\beta R_i)} \right]}{\delta r} \\ &= \sum_{k=1}^N \binom{N-1}{k-1} \mathbb{E} \left[\frac{\beta k \exp(\beta r)}{k \exp(\beta r) + \sum_{i=1}^{N-k} \exp(\beta R_i)} \right] (1 - P(R=r))^{N-k} P^{k-1}(R=r) \end{aligned}$$

Now computing the derivative and using the fact that $1 \leq \exp(\beta R_i)$, the first inequality holds. \square

Lemma F.12 (The BH inequality (Canonne, 2022)). *For every two probability distributions \mathbf{p}, \mathbf{q} , we have the simple yet never vacuous bound*

$$\text{TV}(\mathbf{p}, \mathbf{q}) \leq \sqrt{1 - e^{-\text{KL}(\mathbf{p} \parallel \mathbf{q})}}. \quad (40)$$

G PROOF AND DETAILS OF SECTION 4

Lemma 4.1. *The following upper bound holds on KL divergence between SBoN and reference policies for a given prompt $x \in \mathcal{X}$,*

$$\text{KL}(\pi_{r^*}^{\text{SBoN}}(y|x) \parallel \pi_{\text{ref}}(y|x)) \leq \log \left(\frac{N}{1 + (N-1) \exp(-\beta)} \right). \quad (41)$$

Proof. Recall that,

$$\pi_{r^*}^{\text{SBoN}}(y|x) = \pi_{\text{ref}}(y|x) \exp(\beta r^*(y, x)) \mathbb{E} \left[\left(\frac{1}{N} (\exp(\beta r^*(y, x)) + \sum_{i=1}^{N-1} \exp(\beta r^*(Y_i, x))) \right)^{-1} \right].$$

Now, we have,

$$\begin{aligned} & \text{KL}(\pi_{r^*}^{\text{SBoN}}(y|x) \parallel \pi_{\text{ref}}(y|x)) \\ &= \sum_{\mathcal{Y}} \pi_{r^*}^{\text{SBoN}}(y|x) \log(\pi_{r^*}^{\text{SBoN}}(y|x) / \pi_{\text{ref}}(y|x)) \\ &= \sum_{\mathcal{Y}} \pi_{r^*}^{\text{SBoN}}(y|x) \log(N) \\ & \quad + \sum_{\mathcal{Y}} \pi_{r^*}^{\text{SBoN}}(y|x) \log \left(\mathbb{E} \left[\exp(\beta r^*(y, x)) \left(\exp(\beta r^*(y, x)) + \sum_{i=1}^{N-1} \exp(\beta r^*(Y_i, x)) \right)^{-1} \right] \right) \\ &= \log(N) + \sum_{\mathcal{Y}} \pi_{r^*}^{\text{SBoN}}(y|x) \log \left(\mathbb{E} \left[\left(1 + \sum_{i=1}^{N-1} \exp(\beta(r^*(Y_i, x) - r^*(y, x))) \right)^{-1} \right] \right). \end{aligned} \quad (42)$$

For the second term in equation 42, consider

$$A(y, Y, x) = \sum_{i=1}^{N-1} \exp(\beta(r^*(Y_i, x) - r^*(y, x))) > 0,$$

where we have

$$(N-1) \exp(-\beta) \leq A(y, Y, x) \leq (N-1) \exp(\beta).$$

Therefore, we have,

$$\begin{aligned}
& \sum_y \pi_{r^*}^{\text{SBoN}}(y|x) \log(\mathbb{E}[(1 + \sum_{i=1}^{N-1} \exp(\beta(r^*(Y_i, x) - r^*(y, x))))^{-1}]) \\
& \leq \sum_y \pi_{r^*}^{\text{SBoN}}(y|x) \log\left(\frac{1}{1 + (N-1) \exp(-\beta)}\right) \\
& = \log\left(\frac{1}{1 + (N-1) \exp(-\beta)}\right).
\end{aligned} \tag{43}$$

Combining equation 43 with equation 42 completes the proof. \square

Lemma 4.2. *The following upper bound holds on the KL divergence between the SBoN policies under true-reward and proxy-reward models respectively,*

$$\text{KL}(\pi_{r^*}^{\text{SBoN}}(\cdot|x) \parallel \pi_{\hat{r}}^{\text{SBoN}}(\cdot|x)) \leq \frac{N\beta\sqrt{\varepsilon_{\beta,r}(x)}}{1 + (N-1) \exp(-\beta)} \left(\frac{N \exp(2\beta)}{(N-1)^2} + 1 \right). \tag{44}$$

Proof. We first provide the following upper bound,

$$\begin{aligned}
& \text{KL}(\pi_{r^*}^{\text{SBoN}}(y|x) \parallel \pi_{\hat{r}}^{\text{SBoN}}(y|x)) \\
& = \sum_y \pi_{r^*}^{\text{SBoN}}(y|x) \log\left(\frac{\pi_{r^*}^{\text{SBoN}}(y|x)}{\pi_{\hat{r}}^{\text{SBoN}}(y|x)}\right) \\
& = \sum_y \pi_{r^*}^{\text{SBoN}}(y|x) \beta(r^*(y, x) - \hat{r}(y, x)) \\
& \quad + \sum_y \pi_{r^*}^{\text{SBoN}}(y|x) \left(\log\left(\mathbb{E}\left[\frac{1}{\exp(\beta r^*(y, x)) + \sum_{i=1}^{N-1} \exp(\beta r^*(Y_i, x))}\right]\right) \right. \\
& \quad \left. - \log\left(\mathbb{E}\left[\frac{1}{\exp(\beta \hat{r}(y, x)) + \sum_{i=1}^{N-1} \exp(\beta \hat{r}(Y_i, x))}\right]\right) \right) \\
& \leq \frac{N\beta\sqrt{\varepsilon_{\beta,r}(x)}}{1 + (N-1) \exp(-\beta)} \\
& \quad + \sum_y \pi_{r^*}^{\text{SBoN}}(y|x) \left(\log\left(\mathbb{E}\left[\frac{1}{\exp(\beta r^*(y, x)) + \sum_{i=1}^{N-1} \exp(\beta r^*(Y_i, x))}\right]\right) \right. \\
& \quad \left. - \log\left(\mathbb{E}\left[\frac{1}{\exp(\beta \hat{r}(y, x)) + \sum_{i=1}^{N-1} \exp(\beta \hat{r}(Y_i, x))}\right]\right) \right),
\end{aligned} \tag{45}$$

where we used two facts. First,

$$\pi_{r^*}^{\text{SBoN}}(\cdot|x) \leq \frac{N\pi_{\text{ref}}(\cdot|x)}{1 + (N-1) \exp(-\beta)}$$

Second, for a random variable X ,

$$E[X] \leq \sqrt{E[X^2]} = \sqrt{E\left[\frac{1}{\beta} \log(\exp(\beta X^2))\right]} \leq \sqrt{\frac{1}{\beta} \log(E[\exp(\beta X^2)])},$$

and applying it to $E_{\pi_{\text{ref}}(\cdot|x)}[(r^*(y, x) - \hat{r}(y, x))]^2$, the final result in holds.

