
Theoretical Guarantees on the Best-of- n Alignment Policy

Ahmad Beirami¹ Alekh Agarwal² Jonathan Berant¹ Alexander D’Amour¹ Jacob Eisenstein¹ Chirag Nagpal²
Ananda Theertha Suresh²

Abstract

A simple and effective method for the inference-time alignment and scaling test-time compute of generative models is best-of- n sampling, where n samples are drawn from a reference policy, ranked based on a reward function, and the highest ranking one is selected. A commonly used analytical expression in the literature claims that the KL divergence between the best-of- n policy and the reference policy is equal to $\log(n) - (n - 1)/n$. We disprove the validity of this claim, and show that it is an upper bound on the actual KL divergence. We also explore the tightness of this upper bound in different regimes, and propose a new estimator for the KL divergence and empirically show that it provides a tight approximation. We also show that the win rate of the best-of- n policy against the reference policy is upper bounded by $n/(n + 1)$ and derive bounds on the tightness of this characterization. We conclude with analyzing the tradeoffs between win rate and KL divergence of the best-of- n alignment policy, which demonstrate that very good tradeoffs are achievable with $n < 1000$.

1. Introduction

Generative language models have shown to be effective general purpose tools to solve various problems. While many problems can be solved in a zero-shot manner, the output from the so-called *reference* model may not be outright desirable, e.g., it may violate safety rules or may not solve a math problem correctly. *Alignment* (Christiano et al., 2017; Stiennon et al., 2020; Ouyang et al., 2022; Bai et al., 2022) and *test-time compute scaling* (Brown et al., 2024; Snell et al., 2024) aim at remedying this issue by further nudging the outcome to improve a reward function while not drifting

¹Google DeepMind ²Google Research. Correspondence to: Ahmad Beirami <ahmad.beirami@gmail.com>, Ananda Theertha Suresh <theertha@google.com>.

Proceedings of the 42nd International Conference on Machine Learning, Vancouver, Canada. PMLR 267, 2025. Copyright 2025 by the author(s).

too far from the reference model.

Recently, there has been a proliferation of methods for alignment, which include KL-regularized reinforcement learning (Christiano et al., 2017; Ouyang et al., 2022), controlled decoding (Yang & Klein, 2021; Mudgal et al., 2024), SLiC (Zhao et al., 2022), direct preference optimization (Rafailov et al., 2023), and best-of- n finetuning (Touvron et al., 2023). At their core, these methods try to solve the following regularized optimization problem:¹

$$\max_{\pi(\cdot|\mathbf{x})} E_{\mathbf{y} \sim \pi(\cdot|\mathbf{x})} r(\mathbf{x}, \mathbf{y}) - \beta D_{\text{KL}}(\pi(\cdot|\mathbf{x}) \| \pi_{\text{ref}}(\cdot|\mathbf{x})), \quad (1)$$

where π_{ref} denotes a reference language model; $r(\mathbf{x}, \mathbf{y}) \in \mathbb{R}$ represent a scalar reward associated with response \mathbf{y} for prompt \mathbf{x} ; and the KL divergence $D_{\text{KL}}(q(\cdot|\mathbf{x}) \| p(\cdot|\mathbf{x}))$ is defined as

$$D_{\text{KL}}(q(\cdot|\mathbf{x}) \| p(\cdot|\mathbf{x})) := E_{\mathbf{y} \sim q(\cdot|\mathbf{x})} \log \frac{q(\mathbf{y}|\mathbf{x})}{p(\mathbf{y}|\mathbf{x})}.$$

Note that Equation (1) has a closed-form solution (Korbak et al., 2022b;a):

$$\pi_{\beta}^*(\mathbf{y}|\mathbf{x}) \propto \pi_{\text{ref}}(\mathbf{y}|\mathbf{x}) e^{\frac{1}{\beta} r(\mathbf{x}, \mathbf{y})}, \quad (2)$$

which defines an exponential family of distributions with nice properties (Yang et al., 2024a). We also define the KL divergence averaged over prompts as $D_{\text{KL}}^{\mu}(q \| p) := E_{\mathbf{x} \sim \mu} D_{\text{KL}}(q(\cdot|\mathbf{x}) \| p(\cdot|\mathbf{x}))$, where μ is a distribution over prompts. Notice that a small KL divergence between the aligned policy and the reference policy is desired because it implies that the capabilities of the reference policy are largely preserved (Gao et al., 2023; Coste et al., 2024; Eisenstein et al., 2024), which is also theoretically analyzed by Balashankar et al. (2025, Appendix B).

To compare different alignment techniques, it is customary to produce tradeoff curves that measure expected reward (or win rate) as a function of $D_{\text{KL}}(\pi \| \pi_{\text{ref}})$ for some aligned policy π . Guarantees on the KL divergence capture the preservation of the core capabilities of the model and tighter estimates on the KL divergence help give guarantees that

¹While theoretically we analyze this optimization problem as a function of the prompt \mathbf{x} , in practice we can only solve it by taking another expectation over a set of prompts $\mathbf{x} \sim \mu$.

the model doesn't lose core capabilities that were present in the reference checkpoint. Thus, it is desirable to improve the reward with the least drift measured in KL divergence.

Despite all the advancements in alignment, a simple, popular, and well-performing method for alignment remains to be the *best-of- n* policy (Nakano et al., 2021; Stiennon et al., 2020). In fact, Gao et al. (2023); Mudgal et al. (2024); Eisenstein et al. (2024) show that best-of- n consistently achieves compelling win rate vs KL tradeoff curves, that even dominate those of KL-regularized reinforcement learning and other more involved alignment policies. Llama 2 (Touvron et al., 2023) uses best-of- n as a teacher outcomes to further finetune the base model. Mudgal et al. (2024) extended best-of- n through q -learning to block-wise best-of- n decoding. This has also led to recent research on distilling best-of- n into new models (Gui et al., 2024; Amini et al., 2025; Sessa et al., 2025; Qiu et al., 2024). Hughes et al. (2024); Beetham et al. (2024) use best-of- n as an effective method for jailbreaking. Best-of- n is also used as a strong baseline in scaling inference-time compute (Brown et al., 2024; Snell et al., 2024). This overwhelming empirical success motivates our theoretical investigation of the best-of- n alignment policy.

Subsequent to this work, Yang et al. (2024a) provided theoretical reasoning for the performance of best-of- n by showing it achieves asymptotically optimal reward-KL tradeoffs. Gui et al. (2024) characterized the win rate vs KL gap to be small in the asymptotic regime of a language model whose outcomes have infinitesimally small likelihood. Sun et al. (2024) made best-of- n faster through speculative rejection. Mroueh (2024) provided information-theoretic bounds on reward vs KL tradeoffs for best-of- n .

Best-of- n . Let \mathbf{x} be a given input prompt to the language model. Let $\mathbf{y}_1, \dots, \mathbf{y}_n$ be n i.i.d. samples drawn from $\pi_{\text{ref}}(\cdot|\mathbf{x})$. The best-of- n strategy selects²

$$\mathbf{y} = \mathbf{y}_{k^*} \quad \text{where} \quad k^* := \arg \max_{k \in [n]} r(\mathbf{x}, \mathbf{y}_k). \quad (3)$$

This process inherently leads to sampling from a new policy that is aligned to the reward, denoted by $\pi^{(n)}$. Notice that $\pi^{(1)} = \pi_{\text{ref}}$, and increasing n increases the reward at the cost of drifting away from the base model.

Our goal in this paper is to better understand the best-of- n alignment policy. In particular, we are interested in theoretical guarantees on $D_{\text{KL}}^\mu(\pi^{(n)}|\pi_{\text{ref}})$ for different values of n . A commonly used expression in the literature (Stiennon et al., 2020; Hilton & Gao, 2022; Coste et al., 2024; Gao et al., 2023; Go et al., 2023; Scheurer et al., 2023) claims

$$D_{\text{KL}}^\mu(\pi^{(n)}|\pi_{\text{ref}}) \stackrel{\text{claim}}{=} \widetilde{\text{KL}}_n := \log(n) - (n-1)/n. \quad (4)$$

This formula is commonly used to demonstrate reward-KL

²We define $[n] := \{1, \dots, n\}$.

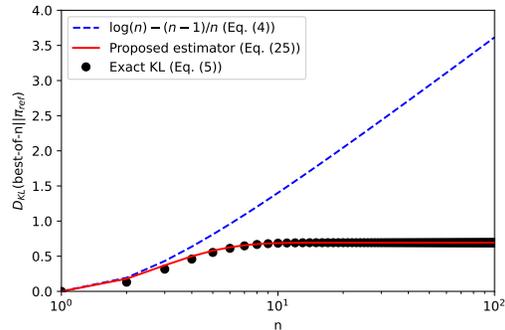


Figure 1. The analytical formula $(\log(n) - (n-1)/n)$ (Equation (4)), the exact KL divergence (Equation (5)), and the proposed estimator (Equation (25)), for Example 1, illustrating a case where the gap between the analytical formula and the exact KL divergence is unbounded.

tradeoffs for the best-of- n policy. Let us further inspect this formula using a toy example.

Example 1. Consider an unprompted model with $\mathbf{x} = \emptyset$ (no input) and binary output, $\mathbf{y} \in \{0, 1\}$. Let the two outcomes be equiprobable, i.e., $\pi_{\text{ref}}(0) = \pi_{\text{ref}}(1) = \frac{1}{2}$. Further, let $r(0) = 0$, and $r(1) = 1$, i.e., outcome 1 is more desirable than outcome 0. In this example, we can compute $\pi^{(n)}$ in closed form. Specifically, we can see that $\pi^{(n)}(0) = \frac{1}{2^n}$ and $\pi^{(n)}(1) = 1 - \frac{1}{2^n}$. Thus,

$$D_{\text{KL}}(\pi^{(n)}|\pi_{\text{ref}}) = \log(2) - h\left(\frac{1}{2^n}\right), \quad (5)$$

where $h(\cdot)$ is the binary entropy function.³ We compare the exact closed-form expression for KL divergence with the analytical formula in Equation (4). As can be seen in Figure 1 (and is evident from Equation (5)), the true KL is upper bounded by $\log(2)$ for all n , whereas $\widetilde{\text{KL}}_n$ grows unbounded as $n \rightarrow \infty$. We also report a new estimator for KL divergence that closely mirrors the true KL divergence.

As we learnt from Example 1, the KL divergence between the best-of- n policy and the reference policy may be quite different from what the analytical formula used in the literature suggests. In the rest of this paper, we shed some light on this formula, derive bounds on the KL divergence, and propose a new estimator for the KL divergence that better captures the behavior of the KL divergence. We also theoretically reason about the win rate vs KL tradeoffs for the best-of- n policy, and justify its widespread use in language model alignment.

2. Derivation of the Best-of- n Policy

Our first step is to provide a derivation for the best-of- n policy under two simplifying assumptions. Let $r(\mathbf{x}, \mathbf{y}) \in \mathbb{R}$ represent the scalar reward of response \mathbf{y} in context \mathbf{x} .

³ $h(x) := -x \log(x) - (1-x) \log(1-x)$, for all $x \in (0, 1)$, and $h(0) = h(1) := 0$. Further, note that all logarithms in this paper are to the base e .

Assumption 2.1. We assume that the reward $r(\mathbf{x}, \mathbf{y})$ is unique for all \mathbf{x}, \mathbf{y} .

Assumption 2.2. Let $\mathcal{Y}^* := \{\mathbf{y} \mid \max_{\mathbf{x} \in \mathcal{X}} \pi_{\text{ref}}(\mathbf{y}|\mathbf{x}) > 0\}$. We assume that the language model is such that $|\mathcal{Y}^*| < \infty$, i.e., there are finite possible outcomes (in each context).

Note that Assumptions 2.1-2.2 are fairly non-restrictive and make the presentation of the results clearer.

The following result gives the probability mass function (PMF) of the best-of- n policy.

Lemma 2.3. Under Assumptions 2.1-2.2, for all $n \in \mathbb{N}$, the PMF of the best-of- n policy is given by

$$\pi^{(n)}(\mathbf{y}|\mathbf{x}) = \mathcal{F}_{\pi_{\text{ref}}}(\mathbf{y}|\mathbf{x})^n - \mathcal{F}_{\pi_{\text{ref}}}^-(\mathbf{y}|\mathbf{x})^n, \quad (6)$$

where for any distribution π ,

$$\mathcal{F}_{\pi}(\mathbf{y}|\mathbf{x}) := P_{\mathbf{z} \sim \pi(\cdot|\mathbf{x})}[r(\mathbf{x}, \mathbf{z}) \leq r(\mathbf{x}, \mathbf{y})], \quad (7)$$

$$\mathcal{F}_{\pi}^-(\mathbf{y}|\mathbf{x}) := P_{\mathbf{z} \sim \pi(\cdot|\mathbf{x})}[r(\mathbf{x}, \mathbf{z}) < r(\mathbf{x}, \mathbf{y})]. \quad (8)$$

Proof. Let $\mathcal{Y}_{\mathbf{x}}$ be the set of all possible outcomes of the language model, given prompt \mathbf{x} , i.e., $\mathcal{Y}_{\mathbf{x}} := \{\mathbf{y} \mid \pi_{\text{ref}}(\mathbf{y}|\mathbf{x}) > 0\}$. Further, let $L_{\mathbf{x}} := |\mathcal{Y}_{\mathbf{x}}| < \infty$ (Assumption 2.2). We order all possible $L_{\mathbf{x}}$ outcomes as $\{\tilde{\mathbf{y}}_i\}_{i \in [L_{\mathbf{x}}]}$ such that if $r(\mathbf{x}, \tilde{\mathbf{y}}_j) > r(\mathbf{x}, \tilde{\mathbf{y}}_i)$, then $j > i$. In other words, $\tilde{\mathbf{y}}_1$ is the least desirable outcome associated with the lowest reward, and $\tilde{\mathbf{y}}_{L_{\mathbf{x}}}$ is the most desirable outcome associated with the highest reward.

