

SELF-IMPROVEMENT IN LANGUAGE MODELS: THE SHARPENING MECHANISM

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent work in language modeling has raised the possibility of “self-improvement,” where an LLM evaluates and refines its own generations to achieve higher performance without external feedback. It is impossible for this self-improvement to create information that is not already in the model, so why should we expect that this will lead to improved capabilities?

We offer a new theoretical perspective on the capabilities of self-improvement through a lens we refer to as “sharpening.” Motivated by the observation that language models are often better at verifying response quality than they are at generating correct responses, we formalize self-improvement as using the model itself as a verifier during post-training in order to ‘sharpen’ the model to one placing large mass on high-quality sequences, thereby amortizing the expensive inference-time computation of generating good sequences. We begin by introducing a new statistical framework for sharpening in which the learner has sample access to a pre-trained base policy. Then, we analyze two natural families of self-improvement algorithms based on SFT and RLHF. We find that (i) the SFT-based approach is minimax optimal whenever the initial model has sufficient coverage, but (ii) the RLHF-based approach can improve over SFT-based self-improvement by leveraging online exploration, bypassing the need for coverage. Finally, we empirically validate the sharpening mechanism via both inference-time and amortization experiments. We view these findings as a starting point toward a foundational understanding that can guide the design and evaluation of self-improvement algorithms.

1 INTRODUCTION

Contemporary language models are remarkably proficient on a wide range of natural language tasks (Brown et al., 2020; Ouyang et al., 2022; Touvron et al., 2023; OpenAI, 2023; Google, 2023), but they inherit shortcomings of the data on which they were trained. A fundamental challenge is to achieve better performance than what is directly induced by the distribution of available, human-generated training data. To this end, recent work (Huang et al., 2022; Wang et al., 2022; Bai et al., 2022b; Pang et al., 2023; Yuan et al., 2024) has raised the possibility of “self-improvement,” where a model—typically through forms of self-play or self-training in which the model critiques its own generations—learns to improve on its own, without external feedback. This phenomenon is somewhat counterintuitive; at first glance it would seem to disagree with the well-known data-processing inequality (Cover, 1999), which asserts that no form of self-training should be able to create information not already in the model, motivating the question of why we should expect such supervision-free interventions will lead to stronger reasoning and planning capabilities.

A dominant hypothesis for why improvement without external feedback might be possible is that models contain “hidden knowledge” (Hinton et al., 2015) that is difficult to access. Self-improvement, rather than creating knowledge from nothing, is a means of extracting and distilling this knowledge into a more accessible form, and thus is a computational phenomenon rather than a statistical one. While there is a growing body of empirical evidence for this hidden-knowledge hypothesis (Furlanello et al., 2018; Gotmare et al., 2019; Dong et al., 2019; Abnar et al., 2020; Allen-Zhu & Li, 2020), particularly in the context of self-distillation, a fundamental understanding of self-improvement remains missing. Concretely, where in the model is this hidden knowledge, and when and how can it be extracted?

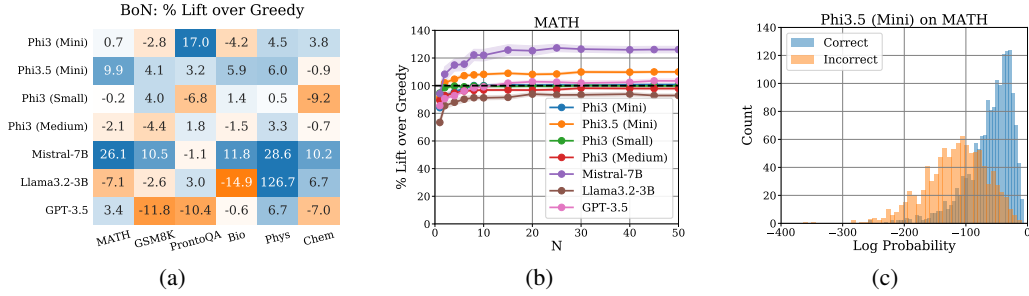


Figure 1: Validation of BoN-Sharpening at inference time. (a) The percent improvement over greedy decoding that BoN for $N = 50$ exhibits on accuracy on 6 tasks and 7 models, colored by performance. (b) The affect that increasing N in BoN has on percent accuracy improvement over greedy for 7 different models. (c) The distribution of sequence-level log probabilities of BoN with $N=1$ sampled completions from Phi3.5-Mini on the MATH dataset, conditioned on whether or not the completion is correct. Correct completions are noticeably higher likelihood than incorrect completions, demonstrating the utility of inference-time sharpening.

1.1 OUR PERSPECTIVE: THE SHARPENING MECHANISM

In this paper, we posit a potential source of hidden knowledge, and offer a formal perspective on how to extract it. Our starting point is the widely observed phenomenon that language models are often better at verifying whether responses are correct than they are at generating correct responses (Huang et al., 2022; Wang et al., 2022; Bai et al., 2022b; Pang et al., 2023; Yuan et al., 2024). This gap may be explained by the theory of computational complexity, which suggests that generating high-quality responses can be less computationally tractable than verification (Cook, 1971; Levin, 1973; Karp, 1972). In autoregressive language modeling, for example, computing the most likely response for a given prompt is NP-hard in the worst case (Appendix D), whereas the model’s likelihood for a given response can be easily evaluated.

We view self-improvement as any attempt to narrow this gap, i.e., use the model as its own verifier to improve generation and *sharpen* the model toward high-quality responses. Formally, consider a learner with access to a base model $\pi_{\text{base}} : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$ representing a conditional distribution that maps mapping a prompt $x \in \mathcal{X}$ to a distribution over responses (i.e., $\pi_{\text{base}}(y | x)$ is the probability that the model generates the response y given the prompt x).¹ In applications, we consider π_{base} to be trained either through next-token prediction, or through additional post-training steps such as SFT or RLHF, with the key feature being that π_{base} is a good verifier, as measured by some *self-reward* function $r_{\text{self}}(y | x; \pi_{\text{base}})$ measuring model certainty. The self-reward function is derived purely from the base model π_{base} , without external supervision or feedback. Examples include normalized and/or regularized sequence likelihood (Meister et al., 2020), models-as-judges (Zheng et al., 2024; Yuan et al., 2024; Wu et al., 2024a; Wang et al., 2024), and model confidence (Wang & Zhou, 2024).

We refer to **sharpening** as any process that tilts π_{base} toward responses that are more certain in the sense that they enjoy greater self-reward r_{self} . More formally, a sharpened model $\hat{\pi}$ is one that (approximately) maximizes the self-reward:

$$\hat{\pi}(x) \approx \arg \max_{y \in \mathcal{Y}} r_{\text{self}}(y | x; \pi_{\text{base}}) \quad (1)$$

Note that, in Eq. (1), y denotes an entire response, rather than a single token. Sharpening may be implemented at inference-time, or **amortized** via self-training (Section 2). Popular decoding strategies such as greedy, low-temperature sampling, and beam-search can all be viewed as instances of the former (albeit at the token-level).² The latter captures many existing self-training schemes

¹Our general results are agnostic to the structure of \mathcal{X} , \mathcal{Y} , and π_{base} , but an important special case for language modeling is the autoregressive setting where $\mathcal{Y} = \mathcal{V}^H$ for a vocabulary space \mathcal{V} and sequence length H , and where π_{base} has the autoregressive structure $\pi_{\text{base}}(y_{1:H} | x) = \prod_{h=1}^H \pi_{\text{base},h}(y_h | y_{1:h-1}, x)$ for $y = y_{1:H} \in \mathcal{Y}$.