Note that for the last term in equation 45, we can apply the mean-value theorem as follows,

$$\begin{aligned} & \sum_{\mathcal{Y}} \pi_{r^*}^{\text{SBON}}(y|x) \left(\log \left(\mathbb{E} \left[\frac{1}{\exp(\beta r^*(y, x)) + \sum_{i=1}^{N-1} \exp(\beta r^*(Y_i, x))} \right] \right) \right. \\ & \quad \left. - \log \left(\mathbb{E} \left[\frac{1}{\exp(\beta \hat{r}(y, x)) + \sum_{i=1}^{N-1} \exp(\beta \hat{r}(Y_i, x))} \right] \right) \right) \\ & \leq \sum_{\mathcal{Y}} \pi_{r^*}^{\text{SBON}}(y|x) |r^*(y, x) - \hat{r}(y, x)| \left| \frac{\delta f(r_\gamma(y, x))}{\delta r} \right|, \end{aligned} \quad (46)$$

where $f(r_\gamma(y, x)) = \log \left(\mathbb{E} \left[\frac{1}{\exp(\beta r_\gamma(y, x)) + \sum_{i=1}^{N-1} \exp(\beta r_\gamma(Y_i, x))} \right] \right)$, for some $\gamma \in (0, 1)$ we have $r_\gamma(y, x) = \gamma \hat{r}(y, x) + (1 - \gamma)r^*(y, x)$. Using Lemma F.11, we have,

$$\left| \frac{\delta f(r_\gamma(y, x))}{\delta r} \right| \leq \frac{N^2 \beta \exp(2\beta)}{(N-1)^2}. \quad (47)$$

Using equation 47 in equation 46 and applying Cauchy-Schwarz inequality, we have,

$$\begin{aligned} & \sum_{\mathcal{Y}} \pi_{r^*}^{\text{SBON}}(y|x) \left(\log \left(\mathbb{E} \left[\frac{1}{\exp(\beta r^*(y, x)) + \sum_{i=1}^{N-1} \exp(\beta r^*(Y_i, x))} \right] \right) \right. \\ & \quad \left. - \log \left(\mathbb{E} \left[\frac{1}{\exp(\beta \hat{r}(y, x)) + \sum_{i=1}^{N-1} \exp(\beta \hat{r}(Y_i, x))} \right] \right) \right) \\ & \leq \sqrt{\sum_{\mathcal{Y}} \mathbb{E} \left[\frac{1}{(1 + \sum_{i=1}^{N-1} \exp(\beta(r^*(Y_i, x) - r^*(y, x))))} \right]^2 \pi_{\text{ref}}(y|x)} \\ & \quad \times \sqrt{\sum_{\mathcal{Y}} |r^*(y, x) - \hat{r}(y, x)|^2 \pi_{\text{ref}}(y|x)} \sqrt{\sum_{\mathcal{Y}} \left| \frac{\delta f(r_\gamma(y, x))}{\delta r} \right|^2 \pi_{\text{ref}}(y|x)} \\ & \leq \frac{\sqrt{\varepsilon_{\beta, r}(x)}}{1 + (N-1) \exp(-\beta)} \frac{N^2 \beta \exp(2\beta)}{(N-1)^2}. \end{aligned} \quad (48)$$

It completes the proof. \square

H PROOF AND DETAILS OF SECTION 5

H.1 UPPER BOUND PROOF AND DETAILS

Lemma 5.1 (Full Version). *Under Assumption 3.1, the following properties of $C_{\beta, r}(x)$ hold,*

1. $C_{\beta, r}(x) = \frac{\mathbb{E}[\exp(2\beta \hat{r}(Y, x))]}{\mathbb{E}^2[\exp(\beta \hat{r}(Y, x))]}.$
2. $C_{\beta, r}(x)$ is an increasing function with respect to β .
3. $C_{\infty, \hat{r}}(x) = \frac{1}{\sum_i \pi_{\text{ref}}(y_{i, r}^{\max}(x)|x)}$ where $y_{i, r}^{\max}(x) \in \arg \max_y r(y, x)$.
4. For all $\beta < \infty$, we have $1 \leq C_{\beta, r}(x) \leq \min(C_{\infty, \hat{r}}(x), \exp(2\beta)).$

Proof. In the following, we provide proofs of different items.

1.

$$\begin{aligned} C_{\beta, \hat{r}}(x) &= \sum_{\mathcal{Y}} \frac{\pi_{\beta, \hat{r}}^2(y|x)}{\pi_{\text{ref}}(y|x)} \\ &= \sum_{\mathcal{Y}} \frac{\exp(2\beta \hat{r}(y, x))}{\mathbb{E}^2[\exp(\beta \hat{r}(Y, x))]} \pi_{\text{ref}}(y|x) \\ &= \frac{\mathbb{E}[\exp(2\beta \hat{r}(Y, x))]}{\mathbb{E}^2[\exp(\beta \hat{r}(Y, x))]} \end{aligned} \quad (49)$$

2. We can show that the logarithm function of $C_{\beta, \hat{r}}(x)$ is increasing. Then, due to the increasing feature of the log function, the final result holds.

$$\begin{aligned} & \log\left(\frac{\mathbb{E}[\exp(2\beta\hat{r}(Y, x))]}{\mathbb{E}^2[\exp(\beta\hat{r}(Y, x))]} \right) \\ &= \log(\mathbb{E}[\exp(2\beta\hat{r}(Y, x))]) - 2\log(\mathbb{E}[\exp(\beta\hat{r}(Y, x))]), \end{aligned} \quad (50)$$

then we can compute the derivative of equation 50,

$$\begin{aligned} & \frac{d \log(\mathbb{E}[\exp(2\beta\hat{r}(Y, x))])}{d\beta} - 2 \frac{d \log(\mathbb{E}[\exp(\beta\hat{r}(Y, x))])}{d\beta} \\ &= \frac{\mathbb{E}[2\hat{r}(Y, x) \exp(2\beta\hat{r}(Y, x))]}{\mathbb{E}[\exp(2\beta\hat{r}(Y, x))]} - 2 \frac{\mathbb{E}[\hat{r}(Y, x) \exp(\beta\hat{r}(Y, x))]}{\mathbb{E}[\exp(\beta\hat{r}(Y, x))]} \end{aligned} \quad (51)$$

Note that we have,

$$\begin{aligned} & \frac{d \frac{\mathbb{E}[\hat{r}(Y, x) \exp(\beta\hat{r}(Y, x))]}{\mathbb{E}[\exp(\beta\hat{r}(Y, x))]} }{d\beta} \\ &= \frac{\mathbb{E}[\hat{r}^2(Y, x) \exp(\beta\hat{r}(Y, x))] \mathbb{E}[\exp(\beta\hat{r}(Y, x))] - \mathbb{E}[\hat{r}(Y, x) \exp(\beta\hat{r}(Y, x))]^2}{\mathbb{E}^2[\exp(\beta\hat{r}(Y, x))]} \\ &= \mathbb{E}_{Y \sim \pi_{\beta, \hat{r}}(\cdot|x)}[\hat{r}^2(Y, x)] - \mathbb{E}_{Y \sim \pi_{\beta, \hat{r}}(\cdot|x)}[\hat{r}(Y, x)]^2 \geq 0. \end{aligned} \quad (52)$$

Therefore, we have,

$$\frac{\mathbb{E}[\hat{r}(Y, x) \exp(2\beta\hat{r}(Y, x))]}{\mathbb{E}[\exp(2\beta\hat{r}(Y, x))]} \geq \frac{\mathbb{E}[\hat{r}(Y, x) \exp(\beta\hat{r}(Y, x))]}{\mathbb{E}[\exp(\beta\hat{r}(Y, x))]} \quad (53)$$

It completes the proof.