First notice that sampling from π_{ref} is equivalent to sampling $u \sim \mathcal{U}[0, 1]$, and returning $\tilde{\mathbf{y}}_i$, such that

$$\mathcal{F}_{\pi_{\text{ref}}}(\tilde{\mathbf{y}}_{i-1}|\mathbf{x}) \leq u < \mathcal{F}_{\pi_{\text{ref}}}(\tilde{\mathbf{y}}_i|\mathbf{x}). \quad (9)$$

Similarly, sampling from the best-of- n strategy is akin to sampling $u_1, \dots, u_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}[0, 1]$, and returning $\tilde{\mathbf{y}}_i$, such that

$$\mathcal{F}_{\pi_{\text{ref}}}(\tilde{\mathbf{y}}_{i-1}|\mathbf{x}) \leq \max_{k \in [n]} u_k < \mathcal{F}_{\pi_{\text{ref}}}(\tilde{\mathbf{y}}_i|\mathbf{x}). \quad (10)$$

On the other hand, we know that the CDF of the maximum of $u_1, \dots, u_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}[0, 1]$, for all $\tau \in [0, 1]$ is given by

$$P \left[\max_{k \in [n]} u_k \leq \tau \right] = \tau^n. \quad (11)$$

Hence, for all $n \in \mathbb{N}$, the PMF of the best-of- n policy, denoted as $\pi^{(n)}$ is given by

$$\pi^{(n)}(\tilde{\mathbf{y}}_i|\mathbf{x}) = \mathcal{F}_{\pi_{\text{ref}}}(\tilde{\mathbf{y}}_i|\mathbf{x})^n - \mathcal{F}_{\pi_{\text{ref}}}(\tilde{\mathbf{y}}_{i-1}|\mathbf{x})^n \quad \forall i \in [L_{\mathbf{x}}], \quad (12)$$

where $\mathcal{F}_{\pi_{\text{ref}}}(\tilde{\mathbf{y}}_0|\mathbf{x}) := 0$, and

$$\mathcal{F}_{\pi_{\text{ref}}}(\tilde{\mathbf{y}}_i|\mathbf{x}) = \sum_{l \in [i]} \pi_{\text{ref}}(\tilde{\mathbf{y}}_l|\mathbf{x}). \quad (13)$$

The proof is completed by noticing that $\mathcal{F}_{\pi_{\text{ref}}}(\tilde{\mathbf{y}}_{i-1}|\mathbf{x}) = \mathcal{F}_{\pi_{\text{ref}}}^-(\tilde{\mathbf{y}}_i|\mathbf{x})$. \square

Notice that if $n = 1$, then $\pi^{(1)}(\mathbf{y}|\mathbf{x}) = \pi_{\text{ref}}(\mathbf{y}|\mathbf{x})$. For any n , Lemma 2.3 gives a closed-form expression for $\pi^{(n)}(\mathbf{y}|\mathbf{x})$, which we will use subsequently to derive theoretical guarantees on the KL divergence and win rate of best-of- n .

We also remark that we can extend this PMF to $\pi^{(\tau)}$ for real $\tau \geq 1$. While it may not be immediately clear how to sample from this extension, it is used for best-of- n distillation (Gui et al., 2024; Amini et al., 2025; Sessa et al., 2025) and we also use it to give bounds in Section 7.

3. Relations Between the KL Divergence and the Analytical Formula

Our first result shows that the analytical formula is an upper bound on the (context-dependent) KL divergence. The proofs for this result and several subsequent results are relegated to Appendix A.

Theorem 3.1. For any $n \in \mathbb{N}$, and any \mathbf{x} , let \widetilde{KL}_n be defined in (4). Then,

$$D_{KL}(\pi^{(n)}(\cdot|\mathbf{x}) \parallel \pi_{\text{ref}}(\cdot|\mathbf{x})) \leq \widetilde{KL}_n = \log(n) - \frac{n-1}{n}.$$

Corollary 3.2. For any n , and any prompt distribution μ ,

$$D_{KL}^{\mu}(\pi^{(n)} \parallel \pi_{\text{ref}}) = E_{\mathbf{x} \sim \mu} D_{KL}(\pi^{(n)}(\cdot|\mathbf{x}) \parallel \pi_{\text{ref}}(\cdot|\mathbf{x})) \leq \widetilde{KL}_n.$$

Proof. This directly follows from Theorem 3.1. \square

Subsequent to this work, Mroueh (2024) has extended this result to a larger class of stochastic processes (with potentially continuous support such as diffusion models) through the application of the strong data processing inequality. In Appendix B, we also extend this result to derive bounds on the KL divergence of the blockwise best-of- n decoding (Mudgal et al., 2024), which generally allows to reach similar reward vs KL tradeoffs with 10x smaller n . In the rest of this section, we characterize the gap defined as follows:

$$G_{KL}^{(n)}(\mathbf{x}) := \widetilde{KL}_n - D_{KL}(\pi^{(n)}(\cdot|\mathbf{x}) \parallel \pi_{\text{ref}}(\cdot|\mathbf{x})) \geq 0. \quad (14)$$

3.1. Upper Bounds on the Gap

We need a definition to state the upper bound results.

Definition 3.3. A model π_{ref} is called δ -bound if $\pi_{\text{ref}}(\mathbf{y}|\mathbf{x}) \leq \delta$ for all $\mathbf{y} \in \mathcal{Y}^*$ and \mathbf{x} .

In particular, we are interested in characterizing the gap for a δ -bound model for a small δ , which is a model with all outcomes having small likelihoods. Next, we state our main upper bound.

Theorem 3.4. *The gap in Equation (14) is upper bounded by*

$$G_{KL}^{(n)}(\mathbf{x}) \leq 2n(n-1)e^{-H_2(\pi_{\text{ref}}|\mathbf{x})}, \quad (15)$$

where $H_2(\pi_{\text{ref}}|\mathbf{x})$ is the conditional Rényi entropy of order 2 of the language model given context \mathbf{x} , and $H_\alpha(\pi)$ for any distribution π is defined as

$$H_\alpha(\pi|\mathbf{x}) := \frac{1}{1-\alpha} \log \left(\sum_{\mathbf{y} \in \mathcal{Y}^*} (\pi(\mathbf{y}|\mathbf{x}))^\alpha \right). \quad (16)$$

Corollary 3.5. *Let $\pi_{\text{ref}}(\mathbf{y}|\mathbf{x}) \leq \delta$ for all $\mathbf{y} \in \mathcal{Y}^*$, i.e., π_{ref} is δ -bound. Then, the gap in Equation (14) is upper bounded by*

$$G_{KL}^{(n)}(\mathbf{x}) \leq 2n(n-1)\delta. \quad (17)$$

Proof. The proof follows by noticing that $H_2(\pi_{\text{ref}}|\mathbf{x}) \geq \log(1/\delta)$ and invoking Theorem 3.4. \square

Intuitively, if the model outcomes are fairly low probability, making it unlikely to get the same sample twice or more in the n outcomes for best-of- n , the analytical formula \widehat{KL}_n in (4) could be relatively accurate, and the gap is bounded above. In other words, if π_{ref} is a δ -bound model, and n is sufficiently small such that $n^2\delta \ll 1$, then

$$D_{KL}(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) \approx \log(n) - \frac{n-1}{n}. \quad (18)$$

This assumption is left implicit in the derivation of Hilton & Gao (2022) for the KL divergence of best-of- n .

3.2. Lower Bounds on the Gap

In this section, we characterize cases where the gap may be large. To this end, let us define

$$\varepsilon_n := \pi_{\text{ref}}(\mathbf{y}|\mathbf{x}), \quad \text{where } \mathbf{y} \sim \pi^{(n)}(\cdot|\mathbf{x}). \quad (19)$$

Note that ε_n is a random function of \mathbf{x} . In the limit as $n \rightarrow \infty$, we define

$$\varepsilon_\infty := \pi_{\text{ref}}(\mathbf{y}_{\max}(\mathbf{x})|\mathbf{x}), \quad (20)$$

where $\mathbf{y}_{\max}(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathcal{Y}^*} r(\mathbf{x}, \mathbf{y})$. Notice that ε_∞ is a deterministic function of \mathbf{x} .

Theorem 3.6. *Let $\varepsilon_\infty > 0$ be defined in Equation (20). For $n \in \mathbb{N}$, the gap between the analytical formula in Equation (4) and KL divergence is lower bounded by*

$$G_{KL}^{(n)}(\mathbf{x}) \geq (1 - (1 - \varepsilon_\infty)^n) \left(\log \frac{n\varepsilon_\infty}{1 - (1 - \varepsilon_\infty)^n} - \frac{n-1}{n} \right) - (n-1)(1 - \varepsilon_\infty)^n \log(1 - \varepsilon_\infty) > 0. \quad (21)$$

Corollary 3.7. *As $n \rightarrow \infty$, the gap is lower bounded by*

$$G_{KL}^{(n)}(\mathbf{x}) \geq \log(n\varepsilon_\infty) + o_n(\log n). \quad (22)$$

In particular, when $n\varepsilon_\infty \gg 1$, then the gap grows unbounded as we already observed in Example 1.

4. Proposed Estimator for KL Divergence

Motivated by the derivation of the best-of- n policy in Lemma 2.3, we propose a new estimator for the KL divergence. As a warm-up, first notice the following upper bound:

Lemma 4.1. *For any $n \in \mathbb{N}$ and any \mathbf{x} ,*

$$D_{KL}(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) \leq E_{\mathbf{y} \sim \pi^{(n)}} \left[\log \left(\frac{1 - (1 - \pi_{\text{ref}}(\mathbf{y}|\mathbf{x}))^n}{\pi_{\text{ref}}(\mathbf{y}|\mathbf{x})} \right) \right].$$

Therefore we may suggest to use the following *alternate estimator* for KL divergence

$$\widehat{D}_{KL, \text{loose}}(\varepsilon_n) := \log \left(\frac{1 - (1 - \varepsilon_n)^n}{\varepsilon_n} \right), \quad (23)$$

where ε_n is defined in Equation (19). Note that the expected value of $\widehat{D}_{KL, \text{loose}}(\varepsilon_n)$ is an upper bound on the KL divergence between the best-of- n policy and the reference policy. However, this estimator is loose by an additive constant of $(n-1)/n$, especially when $n\varepsilon_n \ll 1$.

Here we propose a different estimator to close this gap. To derive the estimator, first notice the following result.

Corollary 4.2. *Let*

$$d_n(\varepsilon) := (1 - \varepsilon)^n \left(\log n + (n-1) \log(1 - \varepsilon) - \frac{n-1}{n} \right) + (1 - (1 - \varepsilon)^n) \log \left(\frac{1 - (1 - \varepsilon)^n}{\varepsilon} \right). \quad (24)$$

Recall the definition of ε_∞ in Equation (19). Then,

$$D_{KL}(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) \leq d_n(\varepsilon_\infty).$$

Note that Corollary 4.2 could not be directly used to derive an estimator for the KL divergence because we do not observe ε_∞ when performing the best-of- n policy. Inspired by this result, and given that we can only observe ε_n , we put forth the following practical estimator on the KL divergence.

Definition 4.3. Let ε_n be defined in (19). Then, we propose the following estimator for the KL divergence of the best-of- n policy and the reference policy:

$$\widehat{D}_{KL}(\varepsilon_n) := d_n(\varepsilon_n). \quad (25)$$

Note that the estimator proposed in Definition 4.3 is a random variable that depends on ε_n . We conjecture that in expectation it provides an upper bound on the true KL divergence.

Conjecture 4.4. *Let ε_n be defined in (19). Then,*

$$D_{KL}(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) \leq E_{\mathbf{y} \sim \pi^{(n)}(\cdot|\mathbf{x})} \left[\widehat{D}_{KL}(\varepsilon_n) \right].$$

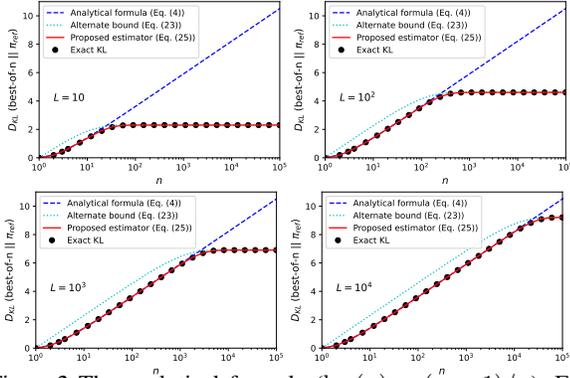


Figure 2. The analytical formula $(\log(n) - (n-1)/n)$, Equation (4), the alternate bound, Equation (23), the proposed estimator, Equation (25), and the exact KL divergence, for uniform distributions supported on alphabets of size $L = 10, 10^2, 10^3, 10^4$ respectively.