²More sophisticated decoding strategies like normalized/regularized sequence likelihood (Meister et al., 2020) or chain-of-thought decoding (Wang & Zhou, 2024) also admit an interpretation as sharpening; see Appendix A.

(Huang et al., 2022; Wang et al., 2022; Bai et al., 2022b; Pang et al., 2023; Yuan et al., 2024), and is the main focus of this paper; we use the term *sharpening* without further qualification to refer to the latter.

We refer to the **sharpening mechanism** as the phenomenon where responses from a model with the highest certainty (in the sense of large self-reward r_{self}) exhibit the greatest performance on a task of interest. Though it is unclear a-priori whether there are self-rewards related to task performance, the successes of self-improvement in prior works (Huang et al., 2022; Wang et al., 2022; Bai et al., 2022b; Pang et al., 2023; Yuan et al., 2024) give strong positive evidence. These works suggest that, in many settings, models do have hidden knowledge: the model’s own self-reward correlates with response quality, but it is computationally challenging to generate high self-rewarding—and thus high quality—responses. It is the role of (algorithmic) sharpening to leverage these verifications to improve the quality of generations, despite computational difficulty.

1.2 CONTRIBUTIONS

We initiate the theoretical study of self-improvement via the sharpening mechanism. We disentangle the choice of self-reward from the algorithms used to optimize it, and aim to understand: (i) When and how does self-training achieve sharpening? (ii) What are the fundamental limits for such algorithms?

Algorithms for sharpening (Section 2). The starting point for our work is to consider two natural families of self-improvement algorithms based on supervised fine-tuning (SFT) and reinforcement learning (RL/RLHF), respectively, SFT-Sharpener and RLHF-Sharpener. Both algorithms **amortize** the sharpening objective (1) into a dedicated post-training/fine-tuning phase:

- SFT-Sharpener filters responses where the self-reward $r_{\text{self}}(y \mid x; \pi_{\text{base}})$ is large and fine-tunes on the resulting dataset, invoking common SFT pipelines (Amini et al., 2024; Sessa et al., 2024).
- RLHF-Sharpener directly applies reinforcement learning techniques (e.g., PPO (Schulman et al., 2017) or DPO (Rafailov et al., 2023)) to optimize the self-reward function $r_{\text{self}}(y \mid x; \pi_{\text{base}})$.

In the remainder of the paper, we introduce a theoretical framework to analyze the performance of these algorithms. Our main contributions are as follows.

Maximum-likelihood sharpening objective (Section 3.1). As a concrete proposal of one source of hidden knowledge, we consider self-rewards defined by the model’s sequence-level log-probabilities:

$$r_{\text{self}}(y \mid x) := \log \pi_{\text{base}}(y \mid x) \quad (2)$$

This is a stylized self-reward function, which offers perhaps the simplest objective for self-improvement in the absence of external feedback (i.e., purely supervision-free), yet also connects self-improvement to a rich body of theoretical computer science literature on computational trade-offs for optimization (inference) versus sampling (Appendix A). In spite of its simplicity, maximum-likelihood sharpening is already sufficient to achieve non-trivial performance gains over greedy decoding on a range of reasoning tasks with several language models; cf. Figure 1. We believe it can serve as a starting point toward understanding forms of self-improvement that use more sophisticated self-rewarding (Huang et al., 2022; Wang et al., 2022; Pang et al., 2023; Yuan et al., 2024).

A statistical framework for sharpening (Sections 3.2 and 3.3). Though the goal of sharpening is computational in nature, we recast self-training according to the maximum-likelihood sharpening objective Eq. (2) as a **statistical** problem where we aim to produce a model approximating (1) using a polynomial number of (i) sample prompts $x \sim \mu$, (ii) sampling queries of the form $y \sim \pi_{\text{base}}(x)$, and (iii) likelihood evaluations of the form $\pi_{\text{base}}(y \mid x)$. Evaluating the efficiency of the algorithm through the number of such queries, this abstraction offers a natural way to evaluate the performance of self-improvement/sharpening algorithms and establish fundamental limits; we use our framework to prove new lower bounds that highlight the importance of the base model’s coverage.

Analysis of sharpening algorithms (Section 4). Within our statistical framework for sharpening, we show that SFT-Sharpener and RLHF-Sharpener provably converge to sharpened models, establishing several results: (i) **SFT-Sharpener is minimax optimal**, and learns a sharpened model whenever π_{base} has sufficient coverage (we also show that a novel variant based on adaptive sampling can sidestep the minimax lower bound); (ii) **RLHF-Sharpener benefits from on-policy exploration**, and can bypass the need for coverage—improving over SFT-Sharpener.

Empirical Investigation (Appendix E). In addition to our theoretical results, we explore empirically the extent to which our theoretical framework can aid language models in a variety of tasks. In Appendix E, we consider three choices of self-reward on an extensive list of model-dataset pairs and conclude that sharpening can often improve performance. We then implement one of our algorithms, SFT-Sharpener, on a subset of these model-dataset pairs and observe a significant positive effect on performance. A summary of our inference-time experiments can be found in Figure 1.

1.3 RELATED WORK

Our work is most directly related to a growing body of empirical research that studies self-improvement/training for language models in a supervision-free setting with no external feedback (Huang et al., 2022; Wang et al., 2022; Bai et al., 2022b; Pang et al., 2023; Yuan et al., 2024). The specific algorithms for self-improvement/sharpening we study can be viewed as applications of standard alignment algorithms (Amini et al., 2024; Sessa et al., 2024; Christiano et al., 2017; Bai et al., 2022a; Ouyang et al., 2022; Rafailov et al., 2023) with a specific choice of reward function. However, note that the maximum likelihood sharpening objective (2) used for our theoretical results has been relatively unexplored within the alignment and self-improvement literature.

On the theoretical side, current understanding of self-training is limited. One line of work, focusing on the *self-distillation* objective (Hinton et al., 2015) for classification and regression, aims to provide convergence guarantees for self-training in stylized setups such as linear models (Mobahi et al., 2020; Frei et al., 2022; Das & Sanghavi, 2023; Das et al., 2024; Pareek et al., 2024), with Allen-Zhu & Li (2020) giving guarantees for feedforward neural networks. To the best of our knowledge, our work is the first to study self-training in a general framework that subsumes language modeling. See Appendix A for a more extensive discussion of related work.

2 SHARPENING ALGORITHMS FOR SELF-IMPROVEMENT

This section introduces the two families of self-improvement algorithms for sharpening that we study. Going forward, we will omit the dependence of r_{self} on π_{base} , when it is clear from context. We will also use the notation $\arg \max_{\pi \in \Pi}$ or $\arg \min_{\pi \in \Pi}$ to denote exact optimization over a user-specified model class Π for theoretical results (Agarwal et al., 2019; Foster & Rakhlin, 2023); empirically, these operations can be implemented by training a neural network to low loss.