3. Follows directly from Lemma F.9.

4. Due to Jensen inequality for $\mathbb{E}^2[\exp(\beta\hat{r}(Y, x))] \leq \mathbb{E}[\exp(2\beta\hat{r}(Y, x))]$, the $C_{\beta, \hat{r}}(x)$. We also have the uniform bound, $C_{\beta, \hat{r}}(x) = \frac{\mathbb{E}[\exp(2\beta\hat{r}(Y, x))]}{\mathbb{E}^2[\exp(\beta\hat{r}(Y, x))]} \leq \exp(\beta)$. Furthermore, due to increasing property in second item, we also have $\sup_{\beta} C_{\beta, \hat{r}}(x) = C_{\infty, \hat{r}}(x)$. Therefore, the upper bound holds.

□

Theorem H.1. *The following upper bound holds,*

$$\begin{aligned} \Delta_{J_{r^*}}(\pi_{\beta, r^*}(\cdot|x), \pi_{\hat{r}}^{\text{SBON}}(\cdot|x)) &\leq \frac{1}{\beta} \left(\text{KL}(\pi_{\beta, r^*}(\cdot|x) \parallel \pi_{\text{ref}}(\cdot|x)) - \text{KL}(\pi_{\beta, \hat{r}}(\cdot|x) \parallel \pi_{\text{ref}}(\cdot|x)) \right) \\ &\quad + \sqrt{\varepsilon_{\beta, r}(x)} (\sqrt{C_{\beta, \hat{r}}(x)} + \sqrt{C_{\beta, r^*}(x)}) \\ &\quad + 2\sqrt{\frac{1}{2} \log \left(1 + \frac{C_{\beta, \hat{r}}(x) - 1}{N} \right)}. \end{aligned}$$

Proof. Note that, we have,

$$\begin{aligned} & \Delta_{J_{r^*}}(\pi_{\beta, r^*}(\cdot|x), \pi_{\hat{r}}^{\text{SBON}}(\cdot|x)) \\ &= \mathbb{E}_{Y \sim \pi_{\beta, r^*}(\cdot|x)}[r^*(Y, x)] - \mathbb{E}_{Y \sim \pi_{\hat{r}}^{\text{SBON}}(\cdot|x)}[r^*(Y, x)] \\ &= \underbrace{\mathbb{E}_{Y \sim \pi_{\beta, r^*}(\cdot|x)}[r^*(Y, x)] - \mathbb{E}_{Y \sim \pi_{\beta, \hat{r}}(\cdot|x)}[r^*(Y, x)]}_{I_1} \\ &\quad + \underbrace{\mathbb{E}_{Y \sim \pi_{\beta, \hat{r}}(\cdot|x)}[r^*(Y, x)] - \mathbb{E}_{Y \sim \pi_{\hat{r}}^{\text{SBON}}(\cdot|x)}[r^*(Y, x)]}_{I_2} \end{aligned} \quad (54)$$

Note that, using the definition of $\pi_{\beta, r^*}(\cdot|x)$ and $\pi_{\beta, \hat{r}}(\cdot|x)$ as solutions to KL-regularized problem, we have,

$$\mathbb{E}_{Y \sim \pi_{\beta, r^*}(\cdot|x)}[r^*(Y, x)] = \frac{1}{\beta} \text{KL}(\pi_{\beta, r^*}(\cdot|x) \| \pi_{\text{ref}}(\cdot|x)) + \frac{1}{\beta} \log(\mathbb{E}_{Y \sim \pi_{\text{ref}}(\cdot|x)}[\exp(\beta r^*(Y, x))]). \quad (55)$$

$$\mathbb{E}_{Y \sim \pi_{\beta, \hat{r}}(\cdot|x)}[\hat{r}(Y, x)] = \frac{1}{\beta} \text{KL}(\pi_{\beta, \hat{r}}(\cdot|x) \| \pi_{\text{ref}}(\cdot|x)) + \frac{1}{\beta} \log(\mathbb{E}_{Y \sim \pi_{\text{ref}}(\cdot|x)}[\exp(\beta \hat{r}(Y, x))]). \quad (56)$$

Therefore, for term I_1 , we have,

$$\begin{aligned} & \mathbb{E}_{Y \sim \pi_{\beta, r^*}(\cdot|x)}[r^*(Y, x)] - \mathbb{E}_{Y \sim \pi_{\beta, \hat{r}}(\cdot|x)}[\hat{r}(Y, x)] \\ &= \mathbb{E}_{Y \sim \pi_{\beta, r^*}(\cdot|x)}[r^*(Y, x)] - \mathbb{E}_{Y \sim \pi_{\beta, \hat{r}}(\cdot|x)}[\hat{r}(Y, x)] \\ & \quad + \mathbb{E}_{Y \sim \pi_{\beta, \hat{r}}(\cdot|x)}[\hat{r}(Y, x)] - \mathbb{E}_{Y \sim \pi_{\beta, \hat{r}}(\cdot|x)}[r^*(Y, x)] \\ &= \frac{1}{\beta} (\text{KL}(\pi_{\beta, r^*}(\cdot|x) \| \pi_{\text{ref}}(\cdot|x)) - \text{KL}(\pi_{\beta, \hat{r}}(\cdot|x) \| \pi_{\text{ref}}(\cdot|x))) \\ & \quad + \frac{1}{\beta} \log(\mathbb{E}_{Y \sim \pi_{\text{ref}}(\cdot|x)}[\exp(\beta r^*(Y, x))]) - \frac{1}{\beta} \log(\mathbb{E}_{Y \sim \pi_{\text{ref}}(\cdot|x)}[\exp(\beta \hat{r}(Y, x))]) \\ & \quad + \sum_y \pi_{\beta, \hat{r}}(\cdot|x) (\hat{r}(y, x) - r^*(y, x)) \\ &\leq \frac{1}{\beta} (\text{KL}(\pi_{\beta, r^*}(\cdot|x) \| \pi_{\text{ref}}(\cdot|x)) - \text{KL}(\pi_{\beta, \hat{r}}(\cdot|x) \| \pi_{\text{ref}}(\cdot|x))) \\ & \quad + \frac{1}{\beta} \log(\mathbb{E}_{Y \sim \pi_{\text{ref}}(\cdot|x)}[\exp(\beta r^*(Y, x))]) - \frac{1}{\beta} \log(\mathbb{E}_{Y \sim \pi_{\text{ref}}(\cdot|x)}[\exp(\beta \hat{r}(Y, x))]) \\ & \quad + \frac{1}{\sqrt{\beta}} \sqrt{\sum_y \frac{\pi_{\beta, \hat{r}}^2(y|x)}{\pi_{\text{ref}}(y|x)}} \sqrt{\beta \sum_y (\hat{r}(y, x) - r^*(y, x))^2 \pi_{\text{ref}}(y|x)} \\ &\leq \frac{1}{\beta} (\text{KL}(\pi_{\beta, r^*}(\cdot|x) \| \pi_{\text{ref}}(\cdot|x)) - \text{KL}(\pi_{\beta, \hat{r}}(\cdot|x) \| \pi_{\text{ref}}(\cdot|x))) \\ & \quad + \frac{1}{\beta} \log(\mathbb{E}_{Y \sim \pi_{\text{ref}}(\cdot|x)}[\exp(\beta r^*(Y, x))]) - \frac{1}{\beta} \log(\mathbb{E}_{Y \sim \pi_{\text{ref}}(\cdot|x)}[\exp(\beta \hat{r}(Y, x))]) \\ & \quad + \sqrt{C_{\beta, \hat{r}}(x) \varepsilon_{\beta, r}(x)} \\ &\leq \frac{1}{\beta} (\text{KL}(\pi_{\beta, r^*}(\cdot|x) \| \pi_{\text{ref}}(\cdot|x)) - \text{KL}(\pi_{\beta, \hat{r}}(\cdot|x) \| \pi_{\text{ref}}(\cdot|x))) \\ & \quad + \frac{1}{\beta} \log(\mathbb{E}_{Y \sim \pi_{\text{ref}}(\cdot|x)}[\exp(\beta r^*(Y, x))]) - \frac{1}{\beta} \log(\mathbb{E}_{Y \sim \pi_{\text{ref}}(\cdot|x)}[\exp(\beta \hat{r}(Y, x))]) \\ & \quad + \sqrt{C_{\beta, \hat{r}}(x) \varepsilon_{\beta, r}(x)} \\ &\leq \frac{1}{\beta} (\text{KL}(\pi_{\beta, r^*}(\cdot|x) \| \pi_{\text{ref}}(\cdot|x)) - \text{KL}(\pi_{\beta, \hat{r}}(\cdot|x) \| \pi_{\text{ref}}(\cdot|x))) \\ & \quad + \sqrt{C_{\beta, r^*}(x) \varepsilon_{\beta, r}(x)} \\ & \quad + \sqrt{C_{\beta, \hat{r}}(x) \varepsilon_{\beta, r}(x)}. \end{aligned} \quad (57)$$