While we don't offer a mathematical proof for Conjecture 4.4, tens of thousands of randomly generated numerical experiments suggest that it holds true.

Let us further inspect the proposed estimator and its variance. We first show that it is strictly upper bounded by \widetilde{KL}_n .

Lemma 4.5. *For any realization of ε_n , we have*

$$0 \leq \widehat{D}_{KL}(\varepsilon_n) \leq \widetilde{KL}_n, \quad (26)$$

and hence

$$E_{\mathbf{y} \sim \pi^{(n)}(\cdot|\mathbf{x})} [\widehat{D}_{KL}(\varepsilon_n)] \leq \widetilde{KL}_n. \quad (27)$$

Given this, we can immediately bound the variance of the estimator too. Lemma 4.5 implies that standard deviation of the estimator is upper bounded by $\log n$, which in turn implies that if the estimator is averaged over $M = O(\log n \log \frac{1}{\delta})$ draws from the best-of- n model, the standard deviation is guaranteed to be smaller than δ . Given that we are generally interested in $n < 1000$, the dependence on n is mild. Having said that, given each of the M batches contains n iid samples (total of $M \times n$ iid samples), one should be able to build a bootstrapped estimator for the variance with better guarantees.

In what follows we numerically inspect the proposed estimator in a few scenarios, and compare it with the analytical formula and the exact KL divergence between the best-of- n policy and the reference policy.

The first set of examples, in Figure 2, are uniform distributions over alphabets of varying sizes. Notice that $\varepsilon_n = \varepsilon_\infty = \frac{1}{L}$ for a uniform distribution, and hence the estimator in Equation (25) and Equation (23) are deterministic. As can be seen KL divergence saturates around

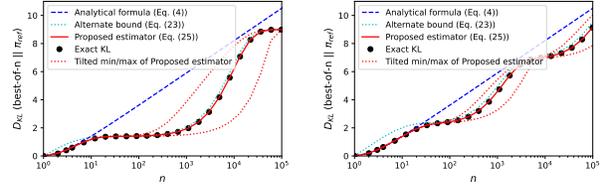


Figure 3. The analytical formula $(\log(n) - (n-1)/n)$, Equation (4), the alternate bound, Equation (23), the proposed estimator, Equation (25), and the exact KL divergence, for two cherry picked examples. In the left panel, the output is supported on an alphabet of size 5, where the highest reward outcome has a probability of 10^{-4} and the rest of the probability mass is uniformly distributed over the rest of the outcomes. In the right panel, the output is supported on an alphabet of size 200, where the three highest reward outcomes have probabilities $10^{-5}, 10^{-3}$ and 10^{-1} respectively. The rest of the probability mass is uniformly distributed over the rest of the outcomes.

$n \approx L$. For $\frac{n}{L} \ll 1$, the analytical formula of Equation (4), $\log(n) - (n-1)/n$, has a small gap with the actual KL divergence (which was also theoretically established in Corollary 3.5). On the other hand, when $\frac{n}{L} \gg 1$, the gap between $\log(n) - (n-1)/n$ and the actual KL divergence becomes large and unbounded (which was also theoretically established in Corollary 3.7). The alternate bound in Equation (23) captures the behavior of the KL divergence for $\frac{n}{L} \gg 1$ and has a finite gap. However, it has a gap of $(n-1)/n$ for $\frac{n}{L} \ll 1$ as previously discussed. Finally, we also observe that the proposed estimator in Equation (25) follows the behavior of the true KL divergence closely in all examples.

In the second set of examples, we cherry pick the probability mass function on the outcome to create different behaviors of the KL divergence, shown in Figure 3. In the left panel, the output is supported on an alphabet of size 5, where the highest reward outcome has a probability of 10^{-4} and the rest of the probability mass is uniformly distributed over the rest of the outcomes. We observe that KL divergence saturates early until the highest reward outcome is discovered with $n \approx 10^4$. In the right panel, the output is supported on an alphabet of size 200, where the highest reward outcome has a probability of 10^{-5} , the second highest reward outcome has a probability of 10^{-3} , and the third highest reward outcome has a probability of 10^{-1} . The rest of the probability mass is uniformly distributed over the rest of the outcomes. As can be seen, the KL divergence starts to saturate until the next high reward is outcome is discovered around $n \approx 10^3$ and $n \approx 10^5$. As can be seen, the analytical formula in Equation (4) does not capture the behavior of the KL divergence at all whereas the alternate bound in Equation (23) is much better aligned with the actual behavior. Finally, we observe that the proposed estimator in Equation (25) closely follows the actual KL divergence. We would like to recall that the proposed estimator is random

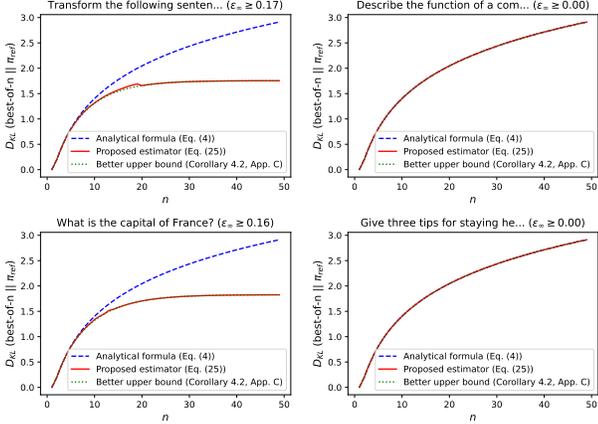


Figure 4. The analytical formula $(\log(n) - (n - 1)/n)$, Equation (4), better upper bound (Corollary 4.2, Appendix C), the proposed estimator, Equation (25), and the exact KL divergence, for four cherry picked examples from the Alpaca dataset (Taori et al., 2023) using Gemma 9B IT model (Gemma et al., 2024) with reward the log-likelihood of response under the reference model.

and we have plotted the expected value of the estimators, whereas in practice both estimators are subject to variance due to the randomness in the value of ε_n (Equation (19)). To capture the deviation from mean, we compute the *tilted min* and *tilted max* (see Appendix C.2 for the definition), which are also plotted in Figure 3. As can be seen, in cases where ε_n could widely vary based on whether a certain outcome appears in the set of n outcomes, the variance could be high.

In Figure 4, we compare the estimates for four cherry picked examples from the Alpaca dataset (Taori et al., 2023) using Gemma 9B IT model (Gemma et al., 2024) with reward being the log-likelihood of the reference model. Note that for two examples where ε_∞ is large, i.e., the prompts that induce less entropy in the response, such as “What is the capital of France?”, the proposed estimator outperforms the analytical formula in Equation (4) considerably and lies very close to the better upper bound in Corollary 4.2, whereas $\widetilde{\text{KL}}_n$ (Equation (4)) is loose even for $n \approx 20$. The details on the prompts used can be found in Appendix C. In Figure 5, we repeat the same experiment but change the reward to the negative of length to prefer more concise responses and see similar trends. We also include experiments with machine translation in Appendix C.3.

5. Win Rate of the Best-of- n Policy

So far, we provided theoretical guarantees on the KL divergence of the best-of- n policy with respect to the reference policy. In this section, we extend our study to characterize the *win rate* of the best-of- n policy against the reference policy. Let *win* of \mathbf{y} against \mathbf{z} in context \mathbf{x} be defined as:

$$w_r(\mathbf{y} \succ \mathbf{z} | \mathbf{x}) := \mathbf{1}(r(\mathbf{x}, \mathbf{y}) > r(\mathbf{x}, \mathbf{z})) + \frac{\mathbf{1}(r(\mathbf{x}, \mathbf{y}) = r(\mathbf{x}, \mathbf{z}))}{2}.$$

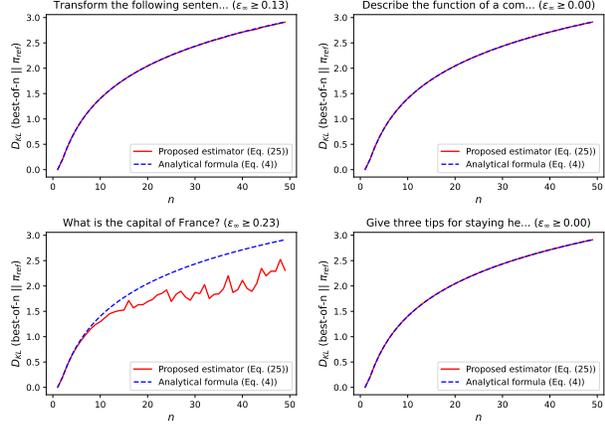


Figure 5. The analytical formula $(\log(n) - (n - 1)/n)$, Equation (4), better upper bound (Corollary 4.2, Appendix C), the proposed estimator, Equation (25), and the exact KL divergence, for four cherry picked examples from the Alpaca dataset (Taori et al., 2023) using Gemma 9B IT model (Gemma et al., 2024) with reward the negative length of the response.

In other words, $w_r(\mathbf{y} \succ \mathbf{z} | \mathbf{x})$ indicates whether response \mathbf{y} wins over response \mathbf{z} in context \mathbf{x} using the judge r . Then, let the *win rate* of policy π over the reference policy π_{ref} for prompt \mathbf{x} be defined as

$$\mathcal{W}_r(\pi(\cdot | \mathbf{x}) || \pi_{\text{ref}}(\cdot | \mathbf{x})) := E_{\mathbf{y} \sim \pi(\cdot | \mathbf{x})} E_{\mathbf{z} \sim \pi_{\text{ref}}(\cdot | \mathbf{x})} w_r(\mathbf{y} \succ \mathbf{z} | \mathbf{x}). \quad (28)$$

It is clear that $\mathcal{W}_r(\pi_{\text{ref}}(\cdot | \mathbf{x}) || \pi_{\text{ref}}(\cdot | \mathbf{x})) = 0.5$ and the goal of alignment is to improve win rate beyond 0.5 (with the lowest KL divergence between the two models). We further define the following, averaged over prompt distribution μ :

$$\mathcal{W}_r^\mu(\pi || \pi_{\text{ref}}) := E_{\mathbf{x} \sim \mu} \mathcal{W}_r(\pi(\cdot | \mathbf{x}) || \pi_{\text{ref}}(\cdot | \mathbf{x})). \quad (29)$$

Note that it is clear that if $\mathbf{y} \sim \pi(\cdot | \mathbf{x})$ and $\mathbf{z} \sim \pi_{\text{ref}}(\cdot | \mathbf{x})$, then $w_r(\mathbf{y} \succ \mathbf{z} | \mathbf{x})$ is an unbiased estimator for win rate of policy π against π_{ref} .

We analyze the win rate and derive theoretical guarantees. Let \mathcal{F}_π and \mathcal{F}_π^- be defined in Equation (7) and Equation (8), respectively. We define *calibrated reward* as:

$$\mathcal{C}_{\pi_{\text{ref}}}(\mathbf{x}, \mathbf{y}) := \frac{\mathcal{F}_{\pi_{\text{ref}}}(\mathbf{y} | \mathbf{x}) + \mathcal{F}_{\pi_{\text{ref}}}^-(\mathbf{y} | \mathbf{x})}{2}. \quad (30)$$

With this definition in place, notice that win rate could be expressed as follows.

Lemma 5.1. *The win rate of any policy π against reference policy π_{ref} could be expressed as*

$$\mathcal{W}_r(\pi(\cdot | \mathbf{x}) || \pi_{\text{ref}}(\cdot | \mathbf{x})) = E_{\mathbf{y} \sim \pi(\cdot | \mathbf{x})} [\mathcal{C}_{\pi_{\text{ref}}}(\mathbf{x}, \mathbf{y})]. \quad (31)$$

Notice that the above result suggests that the win rate of any policy only depends on reward and the reference policy through the calibrated reward function, $\mathcal{C}_{\pi_{\text{ref}}}(\mathbf{x}, \mathbf{y})$. Hence, this notion of calibration may be used as a canonical transformation of the reward for preference optimization against

a given reference policy. In fact, this transformation is theoretically proposed as the objective in IPO (Azar et al., 2024) and is the key to best-of- n distillation (Gui et al., 2024; Amini et al., 2025; Sessa et al., 2025; Yang et al., 2024b) and inference-aware alignment (Balashankar et al., 2025).

Lemma 5.2. *The win rate of best-of- n policy against π_{ref} is given by*

$$\begin{aligned} & \mathcal{W}_r(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) \\ &= \sum_{\mathbf{y}} \left(\mathcal{F}_{\pi_{\text{ref}}}(\mathbf{y}|\mathbf{x})^n - \mathcal{F}_{\pi_{\text{ref}}}^-(\mathbf{y}|\mathbf{x})^n \right) C_{\pi_{\text{ref}}}(\mathbf{x}, \mathbf{y}). \end{aligned}$$

Proof. This is proved by plugging Lemma 2.3 into Lemma 5.1. \square

Our next result is an upper bound on the win rate of best-of- n policy. Intuitively, observe that the win rate could be estimated by drawing $(n+1)$ samples from π_{ref} , associating n to the best-of- n model and the remaining one to the reference model. Hence, the best-of- n model gets n -to-1 chances of winning against the reference model, unless there is a draw due to samples with the same reward value. Thus, intuitively $\mathcal{W}_r(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) \approx \frac{n}{n+1}$. We formalize this as an upper bound on the win rate.