2.1 SELF-IMPROVEMENT THROUGH SFT: SFT-Sharpener

SFT-Sharpener filters responses for which the self-reward $r_{\text{self}}(y \mid x;)$ is large, and applies standard supervised fine-tuning on the resulting dataset (Amini et al., 2024; Sessa et al., 2024; Gui et al., 2024; Pace et al., 2024). This can be viewed as amortizing inference-time sharpening via the effective-but-costly best-of- N sampling approach (Brown et al., 2024; Snell et al., 2024; Wu et al., 2024b). Concretely, suppose we have a collection of prompts x_1, \dots, x_n . For each prompt, we sample N responses $y_{i,1}, \dots, y_{i,N} \sim \pi_{\text{base}}(\cdot \mid x_i)$, then compute the best-of- N response $y_i^{\text{BoN}} = \arg \max_{j \in [N]} \{r_{\text{self}}(y_{i,j} \mid x_i)\}$, scoring via the model’s self-reward function. We compute

$$\hat{\pi}^{\text{BoN}} = \arg \max_{\pi \in \Pi} \sum_{i=1}^n \log \pi(y_i^{\text{BoN}} \mid x_i).$$

This is a simple, flexible self-training scheme, and converges to a sharpened model as $n, N \rightarrow \infty$.

2.2 SELF-IMPROVEMENT THROUGH RLHF: RLHF-Sharpener

A drawback of the SFT-Sharpener algorithm is that it may ignore useful information contained in the self-reward function $r_{\text{self}}(y \mid x)$. Fixing a regularization parameter $\beta > 0$ throughout, our second class of algorithms solve a KL-regularized reinforcement learning problem in the spirit of RLHF and other alignment methods (Christiano et al., 2017; Rafailov et al., 2023). Defining $\mathbb{E}_{\pi}[\cdot] = \mathbb{E}_{x \sim \mu, y \sim \pi_{\text{base}}(\cdot \mid x)}[\cdot]$ and $D_{\text{KL}}(\pi \parallel \pi_{\text{base}}) = \mathbb{E}_{\pi}[\log \frac{\pi(y \mid x)}{\pi_{\text{base}}(y \mid x)}]$, we choose

$$\hat{\pi} \approx \arg \max_{\pi \in \Pi} \{\mathbb{E}_{\pi}[r_{\text{self}}(y \mid x)] - \beta D_{\text{KL}}(\pi \parallel \pi_{\text{base}})\}. \quad (3)$$

The exact optimizer $\pi_{\beta}^* = \arg \max_{\pi \in \Pi} \{\mathbb{E}_{\pi}[r_{\text{self}}(y \mid x)] - \beta D_{\text{KL}}(\pi \parallel \pi_{\text{base}})\}$ for this objective has the form $\pi_{\beta}^*(y \mid x) \propto \pi_{\text{base}}(y \mid x) \cdot \exp(\beta^{-1} r_{\text{self}}(y \mid x))$, which converges to the solution to the sharpening objective in Eq. (1) as $\beta \rightarrow 0$. Thus Eq. (3) can be seen to encourage sharpening.

There are many possible choices for what RLHF/alignment algorithm to use to solve (3). For our theoretical results, we first implement Eq. (3) using an approach inspired by DPO and its reward-based variants (Rafailov et al., 2023; Gao et al., 2024). Given a dataset $\mathcal{D} = \{(x, y, y')\}$ of n examples sampled via $x \sim \mu$ and $y, y' \sim \pi_{\text{base}}(y | x)$, we consider the algorithm that solves

$$\hat{\pi} \in \arg \min_{\pi \in \Pi} \sum_{(x, y, y') \in \mathcal{D}} \left(\beta \log \frac{\pi(y | x)}{\pi_{\text{base}}(y | x)} - \beta \log \frac{\pi(y' | x)}{\pi_{\text{base}}(y' | x)} - (r_{\text{self}}(y | x) - r_{\text{self}}(y' | x)) \right)^2. \quad (4)$$

In the sequel (Section 4), we will show that this approach leads to comparable guarantees to SFT-Sharpener, but that a more sophisticated DPO variant that incorporates *online exploration* (Xie et al., 2024) can offer provable benefits.

3 A STATISTICAL FRAMEWORK FOR SHARPENING

This section introduces the theoretical framework within which we will analyze the SFT-Sharpener and RLHF-Sharpener algorithms. We first introduce the maximum-likelihood sharpening objective as a simple, stylized self-reward function, then introduce our statistical framework for sharpening. We write $f = \tilde{O}(g)$ to denote $f = O(g \cdot \max\{1, \text{polylog}(g)\})$ and $a \lesssim b$ as shorthand for $a = O(b)$.

3.1 MAXIMUM-LIKELIHOOD SHARPENING

Our theoretical results focus on the maximum-likelihood sharpening objective given by

$$r_{\text{self}}(y | x) := \log \pi_{\text{base}}(y | x).$$

This is a simple and stylized self-reward function, but we will show that it already enjoys a rich theory. In particular, we can restate the problem of maximum-likelihood sharpening as follows.

Can we efficiently **amortize maximum likelihood inference (optimization)** for a conditional distribution $\pi_{\text{base}}(y | x)$ given access to a **sampling oracle** that can sample $y \sim \pi_{\text{base}}(\cdot | x)$?

The tacit assumption in this framing is that the maximum-likelihood response constitutes a useful form of hidden knowledge. Maximum-likelihood sharpening connects the study of self-improvement to a large body of research in theoretical computer science demonstrating computational reductions between optimization (inference) and sampling (generation) (Kirkpatrick et al., 1983; Lovász & Vempala, 2006; Singh & Vishnoi, 2014; Ma et al., 2019; Talwar, 2019). Our sharpening framework offers a new learning-theoretic perspective by focusing on the problem of amortizing this type of reduction.

We evaluate the quality of an approximately sharpened model as follows. Let $\mathbf{y}^*(x) := \arg \max_{y \in \mathcal{Y}} \log \pi_{\text{base}}(y | x)$; we interpret $\mathbf{y}^*(x) \subset \mathcal{Y}$ as a set to accommodate non-unique maximizers, and will write $y^*(x)$ to indicate a unique maximizer when it exists (i.e., when $\mathbf{y}^*(x) = \{y^*(x)\}$).

Definition 3.1 (Sharpened model). *We say that a model $\hat{\pi}$ is (ϵ, δ) -sharpened relative to π_{base} if*

$$\mathbb{P}_{x \sim \mu} [\hat{\pi}(\mathbf{y}^*(x) | x) \geq 1 - \delta] \geq 1 - \epsilon.$$

That is, an (ϵ, δ) -sharpened model places at least $1 - \delta$ mass on arg-max responses on all but an ϵ -fraction of prompts under μ . For small δ and ϵ , we are guaranteed that $\hat{\pi}$ is a high-quality generator: sampling from the model will produce an arg-max response with high probability for most prompts.

Maximum-likelihood sharpening for autoregressive models. Though our most general results are agnostic to the structure of \mathcal{X} , \mathcal{Y} , and π_{base} , an important special case is the autoregressive setting in which $\mathcal{Y} = \mathcal{V}^H$ for a *vocabulary space* \mathcal{V} and sequence length H , and where π_{base} has the autoregressive structure $\pi_{\text{base}}(y_{1:H} | x) = \prod_{h=1}^H \pi_{\text{base},h}(y_h | y_{1:h-1}, x)$ for $y = y_{1:H} \in \mathcal{Y}$. We observe that when the response $y = (y_1, \dots, y_H) \in \mathcal{Y} = \mathcal{V}^H$ is a sequence of tokens, the maximum-likelihood sharpening objective (2) sharpens toward the sequence-level arg-max response:

$$\arg \max_{y_{1:H}} \log \pi_{\text{base}}(y_{1:H} | x). \quad (5)$$

Although somewhat stylized, Eq. (5) is a non-trivial (in general, computationally intractable; see Appendix D) solution concept. In particular, we view the sequence-level arg-max as a form of hidden knowledge that cannot necessarily be uncovered through naive sampling or greedy decoding.