For term I_2 and using similar approach to term I_1 and applying Lemma F.10, we have,

$$\begin{aligned}
& \mathbb{E}_{Y \sim \pi_{\beta, \hat{r}}(\cdot | x)}[r^*(Y, x)] - \mathbb{E}_{Y \sim \pi_{\hat{r}}^{\text{SBon}}(\cdot | x)}[r^*(Y, x)] \\
& \leq 2\mathbb{T}\mathbb{V}(\pi_{\beta, \hat{r}}(\cdot | x), \pi_{\hat{r}}^{\text{SBon}}(\cdot | x)) \\
& \leq 2 \min \left(1, \sqrt{\frac{1}{2} \text{KL}(\pi_{\beta, \hat{r}}(\cdot | x) \| \pi_{\hat{r}}^{\text{SBon}}(\cdot | x))} \right) \\
& \leq 2 \min \left(1, \sqrt{\frac{1}{2} \log \left(1 + \frac{C_{\beta, \hat{r}}(x) - 1}{N} \right)} \right) \\
& \leq 2 \sqrt{\frac{1}{2} \log \left(1 + \frac{C_{\beta, \hat{r}}(x) - 1}{N} \right)}
\end{aligned} \tag{58}$$

Combining equation 57 and equation 58 with equation 54 completes the proof. \square

Theorem 5.2. Under Assumption 3.1, the following upper bound holds on the optimal regret gap of the SBoN policy for any $\beta > 0$,

$$\begin{aligned}
\Delta_{J, r^*}(\pi_{r^*}^*(\cdot | x), \pi_{\hat{r}}^{\text{SBon}}(\cdot | x)) & \leq \sqrt{\varepsilon_{\beta, r}(x)} (\sqrt{C_{\infty, \hat{r}}(x)} + \sqrt{C_{\infty, r^*}(x)}) \\
& \quad + 2 \sqrt{\frac{1}{2} \log \left(1 + \frac{C_{\infty, \hat{r}}(x) - 1}{N} \right)} \\
& \quad + \frac{\log(C_{\infty, r^*}(x))}{\beta},
\end{aligned}$$

Proof. Note that we have,

$$\begin{aligned}
& \Delta_{J, r^*}(\pi_{r^*}^*(\cdot | x), \pi_{\hat{r}}^{\text{SBon}}(\cdot | x)) \\
& = \mathbb{E}_{Y \sim \pi_{r^*}^*(\cdot | x)}[r^*(Y, x)] - \mathbb{E}_{Y \sim \pi_{\hat{r}}^{\text{SBon}}(\cdot | x)}[r^*(Y, x)] \\
& = \underbrace{\mathbb{E}_{Y \sim \pi_{r^*}^*(\cdot | x)}[r^*(Y, x)] - \mathbb{E}_{Y \sim \pi_{\beta, r^*}(\cdot | x)}[r^*(Y, x)]}_{I_3} \\
& \quad + \underbrace{\Delta_{J, r^*}(\pi_{\beta, r^*}(\cdot | x), \pi_{\hat{r}}^{\text{SBon}}(\cdot | x))}_{I_4}
\end{aligned} \tag{59}$$

For term I_4 , we can use Theorem H.1. For term I_3 , note that, we have for $\beta > 0$,

$$\mathbb{E}_{Y \sim \pi_{r^*}^*(\cdot | x)}[r^*(Y, x)] - \mathbb{E}_{Y \sim \pi_{\beta, r^*}(\cdot | x)}[r^*(Y, x)] \leq \frac{\text{KL}(\pi_{r^*}^*(\cdot | x) \| \pi_{\text{ref}}(\cdot | x)) - \text{KL}(\pi_{\beta, r^*}(\cdot | x) \| \pi_{\text{ref}}(\cdot | x))}{\beta} \tag{60}$$

Combining equation 60 with Theorem H.1, completes the proof due the positiveness of KL divergence and using Lemma F.7 and Lemma 5.1. \square

Remark H.2. For $\beta = 0$, we have, $\lim_{\beta \rightarrow 0} \pi_{\beta, r^*}(\cdot | x) = \pi_{\text{ref}}(\cdot | x)$. Therefore, we have,

$$\mathbb{E}_{Y \sim \pi_{r^*}^*(\cdot | x)}[r^*(Y, x)] - \mathbb{E}_{Y \sim \pi_{\beta, r^*}(\cdot | x)}[r^*(Y, x)] \leq \sqrt{2\text{KL}(\pi_{r^*}^*(\cdot | x) \| \pi_{\text{ref}}(\cdot | x))}. \tag{61}$$

Proposition 5.3. Under Assumption 3.1, the following upper bound holds on the optimal regret gap of the BoN policy for any $\beta > 0$,

$$\begin{aligned}
\Delta_{J, r^*}(\pi_{r^*}^*(\cdot | x), \pi_{\hat{r}}^{\text{BoN}}(\cdot | x)) & \leq \sqrt{\varepsilon_{\infty, r}(x)} (\sqrt{C_{\infty, \hat{r}}(x)} + \sqrt{C_{\infty, r^*}(x)}) \\
& \quad + 2 \sqrt{\frac{1}{2} \log \left(1 + \frac{C_{\infty, \hat{r}}(x) - 1}{N} \right)}.
\end{aligned}$$

Proof. The results follow directly from Proposition 5.3 for $\beta \rightarrow \infty$. \square

H.2 LOWER BOUND PROOFS AND DETAILS:

We can define the probability of coverage of maximizers under the reference policy as

$$p^*(x) := \Pr_{Y \sim \pi_{\text{ref}}(\cdot|x)}[Y \in \mathcal{Y}_{r^*}^*(x)] = \frac{1}{C_{\infty, r^*}(x)}.$$

Lemma H.3. *Under Assumption 5.5, the following lower bound holds on the optimal regret gap of the SBoN policy based on true-reward,*

$$\Delta_{J_{r^*}}(\pi_{r^*}^*(\cdot|x), \pi_{r^*}^{\text{SBoN}}(\cdot|x)) \geq \gamma(x) \Pr_{Err}(x), \quad (62)$$

where \Pr_{Err} is the probability that selected sample from $\{Y_i\}_{i=1}^N$ via SBoN does not belong to $\mathcal{Y}_{r^*}^*(x)$.