Theorem 5.3. *For all, n , and all \mathbf{x} , the win rate of best-of- n policy is upper bounded by*

$$\mathcal{W}_r(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) \leq \frac{n}{n+1}. \quad (32)$$

In the rest of this section, we derive bounds on the gap between this upper bound and the actual win rate:

$$G_{\mathcal{W}}^{(n)}(\mathbf{x}) := \frac{n}{n+1} - \mathcal{W}_r(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) \geq 0. \quad (33)$$

5.1. Upper Bounds on the Win Rate Gap

Unlike the KL divergence, it is clear that $G_{\mathcal{W}}^{(n)}(\mathbf{x})$ could not grow unbounded, and is upper bounded by $\frac{1}{2}$.

Theorem 5.4. *The win rate gap is upper bounded by*

$$G_{\mathcal{W}}^{(n)}(\mathbf{x}) \leq \frac{n-1}{2} e^{-H_2(\pi_{\text{ref}}|\mathbf{x})}, \quad (34)$$

where $H_2(\cdot)$ denotes the Rényi entropy of order 2 defined in Equation (16).

Corollary 5.5. *Let $\pi_{\text{ref}}(\mathbf{y}|\mathbf{x}) \leq \delta$ for all $\mathbf{y} \in \mathcal{Y}^*$, i.e., π_{ref} is δ -bound. Then,*

$$G_{\mathcal{W}}^{(n)}(\mathbf{x}) \leq \frac{n-1}{2} \delta. \quad (35)$$

Proof. The proof follows by noticing that $H_2(\pi_{\text{ref}}|\mathbf{x}) \geq \log(1/\delta)$ and invoking Theorem 5.4. \square

Given this upper bound, the win rate of best-of- n would be fairly close to $\frac{n}{n+1}$ if π_{ref} is a δ -bound model, and n is sufficiently small such that $n\delta \ll 1$, by combining Theorem 5.3 and Corollary 5.5, we get

$$\mathcal{W}_r(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) \approx \frac{n}{n+1}, \quad (36)$$

which is the claim of Gui et al. (2024, Theorem 2) for the win rate of best-of- n .

5.2. Lower Bounds on the Win Rate Gap

Our next result characterizes the cases where the gap is bounded away from 0.

Theorem 5.6. *Let $\varepsilon_{\infty} > 0$ be defined in Equation (20). Then,*

$$\begin{aligned} G_{\mathcal{W}}^{(n)}(\mathbf{x}) &\geq \frac{n}{n+1} (1 - (1 - \varepsilon_{\infty})^{n+1}) \\ &\quad - (1 - (1 - \varepsilon_{\infty})^n) \left(1 - \frac{\varepsilon_{\infty}}{2}\right) > 0. \end{aligned}$$

Corollary 5.7. *As $n \rightarrow \infty$, we have*

$$G_{\mathcal{W}}^{(n)}(\mathbf{x}) \geq \frac{\varepsilon_{\infty}}{2} (1 + o_n(1)). \quad (37)$$

As $n \rightarrow \infty$, $G_{\mathcal{W}}^{(n)}(\mathbf{x}) \rightarrow \frac{\varepsilon_{\infty}}{2}$, which is bounded from 0.

6. Rewind-and-Repeat: Rejection Sampling Beyond Best-of- n

The best-of- n policy is a form of rejection sampling. Another form is called rewind-and-repeat, where the process of generating a response and scoring it is repeated until a certain threshold on reward is met (Kim et al., 2025). A more involved blockwise variant of this process is recently used by Li et al. (2024). Formally, let \mathbf{x} be a given input prompt to the model, and let $\{\mathbf{y}_k\}_{k=1}^{\infty}$ be a sequence of infinite i.i.d. samples drawn from $\pi_{\text{ref}}(\cdot|\mathbf{x})$. Then, rewind-and-repeat accepts \mathbf{y}_M such that

$$r(\mathbf{x}, \mathbf{y}_M) \geq \Phi \quad \text{and} \quad \forall k < M : r(\mathbf{x}, \mathbf{y}_k) < \Phi, \quad (38)$$

where $\Phi \in \mathbb{R}$ is the threshold on reward. In other words, \mathbf{y}_M is the first draw whose reward reaches a certain threshold Φ . We also call M the (random) number of trials until the threshold is met, which determines the cost of inference from the model. We denote the resulting policy by π_{Φ} .

It is natural to ask: *how do the win rate vs KL tradeoffs of rewind-and-repeat compare with that of best-of- n ?* To answer this question, first we define

$$w_{\Phi}(\mathbf{x}) := E_{\mathbf{y} \sim \pi_{\text{ref}}(\cdot|\mathbf{x})}[\mathbf{1}(r(\mathbf{x}, \mathbf{y}) \geq \Phi)] \quad (39)$$

as the probability of drawing a sample from the reference policy that meets the threshold. Hence, the expected number of trials to output an outcome is $E[M] = 1/w_{\Phi}(\mathbf{x})$.

Next, let us derive the PMF of rewind-and-repeat policy.

Lemma 6.1. *The probability mass function (PMF) of the rewind-and-repeat policy is given by*

$$\pi_{\Phi}(\mathbf{y}|\mathbf{x}) = \begin{cases} \frac{\pi_{\text{ref}}(\mathbf{y}|\mathbf{x})}{w_{\Phi}(\mathbf{x})} & \text{if } r(\mathbf{x}, \mathbf{y}) \geq \Phi \\ 0 & \text{if } r(\mathbf{x}, \mathbf{y}) < \Phi \end{cases}. \quad (40)$$

The proofs are deferred to Appendix A.4. Given its PMF, next we derive KL divergence and win rate of the rewind-and-repeat policy and the reference policy.

Lemma 6.2. *We have*

$$D_{KL}(\pi_{\Phi}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) = \log \frac{1}{w_{\Phi}(\mathbf{x})}, \quad (41)$$

$$\mathcal{W}_r(\pi_{\Phi}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) = 1 - \frac{1}{2}w_{\Phi}(\mathbf{x}). \quad (42)$$

Thus, by sweeping Φ in \mathbb{R} , we will effectively sweep w_{Φ} in $[0, 1]$, and obtain the respective win rate vs KL divergence tradeoff for the rewind-and-repeat policy. Note that the KL divergence was recently derived by Kim et al. (2025, Appendix A.5) and we provide a proof for completeness.

So far, we derived a characterization of the KL divergence of the rewind-and-repeat policy. However, when the number of outcomes of a model is large, similarly to the case of best-of- n , estimating the KL divergence is intractable.

Our main result in this section is an unbiased estimator of the KL divergence of the rewind-and-repeat procedure.

Theorem 6.3. *For $n \geq 1$, let $H_n := \sum_{i=1}^n \frac{1}{i}$ be the n -th Harmonic number and $H_0 = 0$. Further, let M be the number of trials to achieve an outcome in the rewind-and-repeat policy. Then,*

$$D_{KL}(\pi_{\Phi}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) = E[H_{M-1}]. \quad (43)$$

Hence, we propose H_{M-1} as an (unbiased) estimator for the KL divergence of the rewind-and-repeat procedure with respect to the reference policy.

7. Win rate vs KL Divergence Tradeoffs

Thus far, we characterized the KL divergence and win rate of best-of- n . In practice, it is customary to compare different alignment methods based on their win rate at a certain KL divergence from the reference policy. Note that Theorem 3.1 implies that the win rate (or expected reward) vs KL tradeoffs reported in the literature that use the analytical formula in Equation (4) (Gao et al., 2023; Go et al., 2023; Mudgal et al., 2024; Scheurer et al., 2023) are conservative and the actual tradeoff curve of the best-of- n policy is in fact guaranteed to be no worse than what is reported. To further substantiate this point, let us revisit Example 1 and report the win rate vs KL divergence tradeoff curve (Figure 6), where we used the actual win rate in all cases.⁴ The actual

⁴In practice, the win rate could be estimated using the unbiased estimator given by the win random variable in Equation (28).

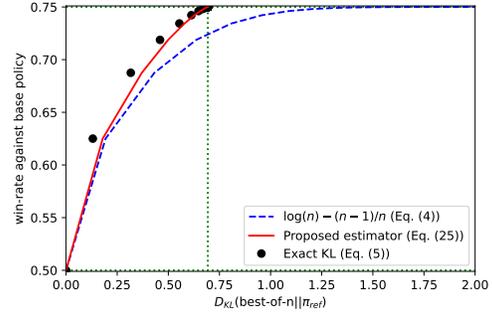


Figure 6. The win rate against reference policy vs KL divergence tradeoff curve for the best-of- n policy. We have used the analytical formula $\log(n) - (n - 1)/n$, Equation (4), the exact KL divergence, Equation (5), and the proposed estimator, Equation (25), for producing the tradeoffs from Example 1, illustrating a case where the actual win rate vs KL divergence tradeoff curve for the best-of- n policy is more favorable than the one predicted if using the upper bound formula, $\log(n) - (n - 1)/n$.

win rate vs KL divergence tradeoff is more favorable than that portrayed by using the formula in Equation (4). In this example, the best-of- n policy in the limit of large n , reaches a KL divergence of $\log(2)$ and a win rate of 0.75.

Definition 7.1. Let $W : \mathbb{R}^+ \rightarrow [0, 1]$ be a function that takes in $D \geq 0$ and outputs $W(D)$. Let $D_n = D_{KL}(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x}))$. We say that W is an upper bound on the tradeoff curve of best-of- n if for all n , $\mathcal{W}_r(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) \leq W(D_n)$. Alternatively, we say that W is a lower bound on the tradeoff curve of best-of- n if for all n , $\mathcal{W}_r(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) \geq W(D_n)$.

We can turn our existing bounds into straightforward lower and upper bounds on the tradeoff curve for best-of- n . (see Appendix A.5). We conjecture a tighter upper bound.

Conjecture 7.2. *For any \mathbf{x} and a given π_{ref} , let the function $W_0(D) = \ell^{-1}(D)$ where for all $\tau \in [0.5, 1)$,*

$$\ell(\tau) = \log \frac{\tau}{1-\tau} + \frac{1}{\tau} - 2.$$

W_0 is an upper bound on the tradeoff curve of best-of- n .

Example 2. We consider a ternary language model with alphabet $\mathcal{X} = \{0, 1, 2\}$, ordered from least preferred to most preferred, with probabilities given by $\pi_{\text{ref}} = (0.3, 0.6, 0.1)$. Hence, the calibrated reward in Equation (30) is given by $(0.3, 0.75, 0.95)$. The set of solutions to the KL-regularized RL problem are given by Equation (2). We also compute the best-of- n solutions (for continuous n) using Lemma 2.3 and rewind-and-repeat using Lemma 6.1. Before discussing the win rate vs KL divergence tradeoffs, we first visualize the set of solutions on the probability simplex in Figure 7 and observe that the solutions could be very different. When we consider the win rate vs KL tradeoffs for this example (Figure 8), we observe that the tradeoffs are strikingly close

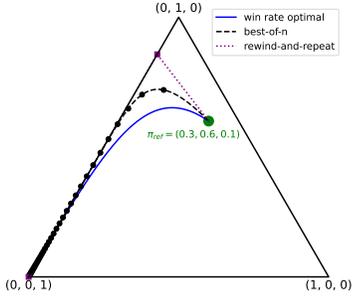


Figure 7. The set of solutions for Example 2: *win rate optimal* is achieved by varying β in Equation (2), *best-of- n* is given by Lemma 2.3, and *rewind-and-repeat* is given by Lemma 6.1.

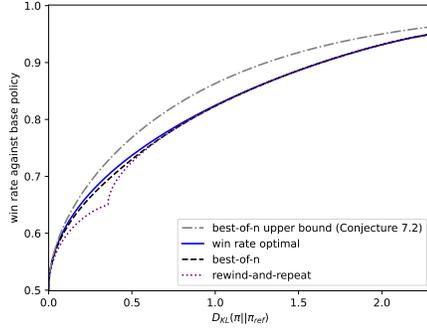


Figure 8. Win rate vs KL tradeoff for Example 2. The tradeoff curve of *Best-of- n* is close to *win rate optimal*, and both are better than *rewind-and-repeat*.

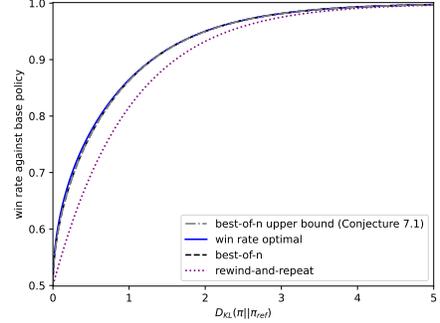


Figure 9. Win rate vs KL tradeoff for Example 3. The tradeoff curve of *Best-of- n* is close to *win rate optimal* and the limit behavior, and both are better than *rewind-and-repeat*.

for best-of- n and win rate optimal models, matching recent findings of Yang et al. (2024a); Gui et al. (2024).

Example 3. We consider a language model with low probability outcomes, and a calibrated reward, and plot win rate vs KL divergence tradeoffs in Figure 9, and observe that in this asymptotic regime where we observe that the tradeoff curve offered by best-of- n is almost optimal and offers better tradeoffs compared to rewind-and-repeat.