Empirical validation of maximum-likelihood sharpening. Empirically, we find that when π_{base} is a pre-trained language model, inference-time maximum-likelihood sharpening leads to a meaningful performance increase over both direct sampling and greedy decoding. We demonstrate this by appealing to a practical approximation, inference-time sharpening via best-of- N sampling: given a prompt $x \in \mathcal{X}$, we draw N responses $y_1, \dots, y_N \sim \pi_{\text{base}}(\cdot | x)$, and return the response $\hat{y} = \arg \max_{y_i} \log \pi_{\text{base}}(y_i | x)$; this is equivalent to [Stiennon et al. \(2020\)](#); [Gao et al. \(2023\)](#); [Yang et al. \(2024\)](#), with reward $r_{\text{self}}(y | x) = \log \pi_{\text{base}}(y | x)$, and is a popular approach in modern deployments.³ An overview of the results can be found in [Figure 1](#) with details provided in [Appendix E](#). Observed improvements suggest that maximum-likelihood sharpening, while stylized, is a desirable criterion.

Role of δ for autoregressive models. As can be verified through simple examples, beam-search and greedy tokenwise decoding do not, in general, return an exact solution to (5). There is one notable exception, which implies that it always suffices to sharpen to level $\delta = 1/2$ (cf. [Definition 3.1](#)).

Proposition 3.1 (Greedy decoding succeeds for sharpened policies). *Let $\pi = \pi_{1:H}$ be an autoregressive model defined over response space $\mathcal{Y} = \mathcal{V}^H$. For a given prompt $x \in \mathcal{X}$, if $y^*(x) = \{y^*(x)\}$ is a singleton and $\pi(y^*(x) | x) > 1/2$, then the greedy decoding strategy that selects $\hat{y}_h = \arg \max_{y_h \in \mathcal{V}} \pi_h(y_h | \hat{y}_1, \dots, \hat{y}_{h-1}, x)$ guarantees that $\hat{y} = y^*(x)$. This result is sharp, in the sense that there exist π with $\pi(y^*(x) | x) \leq 1/2$ for which greedy decoding fails to recover $y^*(x)$.*

3.2 SAMPLE COMPLEXITY FRAMEWORK

As described, sharpening in the sense of [Definition 3.1](#) is a purely computational problem, which makes it difficult to evaluate the quality and optimality of self-improvement algorithms. To address this, we introduce a novel statistical/information-theoretic framework for sharpening, inspired by the success of oracle complexity in optimization ([Nemirovski et al., 1983](#); [Traub et al., 1988](#); [Raginsky & Rakhlin, 2011](#); [Agarwal et al., 2012](#)) and statistical query complexity in computational learning theory ([Blum et al., 1994](#); [Kearns, 1998](#); [Feldman, 2012](#); [2017](#)).

Definition 3.2 (Sample-and-evaluate framework). *In the **Sample-and-Evaluate** framework, the algorithm designer does not have explicit access to the base model π_{base} . Instead, they access π_{base} only through sample-and-evaluate queries. Concretely, the learner is allowed to sample n prompts $x \sim \mu$. For each prompt x , they can sample N responses $y_1, y_2, \dots, y_N \sim \pi_{\text{base}}(\cdot | x)$ and observe the likelihood $\pi_{\text{base}}(y_i | x)$ for each such response. The efficiency, or sample complexity, of the algorithm is measured through the total number of sample-and-evaluate queries $m := n \cdot N$.*

This framework can be seen to capture algorithms like SFT-Sharpener and RLHF-Sharpener (implemented with DPO), which only access the base model π_{base} through i) sampling responses via $y \sim \pi_{\text{base}}(\cdot | x)$ (**generation**), and ii) evaluating the likelihood $\pi_{\text{base}}(y | x)$ (**verification**) for these responses. We view the sample complexity $m = n \cdot N$ as a natural statistical abstraction for the computational complexity of self-improvement (exactly parallel to oracle complexity for optimization algorithms), one which is amenable to information-theoretic lower bounds.⁴ We will aim to show that, under appropriate assumptions, SFT-Sharpener and RLHF-Sharpener can learn an (ϵ, δ) -sharpened model with sample complexity

$$m = \text{poly}(\epsilon^{-1}, \delta^{-1}, C_{\text{prob}})$$

where C_{prob} is a potentially problem-dependent constant.

3.3 FUNDAMENTAL LIMITS

Before diving into our analysis of SFT-Sharpener and RLHF-Sharpener in the sample-and-evaluate framework, let us take a brief detour to give a sense for how sample complexity guarantees in our framework should scale. To this end, we will prove a lower bound or fundamental limit on the sample complexity of any algorithm in the sample-and-evaluate framework.

Intuitively, the performance of any sharpening algorithm based on sampling should depend on how well the base model π_{base} covers the arg-max response $y^*(x)$. To capture this, we define the

³We mention in passing that inference-time best-of- N sampling enjoys provable guarantees for maximizing the maximum-likelihood sharpening objective when N is sufficiently large. See [Appendix B](#) for details.

⁴Concretely, the sample complexity $m = n \cdot N$ is a lower bound on the running time of any algorithm that operates in the sample-and-evaluate framework.

following coverage coefficient:⁵

$$C_{\text{cov}} = \mathbb{E}_{x \sim \mu} \left[\frac{1}{\pi_{\text{base}}(\mathbf{y}^*(x) \mid x)} \right]. \quad (6)$$

Next, for a model π , we define $\mathbf{y}^\pi(x) = \arg \max_{y \in \mathcal{Y}} \pi(y \mid x)$ and $C_{\text{cov}}(\pi) = \mathbb{E}_{x \sim \mu} \left[\frac{1}{\pi(\mathbf{y}^\pi(x) \mid x)} \right]$.

Our main lower bound shows that for worst-case choice of Π , the coverage coefficient acts as a lower bound on the sample complexity of any algorithm.

Theorem 3.1 (Lower bound for sharpening). *Fix an integer $d \geq 1$ and parameters $\epsilon \in (0, 1)$ and $C \geq 1$. There exists a class of models Π such that (i) $\log |\Pi| \approx d(1 + \log(C\epsilon^{-1}))$, (ii) $\sup_{\pi \in \Pi} C_{\text{cov}}(\pi) \lesssim C$, and (iii) $\mathbf{y}^\pi(x)$ is a singleton for all $\pi \in \Pi$, for which any sharpening algorithm $\hat{\pi}$ that achieves $\mathbb{E}[\mathbb{P}_{x \sim \mu}[\hat{\pi}(\mathbf{y}^{\pi_{\text{base}}}(x) \mid x) > 1/2]] \geq 1 - \epsilon$ for all $\pi_{\text{base}} \in \Pi$ must collect a total number of samples $m = n \cdot N$ at least*

$$m \gtrsim \frac{C \log |\Pi|}{\epsilon^2 \cdot (1 + \log(C\epsilon^{-1}))}.$$

This result shows that the complexity of any $(\epsilon, 1/2 - \delta)$ -sharpening algorithm (for $\delta > 0$) in the sample-and-evaluate framework must depend polynomially on the coverage coefficient, as well as the accuracy ϵ . The lower bound also depends on the expressivity of π_{base} , as captured by the model class complexity term $\log |\Pi|$. We will show in the sequel that it is possible to match this lower bound. Note that this result also implies a lower bound for the general sharpening problem (i.e., general r_{self}), since maximum-likelihood sharpening is a special case.