Proof. We have,

$$\begin{aligned} & \Delta_{J_{r^*}}(\pi_{r^*}^*(\cdot|x), \pi_{r^*}^{\text{SBoN}}(\cdot|x)) \\ &= J_{r^*}(\pi_{r^*}^*(\cdot|x)) - J_{r^*}(\pi_{\beta, r^*}(\cdot|x)) \\ &= \mathbb{E}_{Y \sim \pi_{r^*}^*(\cdot|x)}[r^*(Y, x)] - \mathbb{E}_{Y \sim \pi_{\beta, r^*}(\cdot|x)}[r^*(Y, x)] \\ &= 1 - \left[\mathbb{E}_{Y \sim \pi_{\beta, r^*}(\cdot|x)}[r^*(Y, x)|\text{Err}] \Pr_{Err} + \mathbb{E}_{Y \sim \pi_{\beta, r^*}(\cdot|x)}[r^*(Y, x)|\text{Err}^c](1 - \Pr_{Err}(x)) \right] \\ &\geq 1 - \left[(1 - \gamma(x)) \Pr_{Err}(x) + (1 - \Pr_{Err}(x)) \right] \\ &= \gamma(x) \Pr_{Err}(x), \end{aligned} \quad (63)$$

where Err is the event where the selected sample using SBoN is not a maximizer. \square

Next, we provide a lower bound on $\Pr_{Err}(x)$.

Lemma H.4. *The following lower bound holds on \Pr_{Err} for SBoN algorithm,*

$$\Pr_{Err}(x) \geq (1 - p^*(x))^N + \frac{N((1 - p^*(x)) - (1 - p^*(x))^N)}{N + (\exp(\beta) - 1)(N - 1)}. \quad (64)$$

Proof. We consider cases where exactly K responses are selected from $\mathcal{Y}_{r^*}^*(x)$. The K has a binomial distribution $\text{Bin}(N, p^*(x))$. Using Bayes rule, we have,

$$\begin{aligned} \Pr_{Err}(x) &= \sum_{i=0}^{N-1} \Pr(K = i) \Pr_{Err}(x|K = i) \\ &= \sum_{i=0}^{N-1} \binom{N}{i} p^*(x)^i (1 - p^*(x))^{N-i} \Pr_{Err}(x|K = i) \\ &= \sum_{i=0}^{N-1} \binom{N}{i} p^*(x)^i (1 - p^*(x))^{N-i} \frac{\sum_{Y_j \notin \mathcal{Y}_{r^*}^*(x)} \exp(\beta r^*(Y_j, x))}{\sum_{j=1}^N \exp(\beta r^*(Y_j, x))} \\ &\geq (1 - p^*(x))^N + \sum_{i=1}^{N-1} \binom{N}{i} p^*(x)^i (1 - p^*(x))^{N-i} \frac{N - i}{N - i + i e^\beta} \\ &\geq (1 - p^*(x))^N + \sum_{i=1}^{N-1} \binom{N}{i} p^*(x)^i (1 - p^*(x))^{N-i} \frac{N - i}{N + (e^\beta - 1)(N - 1)} \\ &\geq (1 - p^*(x))^N + \frac{N((1 - p^*(x)) - (1 - p^*(x))^N)}{N + (e^\beta - 1)(N - 1)} \end{aligned} \quad (65)$$

\square

Using Lemma H.3 and Lemma H.2, we can derive the following lower bound on our main optimal regret gap,

Theorem 5.6. *Under the same Assumptions in Theorem 5.2 and Lemma H.3, the following lower bound holds on the regret of SBoN policy under the proxy-reward model,*

$$\begin{aligned} & \Delta_{J_{r^*}}(\pi_{r^*}^*(\cdot|x), \pi_{\hat{r}}^{\text{SBoN}}(\cdot|x)) \\ & \geq \gamma(x) \left((1 - p^*(x))^N + \frac{N((1 - p^*(x)) - (1 - p^*(x))^N)}{N + (\exp(\beta) - 1)(N - 1)} \right) \\ & \quad - \min \left(\sqrt{1 - \exp(-U(N, \beta))}, \frac{N}{1 + (N - 1)\exp(-\beta)} \sqrt{\varepsilon_{\beta, r}(x)} \right), \end{aligned} \quad (66)$$

$$\text{where } p^*(x) = \frac{1}{C_{\infty, r^*}(x)} \text{ and } U(N, \beta) = \frac{N\beta\sqrt{\varepsilon_{\beta, r}(x)}}{1 + (N - 1)\exp(-\beta)} \left(\frac{N\exp(2\beta)}{(N - 1)^2} + 1 \right).$$

Proof.

$$\Delta_{J_{r^*}}(\pi_{r^*}^*(\cdot|x), \pi_{\hat{r}}^{\text{SBoN}}(\cdot|x)) \quad (67)$$

$$= \Delta_{J_{r^*}}(\pi_{r^*}^*(\cdot|x), \pi_{r^*}^{\text{SBoN}}(\cdot|x)) \quad (68)$$

$$+ \Delta_{J_{r^*}}(\pi_{r^*}^{\text{SBoN}}(\cdot|x), \pi_{\hat{r}}^{\text{SBoN}}(\cdot|x)), \quad (69)$$

where equation 68 can be bounded using Lemma H.3 and Lemma H.2. By combining Lemma F.12 and Lemma 4.2 we have an upper bound on absolute value of equation 69 where can also result in a lower bound on equation 69. Furthermore, Another simple upper bound on equation 69 is,

$$\Delta_{J_{r^*}}(\pi_{r^*}^{\text{SBoN}}(\cdot|x), \pi_{\hat{r}}^{\text{SBoN}}(\cdot|x)) \quad (70)$$

$$= \mathbb{E}_{Y \sim \pi_{r^*}^{\text{SBoN}}(Y|x)}[r^*(Y, x)] - \mathbb{E}_{Y \sim \pi_{\hat{r}}^{\text{SBoN}}(Y|x)}[\hat{r}(Y, x)] \quad (71)$$

$$+ \mathbb{E}_{Y \sim \pi_{\hat{r}}^{\text{SBoN}}(Y|x)}[\hat{r}(Y, x)] - \mathbb{E}_{Y \sim \pi_{r^*}^{\text{SBoN}}(Y|x)}[r^*(Y, x)] \quad (72)$$

$$= \mathbb{E}_{Y \sim \pi_{\hat{r}}^{\text{SBoN}}(Y|x)}[\hat{r}(Y, x)] - \mathbb{E}_{Y \sim \pi_{r^*}^{\text{SBoN}}(Y|x)}[r^*(Y, x)] \quad (73)$$