8. Conclusion

We studied the best-of- n alignment policy and derived its probability mass function (Lemma 2.3). We proved that an analytical formula used in the literature for the KL divergence of the best-of- n policy with the reference policy, $\log(n) - (n-1)/n$ in Equation (4), is false, and only an upper bound on the KL divergence (Theorem 3.1). We derived bounds on the gap between this formula and the KL divergence where we roughly showed the following: Let \mathbf{y} be a draw from the best-of- n policy. Let ε_n be the probability mass of \mathbf{y} under the base model. Then, if $n\varepsilon_n \ll 1$, the gap between the formula in Equation (4) and the exact KL divergence is small (Theorem 3.4); and if $n\varepsilon_n \gg 1$, the gap between the two may be large and unbounded (Theorem 3.6). We proposed a new estimator for the KL divergence (Definition 4.3), which we demonstrated to capture the behavior of the KL divergence on several numerical experiments. We showed that the win rate of best-of- n against the reference policy is upper bounded by $n/(n+1)$ (Theorem 5.3). Similarly to the KL divergence, we provided upper (Theorem 5.4) and lower (Theorem 5.6) bounds on the gap between the actual win rate and the bound. We compared best-of- n with another form of rejection sampling through rewind-and-repeat and showed its superiority both theoretically and empirically. We also extended the bounds to blockwise best-of- n (Mudgal et al., 2024).

While our results showed that best-of- n offers better trade-

offs on reward vs KL divergence (where KL divergence captures preservation of the core model capabilities) compared to rewind-and-repeat, the latter is more effective in driving reward up for a given compute budget which is important in test-time compute scaling. Hence, it remains to be seen how to best design a method that achieves Pareto optimal tradeoffs between compute, reward, and KL divergence (which captures preservation of capabilities other than what is captured by reward). This might involve combining the rewind-and-repeat with best-of- n and could be an area for future work.

Impact Statement

This paper presents theoretical investigations of best-of- n sampling and other test-time rejection sampling algorithms, which is a simple yet effective method for test-time alignment of generative models. One of the major findings of this paper is that a widely used formula for KL divergence of the best-of- n policy and the reference policy is a theoretical upper bound on the actual KL divergence. This may help ensure that the capabilities of the reference model are largely preserved in the aligned model. For example, this may help compliance or risk-management teams preserve safety by guaranteeing that the policy that is served does not drift too far from a safe reference policy.

Our work also showed that best-of- n is an effective (and almost optimal) test-time alignment method which comes with theoretical guarantees on win rate vs KL divergence tradeoffs motivating its use for improving the capabilities and safety of models. On the flip side, this also shows why best-of- n with an adversarial reward may be used to effectively jailbreak generative models. We hope future work can benefit from our findings in making best-of- n more effective, and can also devise safeguards against best-of- n jailbreaks.

References

- Amini, A., Vieira, T., and Cotterell, R. Variational best-of-n alignment. *International Conference on Learning Representations (ICLR)*, 2025.
- Azar, M. G., Guo, Z. D., Piot, B., Munos, R., Rowland, M., Valko, M., and Calandriello, D. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, pp. 4447–4455. PMLR, 2024.
- Bai, Y., Jones, A., Ndousse, K., Askell, A., Chen, A., DasSarma, N., Drain, D., Fort, S., Ganguli, D., Henighan, T., et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- Balashankar, A., Sun, Z., Berant, J., Eisenstein, J., Collins, M., Hutter, A., Lee, J., Nagpal, C., Prost, F., Sinha, A., Suresh, A. T., and Beirami, A. InfAlign: Inference-aware language model alignment. *International Conference on Machine Learning (ICML)*, 2025.
- Beetham, J., Chakraborty, S., Wang, M., Huang, F., Bedi, A. S., and Shah, M. Liar: Leveraging alignment (best-of-n) to jailbreak llms in seconds. *arXiv preprint arXiv:2412.05232*, 2024.
- Brown, B., Juravsky, J., Ehrlich, R., Clark, R., Le, Q. V., Ré, C., and Mirhoseini, A. Large language monkeys: Scaling inference compute with repeated sampling. *arXiv preprint arXiv:2407.21787*, 2024.
- Christiano, P. F., Leike, J., Brown, T., Martic, M., Legg, S., and Amodei, D. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- Coste, T., Anwar, U., Kirk, R., and Krueger, D. Reward model ensembles help mitigate overoptimization. *International Conference on Learning Representations (ICLR)*, 2024.
- Eisenstein, J., Nagpal, C., Agarwal, A., Beirami, A., D’Amour, A., Dvijotham, D., Fisch, A., Heller, K., Pfohl, S., Ramachandran, D., Shaw, P., and Berant, J. Helping or herding? reward model ensembles mitigate but do not eliminate reward hacking. *Conference on Language Modeling (COLM)*, 2024.
- Gao, L., Schulman, J., and Hilton, J. Scaling laws for reward model overoptimization. In *International Conference on Machine Learning*, pp. 10835–10866. PMLR, 2023.
- Gemma, Riviere, M., Pathak, S., Sessa, P. G., Hardin, C., Bhupatiraju, S., Hussenot, L., Mesnard, T., Shahriari, B., Ramé, A., et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- Go, D., Korbak, T., Kruszewski, G., Rozen, J., and Dymetman, M. Compositional preference models for aligning LMs. *arXiv preprint arXiv:2310.13011*, 2023.
- Gui, L., Gârbasea, C., and Veitch, V. BoNBoN alignment for large language models and the sweetness of best-of-n sampling. *Neural Information Processing Systems (NeurIPS)*, December 2024.
- Hilton, J. and Gao, L. Measuring Goodhart’s law, April 2022. URL <https://openai.com/research/measuring-goodharts-law>. Accessed: 2024-01-03.
- Hughes, J., Price, S., Lynch, A., Schaeffer, R., Barez, F., Koyejo, S., Sleight, H., Jones, E., Perez, E., and Sharma, M. Best-of-n jailbreaking. *arXiv preprint arXiv:2412.03556*, 2024.
- Kim, M., Thonet, T., Rozen, J., Lee, H., Jung, K., and Dymetman, M. Guaranteed generation from large language models. *International Conference on Learning Representations (ICLR)*, 2025.
- Korbak, T., Elsahar, H., Kruszewski, G., and Dymetman, M. 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.
- Korbak, T., Perez, E., and Buckley, C. RL with KL penalties is better viewed as Bayesian inference. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 1083–1091, 2022b.
- Li, T., Beirami, A., Sanjabi, M., and Smith, V. On tilted losses in machine learning: Theory and applications. *Journal of Machine Learning Research*, 24(142):1–79, 2023.
- Li, Y., Wei, F., Zhao, J., Zhang, C., and Zhang, H. Rain: Your language models can align themselves without fine-tuning. In *The Twelfth International Conference on Learning Representations*, 2024.
- Mroueh, Y. Information theoretic guarantees for policy alignment in large language models. *arXiv preprint arXiv:2406.05883*, 2024.
- Mudgal, S., Lee, J., Ganapathy, H., Li, Y., Wang, T., Huang, Y., Chen, Z., Cheng, H.-T., Collins, M., Strohmaier, T., Chen, J., Beutel, A., and Beirami, A. Controlled decoding from language models. *International Conference on Machine Learning (ICML)*, 2024.

- Nakano, R., Hilton, J., Balaji, S., Wu, J., Ouyang, L., Kim, C., Hesse, C., Jain, S., Kosaraju, V., Saunders, W., et al. WebGPT: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- Qiu, J., Lu, Y., Zeng, Y., Guo, J., Geng, J., Wang, H., Huang, K., Wu, Y., and Wang, M. Treebon: Enhancing inference-time alignment with speculative tree-search and best-of-n sampling. *arXiv preprint arXiv:2410.16033*, 2024.
- Rafailov, R., Sharma, A., Mitchell, E., Ermon, S., Manning, C. D., and Finn, C. Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.18290*, 2023.
- Scheurer, J., Campos, J. A., Korbak, T., Chan, J. S., Chen, A., Cho, K., and Perez, E. Training language models with language feedback at scale. *arXiv preprint arXiv:2303.16755*, 2023.
- Sessa, P. G., Dadashi, R., Hussenot, L., Ferret, J., Vieillard, N., Ramé, A., Shariari, B., Perrin, S., Friesen, A., Cideron, G., et al. Bond: Aligning llms with best-of-n distillation. *International Conference on Learning Representations (ICLR)*, 2025.
- Snell, C., Lee, J., Xu, K., and Kumar, A. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- Stiennon, N., Ouyang, L., Wu, J., Ziegler, D., Lowe, R., Voss, C., Radford, A., Amodei, D., and Christiano, P. F. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33: 3008–3021, 2020.
- Sun, H., Haider, M., Zhang, R., Yang, H., Qiu, J., Yin, M., Wang, M., Bartlett, P., and Zanette, A. Fast best-of-n decoding via speculative rejection. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Taori, R., Gulrajani, I., Zhang, T., Dubois, Y., Li, X., Guestrin, C., Liang, P., and Hashimoto, T. B. Stanford alpaca: An instruction-following llama model, 2023.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Yang, J. Q., Salamatian, S., Sun, Z., Suresh, A. T., and Beirami, A. Asymptotics of language model alignment. *International Symposium on Information Theory (ISIT)*, July 2024a.
- Yang, K. and Klein, D. FUDGE: Controlled text generation with future discriminators. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 3511–3535, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.276. URL <https://aclanthology.org/2021.naacl-main.276>.
- Yang, T., Mei, J., Dai, H., Wen, Z., Cen, S., Schuurmans, D., Chi, Y., and Dai, B. Faster wind: Accelerating iterative best-of-n distillation for llm alignment. *arXiv preprint arXiv:2410.20727*, 2024b.
- Zhao, Y., Khalman, M., Joshi, R., Narayan, S., Saleh, M., and Liu, P. J. Calibrating sequence likelihood improves conditional language generation. In *The Eleventh International Conference on Learning Representations*, 2022.

A. Proofs of the main results of the paper

A.1. Proofs of Section 3

We provide the proofs of the main results of the paper. To do so, we need to state two further auxiliary lemmas.

Lemma A.1. For any $0 \leq a < b \leq 1$, and $n \in \mathbb{N}$,

$$(b^n - a^n) \log \frac{b^n - a^n}{b - a} = \int_a^b nv^{n-1} \log(nv^{n-1}) dv - g_n(a, b) \leq \int_a^b nv^{n-1} \log(nv^{n-1}) dv, \quad (44)$$

where $g(a, b) \geq 0$, and is given by

$$g_n(a, b) := (b^n - a^n) D_{\text{KL}}(p_v \| p_u) = (b^n - a^n) \log \frac{n(b-a)}{b^n - a^n} + (n-1)(b^n \log b - a^n \log a) - \frac{n-1}{n}(b^n - a^n), \quad (45)$$

where for $a \leq w \leq b$, we define $p_v(w) = \frac{nw^{n-1}}{b^n - a^n}$ and $p_u(w) = \frac{1}{b-a}$.

Proof. Notice that

$$D_{\text{KL}}(p_v \| p_u) = \int_a^b \frac{nv^{n-1}}{b^n - a^n} \log \frac{\frac{nv^{n-1}}{b^n - a^n}}{\frac{1}{b-a}} dv = \int_a^b \frac{nv^{n-1}}{b^n - a^n} \log \frac{nv^{n-1}(b-a)}{b^n - a^n} dv \geq 0. \quad (46)$$

The inequality is established by algebraic manipulation. The gap $g_n(a, b)$ is obtained by following Lemma A.2 to derive the right-hand-side of the inequality in closed form. \square

Lemma A.2. The following identity holds:

$$\int_a^b nv^{n-1} \log(nv^{n-1}) dv = (b^n - a^n) \log n + (n-1)(b^n \log b - a^n \log a) - \frac{n-1}{n}(b^n - a^n). \quad (47)$$

Proof. The proof is completed by noticing that $\frac{d}{dv}(v^n \log v) = nv^{n-1} \log v + v^{n-1}$. \square

Proof of Theorem 3.1.

$$D_{\text{KL}}(\pi^{(n)}(\cdot | \mathbf{x}) \| \pi_{\text{ref}}(\cdot | \mathbf{x})) = \sum_{\mathbf{y} \in \mathcal{Y}^*} (\mathcal{F}_{\pi_{\text{ref}}}(\mathbf{y} | \mathbf{x})^n - \mathcal{F}_{\pi_{\text{ref}}}^-(\mathbf{y} | \mathbf{x})^n) \log \frac{\mathcal{F}_{\pi_{\text{ref}}}(\mathbf{y} | \mathbf{x})^n - \mathcal{F}_{\pi_{\text{ref}}}^-(\mathbf{y} | \mathbf{x})^n}{\pi_{\text{ref}}(\mathbf{y} | \mathbf{x})} \quad (48)$$

$$\leq \sum_{\mathbf{y} \in \mathcal{Y}^*} \int_{\mathcal{F}_{\pi_{\text{ref}}}^-(\mathbf{y} | \mathbf{x})}^{\mathcal{F}_{\pi_{\text{ref}}}(\mathbf{y} | \mathbf{x})} nv^{n-1} \log(nv^{n-1}) dv \quad (49)$$

$$= \int_0^1 nv^{n-1} \log(nv^{n-1}) dv \quad (50)$$

$$= \log(n) - \frac{n-1}{n}, \quad (51)$$

where Equation (49) follows from Lemma A.1, and Equation (51) follows from Lemma A.2. \square

Lemma A.3. For any $0 \leq a < b \leq 1$, and $n \in \mathbb{N}$,

$$g_n(a, b) \leq 2n(n-1)(b-a)^2. \quad (52)$$

Proof. Recall from Equation (45) that

$$g_n(a, b) := (b^n - a^n)D_{\text{KL}}(p_v \| p_u), \quad (53)$$

where $p_v(w) = \frac{nw^{n-1}}{b^n - a^n}$ and $p_u(w) = \frac{1}{b-a}$ for $a \leq w \leq b$. We can further bound the KL divergence as

$$D_{\text{KL}}(p_v \| p_u) \leq \max_{w \in [a, b]} \log \frac{p_v(w)}{p_u(w)} \quad (54)$$

$$= \max_{w \in [a, b]} \log \frac{nw^{n-1}(b-a)}{b^n - a^n} \quad (55)$$

$$= \log \frac{nb^{n-1}(b-a)}{b^n - a^n} \quad (56)$$

$$= \log \frac{nb^{n-1}}{\sum_{j=0}^{n-1} b^{n-1-j} a^j}. \quad (57)$$

At this point, we will divide the proof into two cases depending on the value of a, b .