Remark 3.1 (Relaxed notions of sharpening and coverage). *The notion of coverage in Eq. (6) is somewhat stringent, since it requires π_{base} place large mass on $\mathbf{y}^*(x)$ on average. In Appendix F, we introduce a more general and permissive notion of approximate sharpening (Definition F.1) which leads to weaker coverage requirements, and use this to give generalized versions of our main results.*

We close this section by noting that numerous recent works—focusing on inference-time computation—show that standard language models exhibit favorable coverage with respect to desirable responses (Brown et al., 2024; Snell et al., 2024; Wu et al., 2024b). We replicate these findings in our experimental setup in Appendix E. These works suggest that the coverage coefficient C_{cov} may be small in practice.

4 ANALYSIS OF SHARPENING ALGORITHMS

Equipped with the sample complexity framework from Section 3, we now prove that the SFT-Sharpener and RLHF-Sharpener families of algorithms provably learn a sharpened model for the maximum-likelihood sharpening objective under natural statistical assumptions.

Throughout this section, we treat the model class Π as a fixed, user-specified parameter. Our results—in the tradition of statistical learning theory—allow for general classes Π , and are agnostic to the structure beyond standard generalization arguments.

4.1 ANALYSIS OF SFT-Sharpener

Recall that when we specialize to the maximum-likelihood sharpening self-reward, the SFT-Sharpener algorithm takes the form $\hat{\pi}^{\text{BoN}} = \arg \max_{\pi \in \Pi} \sum_{i=1}^n \log \pi_{\text{base}}(y_i^{\text{BoN}} \mid x_i)$, where $y_i^{\text{BoN}} = \arg \max_{j \in [N]} \{\log \pi_{\text{base}}(y_{i,j} \mid x_i)\}$ for $y_{i,1}, \dots, y_{i,N} \sim \pi_{\text{base}}(\cdot \mid x_i)$.

To analyze SFT-Sharpener, we first make a realizability assumption. Let $\pi_N^{\text{BoN}}(x)$ be the distribution of the random variable $y_N^{\text{BoN}}(x) \sim \arg \max \{\log \pi_{\text{base}}(y_i \mid x) \mid y_1, \dots, y_N \sim \pi_{\text{base}}(x)\}$.

Assumption 4.1. *The model class Π satisfies $\pi_N^{\text{BoN}} \in \Pi$.*

Our main guarantee for SFT-Sharpener is as follows.

⁵This quantity can be interpreted as a special case of the L_1 -concentrability coefficient (Farahmand et al., 2010; Xie & Jiang, 2020; Zanette et al., 2021) studied in the theory of offline reinforcement learning.

Theorem 4.1 (Sample complexity of SFT-Sharpener). *Let $\epsilon, \delta, \rho \in (0, 1)$ be given, and suppose we set $n = c \cdot \frac{\log(|\Pi|\rho^{-1})}{\delta\epsilon}$ and $N^* = c \cdot \frac{C_{\text{cov}} \log(2\delta^{-1})}{\epsilon}$ for an appropriate constant $c > 0$. Then with probability at least $1 - \rho$, SFT-Sharpener produces a model $\hat{\pi}$ such that that $\mathbb{P}_{x \sim \mu}[\hat{\pi}(\mathbf{y}^*(x) | x) \leq 1 - \delta] \leq \epsilon$, and has total sample complexity⁶*

$$m = O\left(\frac{C_{\text{cov}} \log(|\Pi|\rho^{-1}) \log(\delta^{-1})}{\delta\epsilon^2}\right). \quad (7)$$

This result shows that SFT-Sharpener, via Eq. (7), is minimax optimal in the sample-and-evaluate framework when δ is constant. In particular, the sample complexity bound in Eq. (7) matches the lower bound in Theorem 3.1 up to polynomial dependence on δ and logarithmic factors. Whether the $1/\delta$ factor in Eq. (7) can be removed is an interesting question, but—as discussed in Section 3.2—the regime $\delta = 1/2$ is most meaningful for autoregressive language modeling, rendering such discussion moot.

Remark 4.1 (On realizability and coverage). *Realizability assumptions such as Assumption 4.1 (which asserts that the class Π is powerful enough to model the distribution of the best-of- N responses) are standard in learning theory (Agarwal et al., 2019; Foster & Rakhlin, 2023), though certainly non-trivial (see Appendix D for a natural example where they may not hold). The coverage assumption, while also standard, when combined with the hypothesis that high-likelihood responses are desirable, suggests that π_{base} generates high-quality responses with reasonable probability. In general, doing so may require leveraging non-trivial serial computation at inference time via procedures such as Chain-of-Thought (Wei et al., 2022). Although recent work shows that such serial computation cannot be amortized (Li et al., 2024; Malach, 2023), SFT-Sharpener instead amortizes the parallel computation of best-of- N sampling, and thus has different representational considerations.*

Benefits of adaptive sampling. SFT-Sharpener is optimal in the sample-and-evaluate framework, but we show in Appendix C that a variant which selects the number of responses adaptively based on the prompt x can bypass this lower bound, improving the ϵ -dependence in Eq. (7) from $\frac{1}{\epsilon^2}$ to $\frac{1}{\epsilon}$.

Empirical Validation. In Appendix E, we empirically investigate the benefits of best-of- N on a variety of model-dataset pairs. Our results are summarized in Table 1 and Figs. 7 and 8, and broadly show that the benefits incurred through the inference-time sharpening described above can be, to a certain extent, amortized into training time.

4.2 ANALYSIS OF RLHF-Sharpener

We now turn our attention to theoretical guarantees for the RLHF-Sharpener algorithm family, which uses tools from RL to optimize the self-reward function.

When specialized to maximum-likelihood sharpening, the RL objective used by RLHF-Sharpener takes the form $\hat{\pi} \approx \arg \max_{\pi \in \Pi} \{\mathbb{E}_{\pi}[\log \pi_{\text{base}}(y | x)] - \beta D_{\text{KL}}(\pi \| \pi_{\text{base}})\}$ for $\beta > 0$. The exact optimizer $\pi_{\beta}^* = \arg \max_{\pi \in \Pi} \{\mathbb{E}_{\pi}[\log \pi_{\text{base}}(y | x)] - \beta D_{\text{KL}}(\pi \| \pi_{\text{base}})\}$ for this objective has the form $\pi_{\beta}^*(y | x) \propto \pi_{\text{base}}^{1+\beta^{-1}}(y | x)$, which converges to a sharpened model (per Definition 3.1) as $\beta \rightarrow 0$.

The key challenge we encounter in this section is the mismatch between the RL reward $\log \pi_{\text{base}}(y | x)$ and the sharpening desideratum $\hat{\pi}(\mathbf{y}^*(x) | x)$. For example, suppose a unique argmax—say, $y^*(x)$ —and second-to-argmax—say, $y'(x)$ —are nearly as likely under π_{base} . Then the RL reward $\mathbb{E}_{\hat{\pi}}[\log \pi_{\text{base}}(y | x)]$ must be optimized to extremely high precision before $\hat{\pi}$ can be guaranteed to distinguish the two. To quantify this effect, we introduce a *margin condition*.