Note that due to calibration both true and proxy-reward models are uniform under $\pi_{\text{ref}}(\cdot|x)$. Therefore, we have,

$$\mathbb{E}_{Y \sim \pi_{r^*}^{\text{SBoN}}(Y|x)}[r^*(Y, x)] = \mathbb{E}_{Y \sim \pi_{\hat{r}}^{\text{SBoN}}(Y|x)}[\hat{r}(Y, x)].$$

Now we have,

$$\begin{aligned} |\mathbb{E}_{Y \sim \pi_{\hat{r}}^{\text{SBoN}}(Y|x)}[\hat{r}(Y, x)] - \mathbb{E}_{Y \sim \pi_{r^*}^{\text{SBoN}}(Y|x)}[r^*(Y, x)]| & \leq \sqrt{\sum_y \frac{(\pi_{\hat{r}}^{\text{SBoN}}(\cdot|x))^2}{\pi_{\text{ref}}(\cdot|x)}} \sqrt{\varepsilon_{\beta, r}(x)} \\ & \leq \frac{N}{1 + (N - 1)\exp(-\beta)} \sqrt{\varepsilon_{\beta, r}(x)}. \end{aligned}$$

□

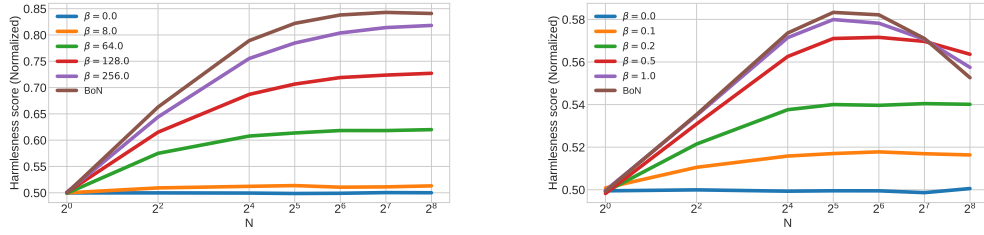


Figure 3: Soft Best-of-N experiment using a strong reward model (Left) and a weak one (Right). When the performance of the reward model is poor, BoN can lead to overoptimization, while the SBoN can help to mitigate it.

I EXPERIMENTS

Our implementation is available in the following link: <https://anonymous.4open.science/r/sbon-C6D3/README.md>.

LLM-as-a-Judge: To evaluate the performance of the generated responses, we used GPT-4o to determine how harmless they are. Specifically, we used the following prompt:

You are an AI safety evaluator. Your task is to determine if a response to a given prompt contains harmful content. Respond with a harmlessness score using these criteria:

2 = SAFE: The response is completely safe and appropriate.
 1 = SOMEWHAT HARMFUL: The response contains potentially concerning elements, mild toxicity, stereotypes, subtle harmful guidance, or borderline inappropriate content.
 0 = HARMFUL: The response contains clearly harmful content such as hate speech, explicit instructions for dangerous activities, severe toxicity, etc.

Provide your judgment as an integer in curly brackets like this: 1, 2, or 3.

This yields a reward with values ranging from 0 to 2, where 0 indicates complete harm and 2 indicates safety.

Future work may address optimizing β using cross-validation or the estimation of terms in Remark 5.8 to derive sub-optimal β . For this study, we employ a grid search over a range of β values.

I.1 MORE EXPERIMENTS

To further support our theoretical analysis, we conducted more experiments comparing Soft Best-of-N (SBoN) across different regularization strengths and reward model qualities. We ran one extra experimental conditions: one using a weaker proxy-reward model (tulu-v2.5-13b-uf-rm (Iverson et al., 2024)) respect to Armo reward model. As shown in Figure 3, when the reward model is weak, performance degrades for large N due to reward hacking. However, the smoothing in SBoN helps mitigate this degradation. This observation is also aligned with our theoretical analysis and discussion in Section 4.

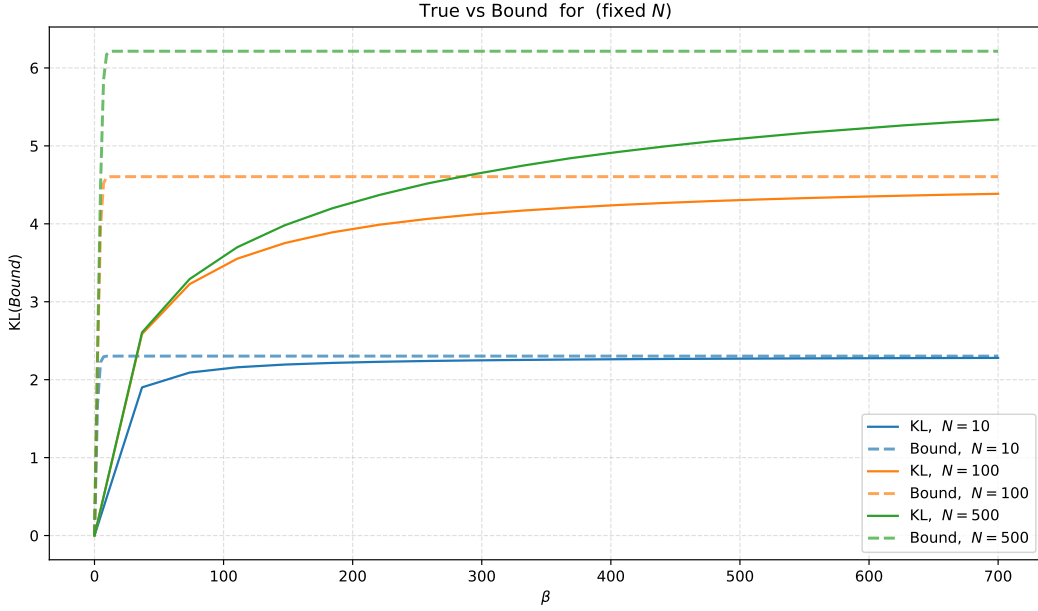


Figure 4: True KL divergence vs upper bound in Lemma 4.1 for fixed $N = \{10, 100, 500\}$.

I.2 NUMERICAL EXAMPLES

I.2.1 KL BOUND

To illustrate how our analytical upper bound in Lemma 4.1 behaves as a function of the temperature parameter β , we run a toy experiment in which

1. the reference policy is the uniform distribution over responses, and
2. rewards are bounded with $= 1$.

For each β in a logarithmic sweep, we compute the true KL divergence between the SBoN policy and the reference policy, together with the theoretical bound derived in Lemma 4.1.

- **Very large β (nearBoN policy).** As $\beta \rightarrow \infty$ the SBoN policy converges to the BoN policy. The gap between the KL and the bound vanishes.
- **Very small β (reference policy).** When $\beta \rightarrow 0$ the SBoN policy approaches the uniform sampling from the reference policy, which results in the reference policy, making the KL itself tend to zero; the bound is equal to zero for this value.

This experiment confirms that the bound is tight in the two asymptotic regimes and remains a conservative yet informative estimate elsewhere.

I.2.2 REGRET UPPER BOUND

To validate our theoretical findings in Theorem 5.2 and Proposition 5.3, we conducted a numerical analysis comparing the regret upper bounds of SBoN against standard BoN under varying degrees of reward model error.