Case I. We first focus on the case when $a < b/2$. In this case,

$$D_{\text{KL}}(p_v \| p_u) \leq \log \frac{nb^{n-1}}{\sum_{j=0}^{n-1} b^{n-1-j} a^j} \quad (58)$$

$$\leq \log n. \quad (59)$$

Hence,

$$g_n(a, b) \leq (b^n - a^n) \log n \quad (60)$$

$$\leq b^n \log n \quad (61)$$

$$\leq b^2 \log n \quad (62)$$

$$\leq 4(b-a)^2 \log n \quad (63)$$

$$\leq 2n(n-1)(b-a)^2. \quad (64)$$

Case II. For $a \geq b/2$, notice that

$$D_{\text{KL}}(p_v \| p_u) \leq \log \frac{nb^{n-1}}{\sum_{j=0}^{n-1} b^{n-1-j} a^j} \quad (65)$$

$$\leq \log \frac{b^{n-1}}{b^{(n-1)/2} a^{(n-1)/2}} \quad (66)$$

$$= \frac{n-1}{2} \log \frac{b}{a} \quad (67)$$

$$\leq \frac{n-1}{2} \frac{(b-a)}{a} \quad (68)$$

$$\leq (n-1) \frac{(b-a)}{b}, \quad (69)$$

where the first inequality follows from AM-GM inequality. Hence,

$$g_n(a, b) \leq (b^n - a^n)(n-1) \frac{(b-a)}{b} \quad (70)$$

$$= (b-a) \left(\sum_{j=0}^{n-1} b^{n-1-j} a^j \right) (n-1) \frac{(b-a)}{b} \quad (71)$$

$$\leq nb^{n-1}(b-a)(n-1) \frac{(b-a)}{b} \quad (72)$$

$$= nb^{n-2}(b-a)(n-1)(b-a) \quad (73)$$

$$\leq 2n(n-1)(b-a)^2. \quad (74)$$

The proof is completed by putting together *Case I* and *Case II*. \square

Proof of Theorem 3.4.

$$\log(n) - \frac{n-1}{n} - D_{\text{KL}}(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) = \sum_{\mathbf{y} \in \mathcal{Y}^*} g_n(\mathcal{F}_{\pi_{\text{ref}}}^-(\mathbf{y}|\mathbf{x}), \mathcal{F}_{\pi_{\text{ref}}}(\mathbf{y}|\mathbf{x})) \quad (75)$$

$$\leq 2n(n-1) \sum_{\mathbf{y} \in \mathcal{Y}^*} (\mathcal{F}_{\pi_{\text{ref}}}(\mathbf{y}|\mathbf{x}) - \mathcal{F}_{\pi_{\text{ref}}}^-(\mathbf{y}|\mathbf{x}))^2 \quad (76)$$

$$= 2n(n-1) \sum_{\mathbf{y} \in \mathcal{Y}^*} (\pi_{\text{ref}}(\mathbf{y}|\mathbf{x}))^2 \quad (77)$$

$$= 2n(n-1)e^{-H_2(\pi_{\text{ref}}|\mathbf{x})}. \quad (78)$$

where Equation (76) follows from Lemma A.3. \square

Proof of Theorem 3.6. Notice that the gap is at least lower bounded by the value of the gap for the highest reward outcome. Hence,

$$\log(n) - \frac{n-1}{n} - D_{\text{KL}}(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) \geq g_n(1 - \varepsilon_\infty, 1), \quad (79)$$

where $g_n(\cdot, \cdot)$ is defined in Equation (45), and is given by

$$g_n(1 - \varepsilon_\infty, 1) = (1 - (1 - \varepsilon_\infty)^n) \log \frac{n\varepsilon_\infty}{1 - (1 - \varepsilon_\infty)^n} - (n-1)(1 - \varepsilon_\infty)^n \log(1 - \varepsilon_\infty) - \frac{n-1}{n}(1 - (1 - \varepsilon_\infty)^n), \quad (80)$$

completing the proof. \square

Proof of Corollary 3.7. The proof is completed by the following inequalities:

$$g_n(1 - \varepsilon_\infty, 1) = (1 - (1 - \varepsilon_\infty)^n) \log \frac{n\varepsilon_\infty}{1 - (1 - \varepsilon_\infty)^n} - (n-1)(1 - \varepsilon_\infty)^n \log(1 - \varepsilon_\infty) - \frac{n-1}{n}(1 - (1 - \varepsilon_\infty)^n) \quad (81)$$

$$\geq (1 - e^{-n\varepsilon_\infty}) \left(\log \frac{n\varepsilon_\infty}{1 - (1 - \varepsilon_\infty)^n} + (n-1) \frac{(1 - \varepsilon_\infty)^n}{1 - (1 - \varepsilon_\infty)^n} \log \frac{1}{1 - \varepsilon_\infty} - \frac{n-1}{n} \right) \quad (82)$$

$$\geq (1 - e^{-n\varepsilon_\infty}) \left(\log \frac{n\varepsilon_\infty}{1 - (1 - \varepsilon_\infty)^n} + (n-1) \frac{\varepsilon_\infty(1 - \varepsilon_\infty)^n}{1 - (1 - \varepsilon_\infty)^n} - \frac{n-1}{n} \right) \quad (83)$$

$$= (1 - e^{-n\varepsilon_\infty}) \left(\log \frac{n\varepsilon_\infty}{1 - (1 - \varepsilon_\infty)^n} + \frac{n-1}{n} \left(\frac{n\varepsilon_\infty}{1 - (1 - \varepsilon_\infty)^n} (1 - \varepsilon_\infty)^n - 1 \right) \right) \quad (84)$$

$$\geq (1 - e^{-n\varepsilon_\infty}) \left(\log(n\varepsilon_\infty) - \frac{n-1}{n} \right) \quad (85)$$

$$\geq (1 - e^{-n\varepsilon_\infty}) \log(n\varepsilon_\infty) - 1. \quad (86)$$

Hence, as $n \rightarrow \infty$,

$$g_n(1 - \varepsilon_\infty, 1) \geq \log(n\varepsilon_\infty) + o_n(\log n), \quad (87)$$

which completes the proof. \square

A.2. Proofs of Section 4

Proof of Lemma 4.1. The proof follows from:

$$D_{\text{KL}}(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) = E_{\mathbf{y}\sim\pi^{(n)}} \left[\log \left(\frac{\mathcal{F}_{\pi_{\text{ref}}}(\mathbf{y}|\mathbf{x})^n - \mathcal{F}_{\pi_{\text{ref}}}^-(\mathbf{y}|\mathbf{x})^n}{\pi_{\text{ref}}(\mathbf{y}|\mathbf{x})} \right) \right] \quad (88)$$

$$\leq E_{\mathbf{y}\sim\pi^{(n)}} \left[\log \left(\frac{1 - (1 - \pi_{\text{ref}}(\mathbf{y}|\mathbf{x}))^n}{\pi_{\text{ref}}(\mathbf{y}|\mathbf{x})} \right) \right]. \quad (89)$$

□

Proof of Corollary 4.2. Recall that $\varepsilon_\infty = \pi_{\text{ref}}(\mathbf{y}_{\text{max}}|\mathbf{x})$ where $\mathbf{y}_{\text{max}} \sim \pi_{\mathbf{y}|\mathbf{x}}^{(\infty)}$. Notice that

$$D_{\text{KL}}(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) \leq \int_0^{1-\varepsilon_\infty} nv^{n-1} \log(nv^{n-1}) dv + (1 - (1 - \varepsilon_\infty)^n) \log \frac{1 - (1 - \varepsilon_\infty)^n}{\varepsilon_\infty}, \quad (90)$$

which follows the same lines as in the proof of Theorem 3.1 except that we singled out the highest reward outcome, and bounded the rest of the terms. The proof is completed by invoking Lemma A.2 to express the integral in closed form. □

Proof of Lemma 4.5. First notice that for any $0 \leq \varepsilon \leq 1$,

$$(1 - (1 - \varepsilon)^n) \log \frac{1 - (1 - \varepsilon)^n}{\varepsilon} \leq \int_{1-\varepsilon}^1 nv^{n-1} \log(nv^{n-1}) dv, \quad (91)$$

which is implied by Lemma A.1 and setting $b = 1$ and $a = 1 - \varepsilon$. The proof is completed by noticing that

$$\widehat{D}_{\text{KL}}(\varepsilon_n) = \int_0^{1-\varepsilon_n} nv^{n-1} \log(nv^{n-1}) dv + (1 - (1 - \varepsilon_n)^n) \log \frac{1 - (1 - \varepsilon_n)^n}{\varepsilon_n}, \quad (92)$$

and

$$\widetilde{\text{KL}}_n = \int_0^{1-\varepsilon_n} nv^{n-1} \log(nv^{n-1}) dv + \int_{1-\varepsilon_n}^1 nv^{n-1} \log(nv^{n-1}) dv. \quad (93)$$

□

A.3. Proofs of Section 5

Proof of Lemma 5.1. Recall the definition of \mathcal{F}_π and \mathcal{F}_π^- , hence

$$\begin{aligned} E_{\mathbf{z} \sim \pi_{\text{ref}}(\cdot|\mathbf{x})} [\mathcal{W}_r(\mathbf{y} \succ \mathbf{z}|\mathbf{x})] &= E_{\mathbf{z} \sim \pi_{\text{ref}}(\cdot|\mathbf{x})} \left\{ \mathbf{1}(r(\mathbf{x}, \mathbf{y}) > r(\mathbf{x}, \mathbf{z})) + \frac{1}{2} \mathbf{1}(r(\mathbf{x}, \mathbf{y}) = r(\mathbf{x}, \mathbf{z})) \right\} \\ &= \frac{1}{2} (\mathcal{F}_{\pi_{\text{ref}}}(\mathbf{y}|\mathbf{x}) + \mathcal{F}^-(\mathbf{y}|\mathbf{x})), \end{aligned} \quad (94)$$

which completes the proof. \square

Lemma A.4. For $0 \leq a \leq b \leq 1$, let

$$h_n(a, b) := \frac{n}{n+1} (b^{n+1} - a^{n+1}) - \frac{1}{2} (b^n - a^n)(b+a). \quad (95)$$

We have

$$h_n(a, b) \geq 0. \quad (96)$$

Equivalently,

$$\frac{1}{2} (b^n - a^n)(b+a) \leq \frac{n}{n+1} (b^{n+1} - a^{n+1}). \quad (97)$$