Assumption 4.2 (Margin). *For a margin parameter $\gamma_{\text{margin}} > 0$, the base model π_{base} satisfies*

$$\max_{y \in \mathcal{Y}} \pi_{\text{base}}(y | x) \geq (1 + \gamma_{\text{margin}}) \cdot \pi_{\text{base}}(y' | x) \quad \forall y' \notin \mathbf{y}^*(x), \quad \forall x \in \text{supp}(\mu).$$

SFT-Sharpener does not suffer from the pathology in the example above, because once $y^*(x)$ and $y'(x)$ are drawn in a batch of N responses, we have $y_i^{\text{BoN}} = y^*(x_i)$ regardless of margin. However, as we shall show in Section 4.2.2, the RLHF-Sharpener algorithm is amenable to online exploration, which may improve dependence on other problem parameters.

⁶We focus on finite classes for simplicity, following a convention in reinforcement learning theory (Agarwal et al., 2019; Foster & Rakhlin, 2023), but our results extend to infinite classes through standard arguments.

4.2.1 GUARANTEES FOR RLHF-Sharpener WITH DIRECT PREFERENCE OPTIMIZATION

The first of our theoretical results for RLHF-Sharpener takes an offline reinforcement learning approach, whereby we implement Eq. (3) using a reward-based variant of Direct Preference Optimization (DPO) (Rafailov et al., 2023; Gao et al., 2024). Let $\mathcal{D}_{\text{pref}} = \{(x, y, y')\}$ be a dataset of n examples sampled via $x \sim \mu, y, y' \sim \pi_{\text{base}}(y | x)$. For a parameter $\beta > 0$, we solve $\hat{\pi} \in \arg \min_{\pi \in \Pi}$

$$\sum_{(x, y, y') \in \mathcal{D}_{\text{pref}}} \left(\beta \log \frac{\pi(y | x)}{\pi_{\text{base}}(y | x)} - \beta \log \frac{\pi(y' | x)}{\pi_{\text{base}}(y' | x)} - (\log \pi_{\text{base}}(y | x) - \log \pi_{\text{base}}(y' | x)) \right)^2. \quad (8)$$

Assumptions. Per Rafailov et al. (2023), the solution to Eq. (8) coincides with that of Eq. (2) asymptotically. To provide finite-sample guarantees, we make a number of statistical assumptions. First, we make a natural realizability assumption (e.g., Zhu et al. (2023); Xie et al. (2024)).

Assumption 4.3 (Realizability). *The model class Π satisfies $\pi_{\beta}^* \in \Pi$.*⁷

Next, we define two concentrability coefficients for a model π :

$$\mathcal{C}_{\pi} = \mathbb{E}_{\pi} \left[\frac{\pi(y | x)}{\pi_{\text{base}}(y | x)} \right], \quad \text{and} \quad \mathcal{C}_{\pi/\pi'; \beta} := \mathbb{E}_{\pi} \left[\left(\frac{\pi(y | x)}{\pi'(y | x)} \right)^{\beta} \right]. \quad (9)$$

The following result shows that both coefficients are bounded for the KL-regularized model π_{β}^* .

Lemma 4.1. *The model π_{β}^* satisfies $\mathcal{C}_{\pi_{\beta}^*} \leq C_{\text{cov}}$ and $\mathcal{C}_{\pi_{\text{base}}/\pi_{\beta}^*; \beta} \leq |\mathcal{Y}|$.*

Motivated by this result, we assume the coefficients in Eq. (9) are bounded for all $\pi \in \Pi$.

Assumption 4.4 (Concentrability). *All $\pi \in \Pi$ satisfy $\mathcal{C}_{\pi} \leq C_{\text{conc}}$ for a parameter $C_{\text{conc}} \geq C_{\text{cov}}$, and $\mathcal{C}_{\pi_{\text{base}}/\pi; \beta} \leq C_{\text{loss}}$ for a parameter $C_{\text{loss}} \geq |\mathcal{Y}|$.*

Per Lemma 4.1, this assumption is consistent with Assumption 4.3 for reasonable bounds on C_{conc} and C_{loss} ; note that our sample complexity bounds will only incur logarithmic dependence on C_{loss} .

Main result. Our sample complexity guarantee for RLHF-Sharpener (via Eq. (8)) is as follows.

Theorem 4.2. *Let $\epsilon, \delta, \rho \in (0, 1)$ be given. Set $\beta \lesssim \gamma_{\text{margin}} \delta \epsilon$, and suppose that Assumptions 4.2 to 4.4 hold with parameters C_{conc} , C_{loss} , and $\gamma_{\text{margin}} > 0$. For an appropriate choice for n , the DPO algorithm (Eq. (8)) ensures that with probability at least $1 - \rho$, $\mathbb{P}_{x \sim \mu}[\hat{\pi}(\mathbf{y}^*(x) | x) \leq 1 - \delta] \leq \epsilon$, and has sample complexity*

$$m = \tilde{O} \left(\frac{C_{\text{conc}} \log^3(C_{\text{loss}} |\Pi| \rho^{-1})}{\gamma_{\text{margin}}^2 \delta^2 \epsilon^2} \right).$$

Compared to the guarantee for SFT-Sharpener, RLHF-Sharpener learns a sharpened model with the same dependence on the accuracy ϵ , but a worse dependence on δ ; as we primarily consider δ constant (cf. Proposition 3.1), we view this as relatively unimportant. We further remark that RLHF-Sharpener uses $N = 2$ responses per prompt, while SFT-Sharpener uses many ($N = 1/\epsilon$) responses (but fewer prompts). Other differences include:

- RLHF-Sharpener requires the margin condition in Assumption 4.2, and has sample complexity scaling with $\gamma_{\text{margin}}^{-1}$. We believe this dependence is fundamental for algorithms based on reinforcement learning, as it is needed to translate bounds on suboptimality with respect to the reward function $r_{\text{self}}(y | x) = \log \pi_{\text{base}}(y | x)$ (i.e., $\mathbb{E}_{x \sim \mu}[\max_{y \in \mathcal{Y}} \log \pi_{\text{base}}(y | x) - \mathbb{E}_{y \sim \hat{\pi}(x)}[\log \pi_{\text{base}}(y | x)]] \leq \epsilon$, the objective minimized by reinforcement learning) into bounds on the approximate sharpening error $\mathbb{P}_{x \sim \mu}[\hat{\pi}(\mathbf{y}^*(x) | x) \leq 1 - \delta]$.
- RLHF-Sharpener requires a bound on the uniform coverage parameter C_{conc} , which is larger than the parameter C_{cov} required by SFT-Sharpener in general. We expect that this assumption can be removed by incorporating pessimism in the vein of (Liu et al., 2024; Huang et al., 2024). Also, RLHF-Sharpener requires a bound on the parameter C_{loss} . This grants control over the range of the reward function $\log \pi_{\text{base}}(y | x)$, which can otherwise be unbounded. Since the dependence on C_{loss} is only logarithmic, we view this as a fairly mild assumption. Overall, the guarantee in Theorem 4.2 may be somewhat pessimistic in practice; it would be interesting if the result can be improved to match the sample complexity of SFT-Sharpener whenever γ_{margin} is held constant.

⁷See Remark 4.1 for a discussion of this assumption.