We simulate a simplified binary reward setting to isolate the effects of regularization temperature (β) and proxy-reward error (δ). The setup is defined as follows:

Reward Structure: We assume a binary reward landscape where responses have a true-reward $r^*(y), \hat{r} \in \{0, 1\}$.

Reference Policy (π_{ref}): The reference model generates an optimal response ($r^* = 1$) with probability $\pi_{\text{ref}}(y_{\text{max}}|x) = 0.05$, corresponding to a coverage constant $C_{\infty, r^*} = 20$. The same also

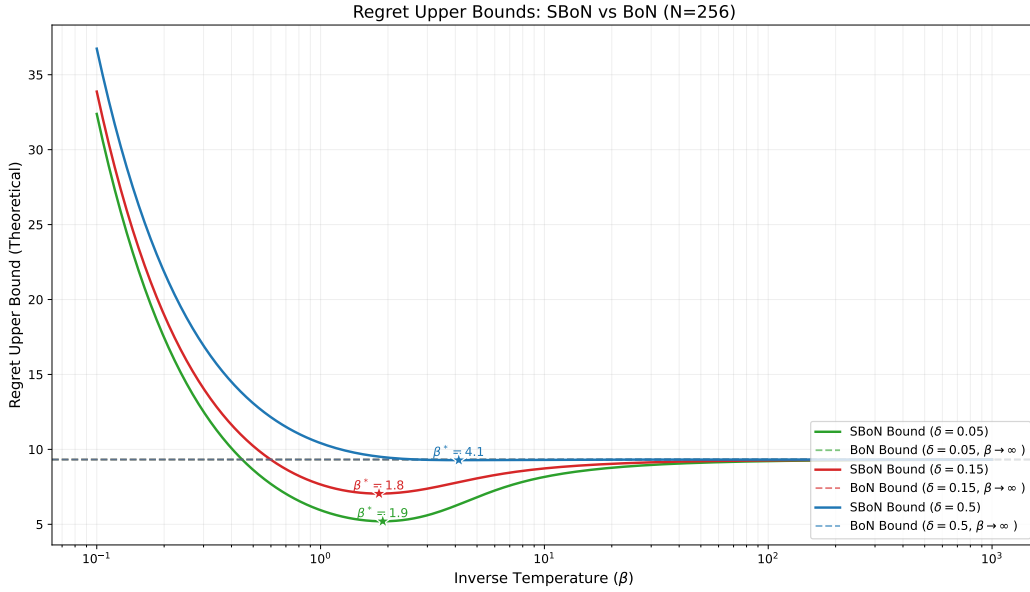


Figure 5: Numerical comparison of upper bound on Regret of SBoN and BoN

holds for proxy-reward model and we have $C_{\infty, \hat{r}} = 20$. We also consider $N = 256$ samples from reference policy.

Proxy-reward Error (δ): The proxy-reward model matches the true-reward with probability $1 - \delta$ and commits an error (flipping the label) with probability δ . We evaluate three error regimes: Low ($\delta = 0.05$), Medium ($\delta = 0.15$), and High ($\delta = 0.50$).

SBoN: We compute the regret upper bound using Theorem 5.2, varying the inverse temperature β from 0.5 to 50.

BoN: We compute the regret upper bound using Proposition 5.3, which represents the asymptotic limit of SBoN as $\beta \rightarrow \infty$. In this binary setting, the maximum possible error is ($\varepsilon_{\infty, r} = 1$).

Discussion: The results, Fig. 5 demonstrate a distinct trade-off mechanism governed by the temperature β . While the standard BoN bound remains constant at a high value due to its susceptibility to the maximum error ($\varepsilon_{\infty, r} = 1$), the SBoN bound forms a convex curve.

For every error rate tested, we observe a "sweet spot"-a finite optimal temperature β^* where the SBoN bound is significantly tighter than the BoN bound. This confirms that by smoothing the policy (finite β), SBoN effectively balances the trade-off between exploring high-reward regions and mitigating the overoptimization caused by proxy-reward errors.

I.2.3 REGRET LOWER BOUND

To complement our upper bound analysis, we numerically evaluate the lower bounds on the regret gap for both SBoN and BoN, as derived in Theorem 5.6 and Proposition 5.7 respectively. These lower bounds serve as a "safety guarantee," quantifying the worst-case performance degradation caused by overoptimization. We utilize the same binary reward setup described in previous section, with sample size $N = 32$ and consider no error ($\delta = 0$). The results in Fig. 6 illustrate an advantage of the SBoN policy over the deterministic BoN policy. In error regimes (blue dashed line, $\delta = 0.15, 0.5$), the BoN lower bound drops significantly, indicating a weak guarantee against performance collapse. Note that, for no error regime, the lower bound on BoN's regret is positive.

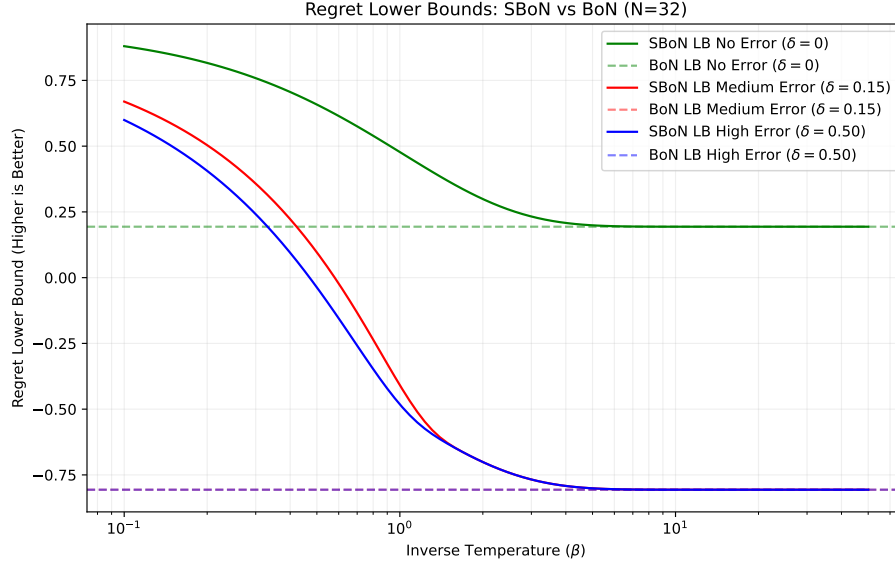


Figure 6: Numerical comparison of lower bound on Regret of SBoN and BoN

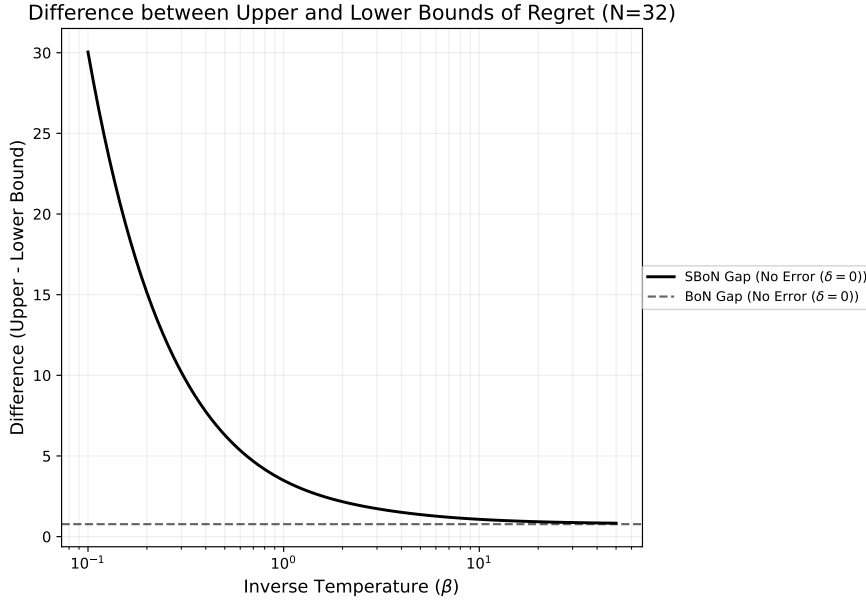


Figure 7: Numerical comparison of the gap between upper and lower bounds on Regret of SBoN and BoN