Proof. Let us consider the function $h_n(a, v)$ for fixed a as a function of v . Given that $h_n(a, a) = 0$, it suffices to show that $\frac{\partial}{\partial v} h_n(a, v) \geq 0$ for all $0 \leq a \leq v \leq 1$. We have

$$\frac{\partial}{\partial v} h_n(a, v) = nv^n - \frac{n}{2} v^{n-1} (v+a) - \frac{1}{2} (v^n - a^n) \quad (98)$$

$$= \frac{1}{2} (nv^{n-1} (v-a) - (v^n - a^n)) \quad (99)$$

$$\geq 0, \quad (100)$$

completing the proof. \square

Lemma A.5. The following is an upper bound on $h_n(a, b)$:

$$h_n(a, b) \leq \frac{n-1}{2} (b-a)^2. \quad (101)$$

Proof. We have

$$h_n(a, b) = \frac{n}{n+1} (b^{n+1} - a^{n+1}) - \frac{1}{2} (b^n - a^n)(b+a) \quad (102)$$

$$= (b-a) \left(\frac{n}{n+1} \sum_{i=0}^n b^i a^{n-i} - \frac{1}{2} (a+b) \sum_{i=0}^{n-1} b^i a^{n-1-i} \right) \quad (103)$$

$$= (b-a) \left(\left(\frac{n}{n+1} - \frac{1}{2} \right) (a^n + b^n) + \left(\frac{n}{n+1} - 1 \right) \sum_{i=1}^{n-1} a^i b^{n-i} \right) \quad (104)$$

$$\leq (b-a) \left(\frac{n-1}{2(n+1)} (a^n + b^n) - \frac{n-1}{n+1} a^n \right) \quad (105)$$

$$= (b-a) \frac{n-1}{2(n+1)} (b^n - a^n) \quad (106)$$

$$\leq (b-a)^2 \frac{n(n-1)}{2(n+1)} b^{n-1} \quad (107)$$

$$\leq \frac{n-1}{2} (b-a)^2, \quad (108)$$

completing the proof. \square

Proof of Theorem 5.3. Note that

$$\mathcal{W}_r(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) = \frac{1}{2} \sum_{\mathbf{y} \in \mathcal{Y}^*} (\mathcal{F}_{\pi_{\text{ref}}}(\mathbf{y}|\mathbf{x})^n - \mathcal{F}_{\pi_{\text{ref}}}^-(\mathbf{y}|\mathbf{x})^n) (\mathcal{F}_{\pi_{\text{ref}}}(\mathbf{y}|\mathbf{x}) + \mathcal{F}_{\pi_{\text{ref}}}^-(\mathbf{y}|\mathbf{x})) \quad (109)$$

$$\leq \frac{n}{n+1} \sum_{\mathbf{y} \in \mathcal{Y}^*} (\mathcal{F}_{\pi_{\text{ref}}}(\mathbf{y}|\mathbf{x})^{n+1} - \mathcal{F}_{\pi_{\text{ref}}}^-(\mathbf{y}|\mathbf{x})^{n+1}) \quad (110)$$

$$= \frac{n}{n+1}, \quad (111)$$

where Equation (109) follows from Lemma 5.2 and Equation (110) follows from Lemma A.4. \square

Proof of Theorem 5.4. We have

$$\frac{n}{n+1} - \mathcal{W}_r(\pi^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) = \sum_{\mathbf{y} \in \mathcal{Y}^*} h_n(\mathcal{F}_{\pi_{\text{ref}}}^-(\mathbf{y}|\mathbf{x}), \mathcal{F}_{\pi_{\text{ref}}}(\mathbf{y}|\mathbf{x})) \quad (112)$$

$$\leq \frac{n-1}{2} \sum_{\mathbf{y} \in \mathcal{Y}^*} (\mathcal{F}_{\pi_{\text{ref}}}(\mathbf{y}|\mathbf{x}) - \mathcal{F}_{\pi_{\text{ref}}}^-(\mathbf{y}|\mathbf{x}))^2 \quad (113)$$

$$= \frac{n-1}{2} \sum_{\mathbf{y} \in \mathcal{Y}^*} (\pi_{\text{ref}}(\mathbf{y}|\mathbf{x}))^2 \quad (114)$$

$$= \frac{n-1}{2} e^{-H_2(\pi_{\text{ref}}|\mathbf{x})}, \quad (115)$$

where Equation (113) follows from Lemma A.5. \square

Proof of Theorem 5.6. Recall Lemma A.4. Therefore,

$$G_{\mathcal{W}}^{(n)}(\mathbf{x}) \geq h_n(1 - \varepsilon_\infty, 1) \geq 0. \quad (116)$$

The proof is completed by noticing from Equation (95) that

$$h_n(1 - \varepsilon_\infty, 1) = \frac{n}{n+1} (1 - (1 - \varepsilon_\infty)^{n+1}) - (1 - (1 - \varepsilon_\infty)^n) \left(1 - \frac{\varepsilon_\infty}{2}\right). \quad (117)$$

\square

Proof of Corollary 5.7. As $n \rightarrow \infty$,

$$\frac{n}{n+1} (1 - (1 - \varepsilon_\infty)^{n+1}) - (1 - (1 - \varepsilon_\infty)^n) \left(1 - \frac{\varepsilon_\infty}{2}\right) = \frac{\varepsilon_\infty}{2} (1 + o_n(1)), \quad (118)$$

which completes the proof. \square

A.4. Proofs of Section 6

Proof of Lemma 6.1. The proof is completed by observing that with probability $w_\Phi(\mathbf{x})$ a good outcome is obtained, and with probability $(1 - w_\Phi(\mathbf{x}))$ the trial is repeated. \square

Proof of Lemma 6.2. First consider the KL divergence as follows:

$$D_{\text{KL}}(\pi_\Phi(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) = E_{\mathbf{y}\sim\pi_\Phi(\cdot|\mathbf{x})} \log \frac{\pi_\Phi(\mathbf{y}|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}|\mathbf{x})} \quad (119)$$

$$= E_{\mathbf{y}\sim\pi_\Phi(\cdot|\mathbf{x})} \log \frac{1}{w_\Phi(\mathbf{x})} \quad (120)$$

$$= \log \frac{1}{w_\Phi(\mathbf{x})}. \quad (121)$$

Next, let us consider the win rate:

$$\mathcal{W}_r(\pi_\Phi(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) = (1 - w_\Phi(\mathbf{x})) \times 1 + w_\Phi(\mathbf{x}) \times \frac{1}{2} \quad (122)$$

$$= 1 - \frac{1}{2}w_\Phi(\mathbf{x}), \quad (123)$$

completing the proof. \square

Proof of Theorem 6.3. Let $w_\Phi(\mathbf{x})$ be the probability of selecting a good sequence. Then, we note that

$$\frac{-\log(w_\Phi(\mathbf{x}))}{w_\Phi(\mathbf{x})} = \sum_{k=1}^{\infty} H_k (1 - w_\Phi(\mathbf{x}))^k,$$

where for $k \geq 1$ $H_k = \sum_{i=1}^k \frac{1}{i}$ is the k -th Harmonic number and $H_0 = 0$. Hence,

$$-\log(w_\Phi(\mathbf{x})) = \sum_{k=1}^{\infty} H_k (1 - w_\Phi(\mathbf{x}))^k w_\Phi(\mathbf{x}).$$

Let M be the location of the first one (the number of trials). Then,

$$P(M = k) = (1 - w_\Phi(\mathbf{x}))^{k-1} w_\Phi(\mathbf{x}).$$

Hence, the following is an unbiased estimator of $-\log(w_\Phi(\mathbf{x}))$:

$$H_{M-1}.$$

\square

A.5. Proofs of Section 7

Theorem A.6. For any \mathbf{x} and a given π_{ref} , let the function $W_L(D_L)$ be implicitly defined for $\tau \geq 1$:

$$W_L(\tau) = \frac{\tau}{\tau + 1} - \frac{\tau - 1}{2} e^{-H_2(\pi_{\text{ref}}|\mathbf{x})},$$

$$D_L(\tau) = \log(\tau) - \frac{\tau - 1}{\tau}.$$

W_L is a lower bound on the tradeoff curve of best-of-n.

Proof. This is obtained by putting together Theorem 3.1 and Theorem 5.4. □

Theorem A.7. For any \mathbf{x} and a given π_{ref} , let the function $W_U(D_U)$ be implicitly defined for $\tau \geq 1$:

$$W_U(\tau) = \frac{\tau}{\tau + 1},$$

$$D_U(\tau) = \log(\tau) - \frac{\tau - 1}{\tau} - 2\tau(\tau - 1)e^{-H_2(\pi_{\text{ref}}|\mathbf{x})}.$$

W_U is an upper bound on the tradeoff curve of best-of-n.

Proof. This is obtained by combining Theorem 5.3 and Theorem 3.4. □

B. KL divergence of blockwise best-of- n

While best-of- n could be viewed as a sequence-level rejection sampling, most generative models perform decoding step by step. For example, language models perform generation token-by-token or diffusion models perform generation through several denoising steps. Best-of- n could be extended to exert control at different levels of granularity. In particular, best-of- n has been extended to provide control through blockwise decoding (Mudgal et al., 2024), where a block of length B is decoded n times and one with highest token-level reward is selected, and decoding is continued as such. Blockwise decoding could also be viewed as a simple form of tree search and also beam search. Our main result in this section characterizes the KL divergence of blockwise best-of- n with respect to the reference policy.

For the purpose of this section, we assume that the fully decoded response $\mathbf{y} := (y_1, \dots, y_T)$ consist of T intermediate steps (which are tokens in the context of language model and denoising steps in the context of diffusion). For the purposes of this section, we focus on the presentation using generative language models. We let $y_T = \text{EOS}$ be a special token that determines the end of sequence. The autoregressive decoding could be further explained as follows. In the first decoding step, the model π assigns a probability distribution over the first token given by $\pi(\cdot|\mathbf{x})$, and one token y_1 is drawn from this distribution. Next, the second token, y_2 is drawn from $\pi(\cdot|\mathbf{x}, y_1)$, and so on. In short, in each step a token is drawn from $\pi(\cdot|\mathbf{x}, y^{t-1})$ where $y^t := (y_1, \dots, y_t)$. Note that through the chain rule, any language model could be equivalently expressed as a next-token predictor.

Blockwise best-of- n decoding rule was recently proposed by Mudgal et al. (2024). Here, the decoder samples n blocks of length B tokens from the reference model and accepts one with highest reward. Formally, if the decoder has already decoded y^t and the context is \mathbf{x} , then

$$z^B := \arg \max_{\{z_{(k)}^B\}_{k \in [n]}} r(\mathbf{x}, y^t, z_{(k)}^B), \quad (124)$$

where $z_{(k)}^B$ is the k -th continuation of length B drawn independently from the reference model:

$$\{z_{(k)}^B\}_{k \in [n]} \stackrel{i.i.d.}{\sim} \pi_{\text{ref}}(\cdot|\mathbf{x}, y^t), \quad (125)$$

and decoding continues until EOS is reached. We call the resulting distribution $\pi_B^{(n)}$. Note that when $B \rightarrow \infty$, $\pi_B^{(n)}$ becomes the sequence level best-of- n distribution, $\pi^{(n)}$.

Here, we assume we have access to a reward model can produce scalar values for $r(\mathbf{x}, y^t)$ for any partially decoded sequence y^t , which extends the definition of a reward model to also score partial sequences in addition to fully decoded sequences. Mudgal et al. (2024) argue that for this extension to be meaningful the extended function should encode the *value function* for the sequence-level reward model. However, for the purposes of bounding the KL divergence, we don't care how the token-level reward function is obtained.

In the blockwise best-of- n decoding, control is applied more frequently than the sequence-level best-of- n , and the blockwise decoding enables to effectively score an exponentially large number of fully decoded sequences (similarly to how beam search works). Thus, intuitively the effective n would be exponentially larger and the KL divergence would be expected to grow linearly with the number of times control is applied, i.e., $|\mathbf{y}|/B$ where $|\mathbf{y}|$ is the length of the decoded sequence and B is the block length. Next, we formalize this intuition by recalling Theorem B.1.

Theorem B.1. *Let a decoder, $\pi_B^{(n)}$, perform block-wise best-of- n with blocks of length B steps. Further, let \mathbf{y} be draw from this decoder in context \mathbf{x} such that $|\mathbf{y}|$ is the length of \mathbf{y} , i.e., the total number of decoding steps. Then,*

$$D_{KL}(\pi_B^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) \leq E_{\mathbf{y} \sim \pi_B^{(n)}(\cdot|\mathbf{x})} \left[\frac{|\mathbf{y}|}{B} \right] \widetilde{KL}_n,$$

where $\lceil \cdot \rceil$ is the ceiling operator (smallest larger integer), and \widetilde{KL}_n is defined in Equation (4).