4.2.2 BENEFITS OF EXPLORATION

The sample complexity guarantees we have presented scale with the coverage parameter $C_{\text{cov}} = \mathbb{E}[1/\pi_{\text{base}}(\mathbf{y}^*(x)|x)]$, which is unavoidable in general in the sample-and-evaluate framework via our lower bound, [Theorem 3.1](#). Although C_{cov} is a problem-dependent parameter, in the worst case it can be as large as $|\mathcal{Y}|$ (which is exponential in sequence length for autoregressive models). Luckily, unlike SFT-Sharpener, the RLHF-Sharpener objective [\(3\)](#) is amenable to RL algorithms employing active exploration, leading to improved sample complexity when the class Π has additional structure.

Our below guarantees for RLHF-Sharpener replace the assumption of bounded coverage with boundedness of a structural parameter for the model class Π known as the “sequential extrapolation coefficient” (SEC) ([Xie et al., 2023; 2024](#)), which we denote by $\text{SEC}(\Pi)$. The formal definition is deferred to [Appendix J.2](#). Conceptually, $\text{SEC}(\Pi)$ may thought of as a generalization of the eluder dimension ([Russo & Van Roy, 2013; Jin et al., 2021](#)), and can always be bounded by the coverability coefficient of the model class ([Xie et al., 2024](#)). Beyond boundedness of the SEC, we require a bound on the range of the log-probabilities of π_{base} .

Assumption 4.5 (Bounded log-probabilities). *For all $\pi \in \Pi$, $(x, y) \in \mathcal{X} \times \mathcal{Y}$, $|\log \frac{1}{\pi_{\text{base}}(y|x)}| \leq R_{\text{max}}$.*

We expect that the dependence on R_{max} in our result can be replaced with $\log(C_{\text{loss}})$ ([Assumption 4.4](#)), but we omit this extension to simplify presentation as much as possible.

We appeal to (a slight modification of) XPO, an iterative language model alignment algorithm due to [Xie et al. \(2024\)](#). XPO is based on the objective in [Eq. \(8\)](#), but unlike DPO, incorporates a bonus term to encourage exploration to leverage **online** interaction. See [Appendix J.2](#) for a detailed overview.

Theorem 4.3 (Informal version of [Theorem J.2](#)). *Suppose that [Assumptions 4.2 and 4.5](#) hold with parameters $\gamma_{\text{margin}}, R_{\text{max}} > 0$, and that [Assumption 4.3](#) holds with $\beta = \gamma_{\text{margin}}/(2 \log(2|\mathcal{Y}|/\delta))$. For any $m \in \mathbb{N}$ and $\rho \in (0, 1)$, XPO ([Algorithm 1](#)), when configured appropriately, produces an (ϵ, δ) -sharpened model $\hat{\pi} \in \Pi$ with probability at least $1 - \rho$, and uses sample complexity $m = \tilde{O}((\gamma_{\text{margin}}\delta\epsilon)^{-2}\text{SEC}(\Pi) \cdot \log(|\Pi|\rho^{-1}))$.⁸*

The takeaway from [Theorem 4.3](#) is that there is no dependence on the coverage coefficient for π_{base} . Instead, the rate depends on the complexity of exploration, as governed by the sequential extrapolation coefficient $\text{SEC}(\Pi)$. We expect similar guarantees can be derived for other active exploration algorithms and complexity measures ([Jiang et al., 2017; Foster et al., 2021; Jin et al., 2021; Xie et al., 2023](#)).

5 CONCLUSION

We view our theoretical framework for sharpening as a starting point toward a foundational understanding of self-improvement that can guide the design and evaluation of algorithms. To this end, we raise several directions for future research.

- *Representation learning.* A conceptually appealing feature of our framework is that it is agnostic to the structure of the model under consideration, but an important direction for future work is to study the dynamics of self-improvement for specific models (e.g. transformers), and understand the representations these models learn under self-training.
- *Richer forms of self-reward.* Our theoretical results study the dynamics of self-training in a stylized framework where the model uses its own logits for self-reward. Empirical research on self-improvement leverages more sophisticated approaches (e.g. specific prompting techniques) ([Huang et al., 2022; Wang et al., 2022; Bai et al., 2022b; Pang et al., 2023; Yuan et al., 2024](#)) and it is important to understand when and how these forms of self-improvement are beneficial.

⁸Technically, [Algorithm 1](#) operates in a slight generalization of the sample-and-evaluate framework for accessing π_{base} ([Definition 3.2](#)), where the algorithm is allowed to query $\pi_{\text{base}}(y|x)$ for arbitrary x, y . We expect that our lower bound ([Theorem 3.1](#)) can be extended to this more general framework, in which case [Algorithm 1](#) is fundamentally using additional structure of Π (via the SEC) to avoid dependence on C_{cov} .

REFERENCES

- Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv:2404.14219*, 2024.
- Samira Abnar, Mostafa Dehghani, and Willem Zuidema. Transferring inductive biases through knowledge distillation. *arXiv preprint arXiv:2006.00555*, 2020.
- Alekh Agarwal, Peter L Bartlett, Pradeep Ravikumar, and Martin J Wainwright. Information-theoretic lower bounds on the oracle complexity of stochastic convex optimization. *IEEE Transactions on Information Theory*, 5(58):3235–3249, 2012.
- Alekh Agarwal, Daniel Hsu, Satyen Kale, John Langford, Lihong Li, and Robert Schapire. Taming the monster: A fast and simple algorithm for contextual bandits. In *International Conference on Machine Learning*, pp. 1638–1646, 2014.
- Alekh Agarwal, Nan Jiang, and Sham M Kakade. Reinforcement learning: Theory and algorithms. <https://rltheorybook.github.io/>, 2019. Version: January 31, 2022.
- Zeyuan Allen-Zhu and Yuanzhi Li. Towards understanding ensemble, knowledge distillation and self-distillation in deep learning. *arXiv preprint arXiv:2012.09816*, 2020.
- Afra Amini, Tim Vieira, and Ryan Cotterell. Variational best-of-n alignment. *arXiv preprint arXiv:2407.06057*, 2024.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv:2204.05862*, 2022a.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022b.
- Francisco Barahona. On the computational complexity of ising spin glass models. *Journal of Physics A: Mathematical and General*, 15(10):3241, 1982.
- Matthew James Beal. *Variational algorithms for approximate Bayesian inference*. University of London, University College London (United Kingdom), 2003.
- Emmanuel Bengio, Moksh Jain, Maksym Korablyov, Doina Precup, and Yoshua Bengio. Flow network based generative models for non-iterative diverse candidate generation. *Advances in Neural Information Processing Systems*, 34:27381–27394, 2021.
- Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, 57(1):289–300, 1995.
- Avrim Blum, Merrick Furst, Jeffrey Jackson, Michael Kearns, Yishay Mansour, and Steven Rudich. Weakly learning dnf and characterizing statistical query learning using fourier analysis. In *Proceedings of the twenty-sixth annual ACM symposium on Theory of computing*, pp. 253–262, 1994.
- 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.