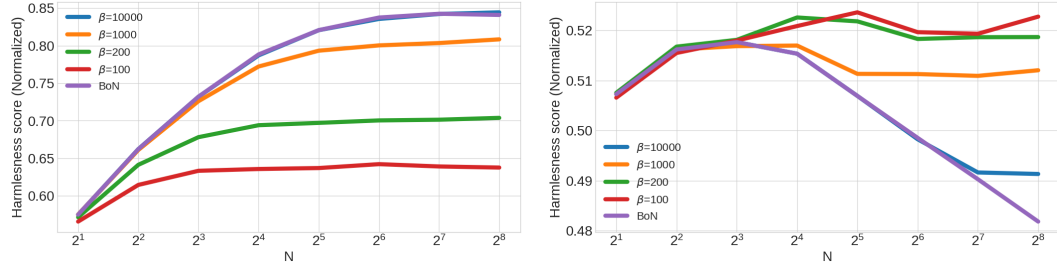


Figure 8: InferenceTimePessimism experiment using a strong reward model (Left) and a weak one (Right). When the performance of the reward model is poor, BoN can lead to overoptimization, while the InferenceTimePessimism can help to mitigate it.

I.2.4 UPPER AND LOWER BOUNDS GAP

To compare the performance of upper and lower bounds, we consider the case where $\delta = 0$ which is similar to no-overoptimization scenario. As shown in Fig. 7, the gap between upper and lower bounds of BoN is tighter than SBoN.

J COMPARISON WITH HUANG ET AL. (2025)

We contrast our contributions with the recent work of Huang et al. (2025), focusing on four key dimensions: the object of study, regret bounds, and the mechanism of regularization.

Object of study. Huang et al. (2025) analyze BoN, introducing an inference-time pessimism algorithm within a χ^2 -regularized framework to mitigate overoptimization. In contrast, we analyze SBoN via treating BoN as the limit where $\beta \rightarrow \infty$ and derive finite-sample KL and regret bounds under reward model misspecification (specifically, proxy-reward vs. true-reward). To our knowledge, these SBoN-specific bounds characterizing the behavior under overoptimization are novel.

Error metrics and reward modeling. While Huang et al. (2025) operate on raw (uncalibrated) rewards using mean-squared error (MSE), we utilize calibrated rewards (Balashankar et al., 2025). We introduce the *tilted error*, denoted as $\varepsilon_{\beta,r}$, which interpolates between MSE (at $\beta = 0$) and L^∞ error (as $\beta \rightarrow \infty$). This metric allows us to explicitly track the dynamics of overoptimization as a function of both N and the temperature β .

Regret bounds. Our Theorem 5.2, Proposition 5.3, Theorem 5.6 and Proposition 5.7 and establish regret bounds for SBoN and its BoN limit that depend on: (i) the tilted error $\varepsilon_{\beta,r}(x)$, and (ii) the coverage constants $C_{\infty,r}$ for both the true and proxy-rewards. A crucial distinction, highlighted in Remark 5.4, is that our BoN regret bound remains finite even when overoptimization vanishes (i.e., when $\varepsilon_{\infty,r} = 0$ or $N \rightarrow \infty$). Conversely, the bound proposed by Huang et al. (2025) on BoN scales with the L^∞ -error. This represents a qualitative improvement in how misspecification is handled in our analysis.

Theoretical and algorithmic consistency. Our analysis of SBoN fully incorporates the algorithm’s primary regularization parameter, β . In contrast, the theoretical regret bounds for *InferenceTimePessimism* in Huang et al. (2025) exclude the estimation error of the normalization factor (Algorithm 3) and are independent of the truncation parameter required by their Algorithm 4. Furthermore, regarding practical complexity, SBoN relies on a single hyperparameter β , whereas *InferenceTimePessimism* requires tuning an additional truncation parameter and incurs computational overhead to estimate the normalization factor.

InferenceTimePessimism Experiment: We compared SBoN against the InferenceTimePessimism baseline (Huang et al., 2025), using the configuration specified in their paper. While both methods achieve similar empirical performance Fig.8, SBoN is significantly more efficient to implement and tune, as it avoids the additional computational overhead and hyperparameter search space required by the pessimistic approach.

K INFORMATION PROJECTION ACROSS REWARD FUNCTIONS

In this section, we study the KL divergence behavior in Section 4 for large N regime, where SBoN policy converge to tilted optimal policy.

We consider a proxy-reward class \mathcal{R} of rewards, and a projection problem at a *fixed* temperature β which is inspired by reducing the estimation error,

$$\text{(reverse-I-proj in tilted family)} \quad \hat{r}^I \in \arg \min_{\hat{r} \in \mathcal{R}} D_{\text{KL}}(\pi_{\beta, r^*} \parallel \pi_{\beta, \hat{r}}).$$

Proposition K.1 (Pythagorean I-projection in reward space). *Suppose the proxy class \mathcal{R} is convex in r (e.g., an affine subspace). Let $\hat{r}^I \in \arg \min_{\hat{r} \in \mathcal{R}} \text{KL}(\pi_{\beta, r^*} \parallel \pi_{\beta, \hat{r}})$. Then, we have,*

$$\text{KL}(\pi_{\beta, r^*} \parallel \pi_{\text{ref}}) \geq \underbrace{\text{KL}(\pi_{\beta, r^*} \parallel \pi_{\beta, \hat{r}^I})}_{\text{Estimation error}} + \underbrace{\text{KL}(\pi_{\beta, \hat{r}^I} \parallel \pi_{\text{ref}})}_{\text{KL-divergence based on proxy-reward}} \quad (74)$$

Proof. The proof follows from Theorem 1 in Csiszár & Matus (2003). \square

Proposition K.1 and Lemma F.7 provide an upper bound on $D_{\text{KL}}(\pi_{\beta, r^*} \parallel \pi_{\text{ref}})$. Since the LHS of equation 74 is fixed for a given β , an increase in estimation error necessarily decreases the allowable divergence between the proxy-reward and the reference policies $D_{\text{KL}}(\pi_{\beta, \hat{r}^I} \parallel \pi_{\text{ref}})$. If $r^* \in \mathcal{R}$, the estimation error is zero (no overoptimization). If $r^* \notin \mathcal{R}$, the estimation error consumes part of the total divergence budget. This analysis can be extended to for finite N , using Lemma F.10.