Proof. Let us consider all possible outcomes for decoding a block of length B . The outcomes either finish and we reach EOS withing the block. We call this event $\mathcal{E} = 1$, or leads to a partially decoded sequence $\mathbf{z} = z^B$, and we call this outcome

$\mathcal{E} = 0$. In this case, we write $\mathbf{y} = (z^B, \mathbf{s})$.

$$D_{\text{KL}}(\pi_B^{(n)}(\cdot|\mathbf{x})\|\pi_{\text{ref}}(\cdot|\mathbf{x})) = E_{\mathbf{y}\sim\pi_B^{(n)}(\cdot|\mathbf{x})} \log \frac{\pi_B^{(n)}(\mathbf{y}|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}|\mathbf{x})} \quad (126)$$

$$= E_{\mathbf{y}\sim\pi_B^{(n)}(\cdot|\mathbf{x})} \left\{ \mathbf{1}(\mathcal{E} = 1) \log \frac{\pi_B^{(n)}(\mathbf{y}|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}|\mathbf{x})} \right\} + E_{\mathbf{y}\sim\pi_B^{(n)}(\cdot|\mathbf{x})} \left\{ \mathbf{1}(\mathcal{E} = 0) \log \frac{\pi_B^{(n)}(\mathbf{y}|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}|\mathbf{x})} \right\} \quad (127)$$

$$= E_{\mathbf{y}\sim\pi_B^{(n)}(\cdot|\mathbf{x})} \left\{ \mathbf{1}(\mathcal{E} = 1) \log \frac{\pi_B^{(n)}(\mathbf{y}|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}|\mathbf{x})} \right\} + E_{z^B\sim\pi_B^{(n)}(\cdot|\mathbf{x})} E_{\mathbf{s}\sim\pi_B^{(n)}(\cdot|\mathbf{x}, z^B)} \left\{ \mathbf{1}(\mathcal{E} = 0) \left(\log \frac{\pi_B^{(n)}(z^B|\mathbf{x})}{\pi_{\text{ref}}(z^B|\mathbf{x})} + \log \frac{\pi_B^{(n)}(\mathbf{s}|\mathbf{x}, z^B)}{\pi_{\text{ref}}(\mathbf{s}|\mathbf{x}, z^B)} \right) \right\} \quad (128)$$

$$= E_{\mathbf{y}\sim\pi_B^{(n)}(\cdot|\mathbf{x})} \left\{ \mathbf{1}(\mathcal{E} = 1) \log \frac{\pi_B^{(n)}(\mathbf{y}|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}|\mathbf{x})} \right\} + E_{z^B\sim\pi_B^{(n)}(\cdot|\mathbf{x})} \left\{ \mathbf{1}(\mathcal{E} = 0) \log \frac{\pi_B^{(n)}(z^B|\mathbf{x})}{\pi_{\text{ref}}(z^B|\mathbf{x})} \right\} + E_{z^B\sim\pi_B^{(n)}(\cdot|\mathbf{x})} E_{\mathbf{s}\sim\pi_B^{(n)}(\cdot|\mathbf{x}, z^B)} \left\{ \mathbf{1}(\mathcal{E} = 0) \log \frac{\pi_B^{(n)}(\mathbf{s}|\mathbf{x}, z^B)}{\pi_{\text{ref}}(\mathbf{s}|\mathbf{x}, z^B)} \right\} \quad (129)$$

$$\leq \widetilde{\text{KL}}_n + E_{z^B\sim\pi_B^{(n)}(\cdot|\mathbf{x})} \left\{ \mathbf{1}(\mathcal{E} = 0) D_{\text{KL}}(\pi_B^{(n)}(\cdot|\mathbf{x}, z^B)\|\pi_{\text{ref}}(\cdot|\mathbf{x}, z^B)) \right\} \quad (130)$$

$$\leq \widetilde{\text{KL}}_n + E_{z^B\sim\pi_B^{(n)}(\cdot|\mathbf{x})} \mathbf{1}(\mathcal{E} = 0) E_{\mathbf{s}\sim\pi_B^{(n)}(\cdot|\mathbf{x}, z^B)} \left\lceil \frac{|\mathbf{s}|}{B} \right\rceil \widetilde{\text{KL}}_n \quad (131)$$

$$= E_{\mathbf{y}\sim\pi_B^{(n)}(\cdot|\mathbf{x})} \left\lceil \frac{|\mathbf{y}|}{B} \right\rceil \widetilde{\text{KL}}_n, \quad (132)$$

where Equation (131) follows from recursively applying the same procedure to the subsequent blocks. \square

This theorem immediately suggests that $\lceil \frac{|\mathbf{y}|}{B} \rceil \widetilde{\text{KL}}_n$ could be used as an estimator for the KL divergence of block-wise best-of- n , and in expectation the estimator provides an upper bound on the KL divergence of blockwise best-of- n and the reference model. Sequence-level best-of- n could be viewed as blockwise best-of- n with $B \rightarrow \infty$, and Theorem B.1 simplifies to Theorem 3.1 in this asymptotic regime.

C. Experiments

C.1. Details of experiments on Alpaca dataset

We select the following four prompts from the Alpaca dataset (Taori et al., 2023).

- P1 Transform the following sentence into a yes/no question. It is going to rain tomorrow.
- P2 Describe the function of a computer motherboard.
- P3 What is the capital of France?
- P4 Give three tips for staying healthy.

We plot results based on Gemma 2 9B parameter instruction tuned model (Gemma et al., 2024) with temperature one. We further modify the prompts as per the instructions in <https://ai.google.dev/gemma/docs/formatting>. We use log-likelihood of the reference model as the reward. To get the better upper bound we use Corollary 4.2, where we bound ε_∞ using 100 samples.

C.2. Computation of tilted min/max

We compute the tilted average of the estimator defined by (Li et al., 2023, Equation (2)):

$$\frac{1}{t} \log E_{\mathbf{y} \sim \pi^{(n)}(\cdot|\mathbf{x})} \left[e^{t d_n(\varepsilon_n)} \right]. \quad (133)$$

Note that Equation (133) recovers $\min_{\mathbf{y} \sim \pi^{(n)}(\cdot|\mathbf{x})} [d_n(\varepsilon_n)]$ for $t \rightarrow -\infty$, and $\max_{\mathbf{y} \sim \pi^{(n)}(\cdot|\mathbf{x})} [d_n(\varepsilon_n)]$ for $t \rightarrow \infty$ (Li et al., 2023). To capture the variance, we call the value for $t = -1$ the *tilted min*, and the value for $t = 1$ the *tilted max*. The tilted min/max capture the low/high quantiles of the value of the estimator and their difference portrays the deviation from the mean.

C.3. Experiments with machine translation prompts

In this section, we demonstrate the value of the new estimator using machine translation prompts.

- P1 Translate the next sentences to German. I want to buy bread.
- P2 Translate the next sentences to French. I want to buy bread.
- P3 Translate the next sentences to German. A simple and effective method for the inference-time alignment and scaling test-time compute of generative models is the best-of- n policy, where n samples are drawn from a reference policy, ranked based on a reward function, and the highest ranking one is selected. A commonly used analytical expression in the literature claims that the KL divergence between the best-of- n policy and the reference policy is equal to $\log(n) - (n - 1)/n$. We disprove the validity of this claim, and show that it is an upper bound on the actual KL divergence. We also explore the tightness of this upper bound in different regimes, and propose a new estimator for the KL divergence and empirically show that it provides a tight approximation. We also show that the win rate of the best-of- n policy against the reference policy is upper bounded by $n/(n + 1)$ and derive bounds on the tightness of this characterization. We conclude with analyzing the tradeoffs between win rate and KL divergence of the best-of- n alignment policy, which demonstrate that very good tradeoffs are achievable with $n < 1000$.
- P4 Translate the next sentences to French. A simple and effective method for the inference-time alignment and scaling test-time compute of generative models is the best-of- n policy, where n samples are drawn from a reference policy, ranked based on a reward function, and the highest ranking one is selected. A commonly used analytical expression in the literature claims that the KL divergence between the best-of- n policy and the reference policy is equal to $\log(n) - (n - 1)/n$. We disprove the validity of this claim, and show that it is an upper bound on the actual KL divergence. We also explore the tightness of this upper bound in different regimes, and propose a new estimator for the KL divergence and empirically show that it provides a tight approximation. We also show that the win rate of the best-of- n policy against the reference policy is upper bounded by $n/(n + 1)$ and derive bounds on the tightness of this characterization. We conclude with analyzing the tradeoffs between win rate and KL divergence of the best-of- n alignment policy, which demonstrate that very good tradeoffs are achievable with $n < 1000$.

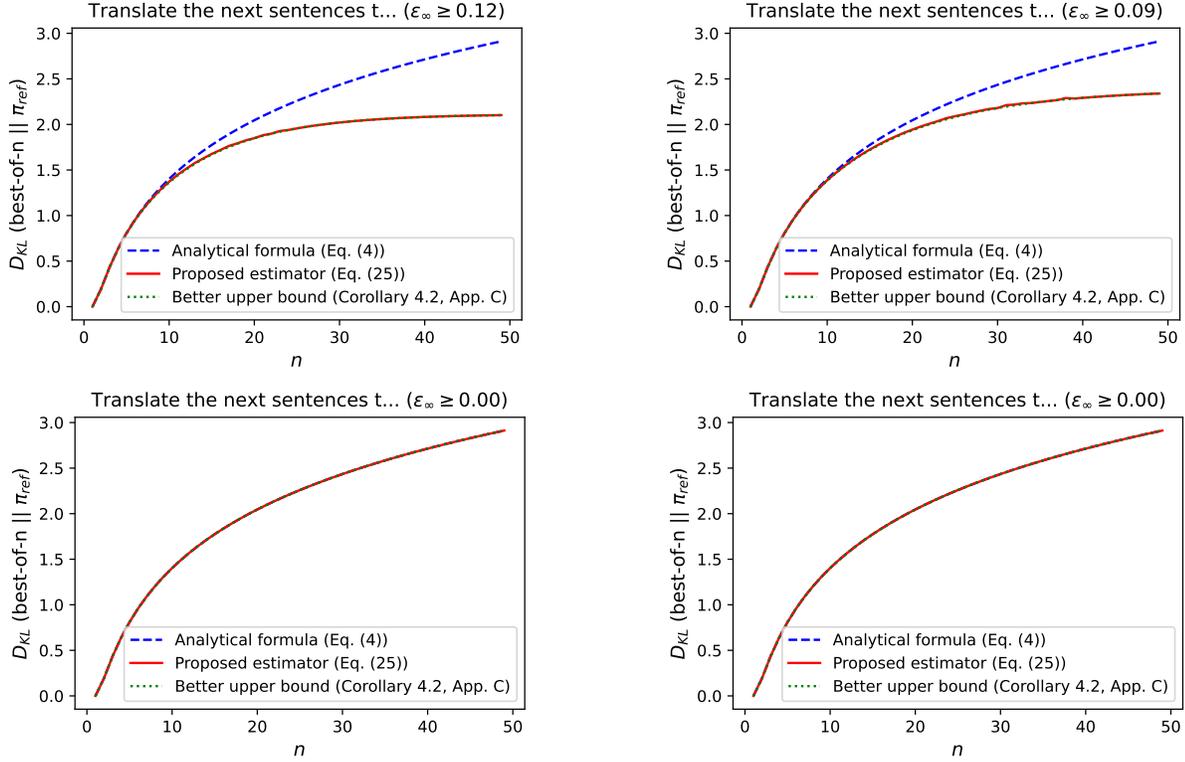


Figure 10. The analytical formula ($\log(n) - (n - 1)/n$, Equation (4), better upper bound (Corollary 4.2, Appendix C), the proposed estimator, Equation (25), and the exact KL divergence, for four machine translation examples using Gemma 9B IT model (Gemma et al., 2024) with reward the log-likelihood of response under the reference model.

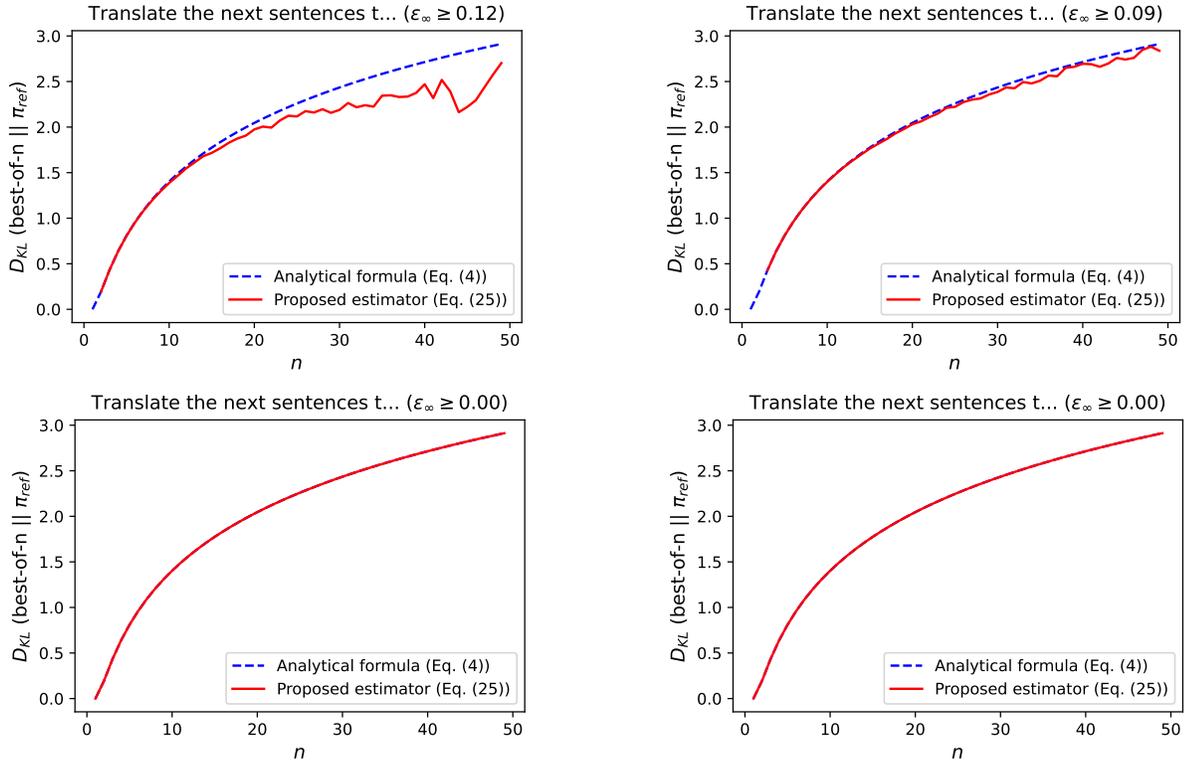


Figure 11. The analytical formula ($\log(n) - (n - 1)/n$, Equation (4), better upper bound (Corollary 4.2, Appendix C), the proposed estimator, Equation (25), and the exact KL divergence, for four machine translation examples using Gemma 9B IT model (Gemma et al., 2024) with reward the negative length of response under the reference model.