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, 2020.
- Cristian Buciluă, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 535–541, 2006.
- Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. Self-play fine-tuning converts weak language models to strong language models. *arXiv preprint arXiv:2401.01335*, 2024.
- 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*, 2017.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv:2110.14168*, 2021.
- Stephen A Cook. The complexity of theorem-proving procedures. In *Proceedings of the third annual ACM symposium on Theory of computing*, pp. 151–158, 1971.
- Thomas M Cover. *Elements of information theory*. John Wiley & Sons, 1999.
- Rudrajit Das and Sujay Sanghavi. Understanding self-distillation in the presence of label noise. In *International Conference on Machine Learning*, pp. 7102–7140. PMLR, 2023.
- Rudrajit Das, Inderjit S Dhillon, Alessandro Epasto, Adel Javanmard, Jieming Mao, Vahab Mirrokni, Sujay Sanghavi, and Peilin Zhong. Retraining with predicted hard labels provably increases model accuracy. *arXiv preprint arXiv:2406.11206*, 2024.
- Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Bin Dong, Jikai Hou, Yiping Lu, and Zhihua Zhang. Distillation \approx early stopping? harvesting dark knowledge utilizing anisotropic information retrieval for overparameterized neural network. *arXiv preprint arXiv:1910.01255*, 2019.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv:2407.21783*, 2024.
- Ronen Eldan, Frederic Koehler, and Ofer Zeitouni. A spectral condition for spectral gap: fast mixing in high-temperature ising models. *Probability theory and related fields*, 182(3):1035–1051, 2022.
- Amir-massoud Farahmand, Csaba Szepesvári, and Rémi Munos. Error propagation for approximate policy and value iteration. *Advances in Neural Information Processing Systems*, 2010.
- Vitaly Feldman. A complete characterization of statistical query learning with applications to evolvability. *Journal of Computer and System Sciences*, 78(5):1444–1459, 2012.
- Vitaly Feldman. A general characterization of the statistical query complexity. In *Conference on Learning Theory*, pp. 785–830. PMLR, 2017.
- Dylan J Foster and Alexander Rakhlin. Foundations of reinforcement learning and interactive decision making. *arXiv preprint arXiv:2312.16730*, 2023.
- Dylan J Foster, Sham M Kakade, Jian Qian, and Alexander Rakhlin. The statistical complexity of interactive decision making. *arXiv preprint arXiv:2112.13487*, 2021.

- Spencer Frei, Difan Zou, Zixiang Chen, and Quanquan Gu. Self-training converts weak learners to strong learners in mixture models. In *International Conference on Artificial Intelligence and Statistics*, pp. 8003–8021. PMLR, 2022.
- Tommaso Furlanello, Zachary Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar. Born again neural networks. In *International conference on machine learning*, pp. 1607–1616. PMLR, 2018.
- Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In *International Conference on Machine Learning*, pp. 10835–10866. PMLR, 2023.
- Zhaolin Gao, Jonathan D Chang, Wenhao Zhan, Owen Oertell, Gokul Swamy, Kianté Brantley, Thorsten Joachims, J Andrew Bagnell, Jason D Lee, and Wen Sun. REBEL: Reinforcement learning via regressing relative rewards. *arXiv:2404.16767*, 2024.
- Samuel Gershman and Noah Goodman. Amortized inference in probabilistic reasoning. In *Proceedings of the annual meeting of the cognitive science society*, number 36, 2014.
- Google. Palm 2 technical report. *arXiv:2305.10403*, 2023.
- Akhilesh Gotmare, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. A closer look at deep learning heuristics: Learning rate restarts, warmup and distillation. In *International Conference on Learning Representations*, 2019.
- 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.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv:2009.03300*, 2020.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Edward J Hu, Moksh Jain, Eric Elmoznino, Younesse Kaddar, Guillaume Lajoie, Yoshua Bengio, and Nikolay Malkin. Amortizing intractable inference in large language models. *arXiv preprint arXiv:2310.04363*, 2023.
- 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 overparameterization via chi-squared preference optimization. *arXiv:2407.13399*, 2024.
- Jiixin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. Large language models can self-improve. *arXiv preprint arXiv:2210.11610*, 2022.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv:2310.06825*, 2023.
- Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Contextual decision processes with low Bellman rank are PAC-learnable. In *International Conference on Machine Learning*, pp. 1704–1713, 2017.
- Chi Jin, Qinghua Liu, and Sobhan Miryoosefi. Bellman eluder dimension: New rich classes of RL problems, and sample-efficient algorithms. *Neural Information Processing Systems*, 2021.
- Richard M Karp. *Reducibility among combinatorial problems*. Springer, 1972.

- Michael Kearns. Efficient noise-tolerant learning from statistical queries. *Journal of the ACM (JACM)*, 45(6):983–1006, 1998.
- Scott Kirkpatrick, C Daniel Gelatt Jr, and Mario P Vecchi. Optimization by simulated annealing. *science*, 220(4598):671–680, 1983.
- Leonid Anatolevich Levin. Universal sequential search problems. *Problemy peredachi informatsii*, 9(3):115–116, 1973.
- Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of thought empowers transformers to solve inherently serial problems. *arXiv:2402.12875*, 2024.
- Zhihan Liu, Miao Lu, Shenao Zhang, Boyi Liu, Hongyi Guo, Yingxiang Yang, Jose Blanchet, and Zhaoran Wang. Provably mitigating overoptimization in rlhf: Your sft loss is implicitly an adversarial regularizer. *arXiv:2405.16436*, 2024.
- László Lovász and Santosh Vempala. Fast algorithms for logconcave functions: Sampling, rounding, integration and optimization. In *2006 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS'06)*, pp. 57–68. IEEE, 2006.
- Yi-An Ma, Yuansi Chen, Chi Jin, Nicolas Flammarion, and Michael I Jordan. Sampling can be faster than optimization. *Proceedings of the National Academy of Sciences*, 116(42):20881–20885, 2019.
- Eran Malach. Auto-regressive next-token predictors are universal learners. *arXiv:2309.06979*, 2023.
- Clara Meister, Tim Vieira, and Ryan Cotterell. If beam search is the answer, what was the question? *arXiv preprint arXiv:2010.02650*, 2020.
- Hossein Mobahi, Mehrdad Farajtabar, and Peter Bartlett. Self-distillation amplifies regularization in hilbert space. *Advances in Neural Information Processing Systems*, 33:3351–3361, 2020.
- 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. *arXiv preprint arXiv:2310.17022*, 2023.
- Arkadii Nemirovski, David Borisovich Yudin, and Edgar Ronald Dawson. Problem complexity and method efficiency in optimization. 1983.
- OpenAI. Gpt-4 technical report. *arXiv:2303.08774*, 2023.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 2022.
- Alizée Pace, Jonathan Mallinson, Eric Malmi, Sebastian Krause, and Aliaksei Severyn. West-of-n: Synthetic preference generation for improved reward modeling. *arXiv preprint arXiv:2401.12086*, 2024.
- Jing-Cheng Pang, Pengyuan Wang, Kaiyuan Li, Xiong-Hui Chen, Jiacheng Xu, Zongzhang Zhang, and Yang Yu. Language model self-improvement by reinforcement learning contemplation. *arXiv preprint arXiv:2305.14483*, 2023.
- Divyansh Pareek, Simon S Du, and Sewoong Oh. Understanding the gains from repeated self-distillation. *arXiv preprint arXiv:2407.04600*, 2024.
- Hieu Pham, Zihang Dai, Qizhe Xie, and Quoc V Le. Meta pseudo labels. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11557–11568, 2021.
- Yuxiao Qu, Tianjun Zhang, Naman Garg, and Aviral Kumar. Recursive introspection: Teaching language model agents how to self-improve. *arXiv preprint arXiv:2407.18219*, 2024.

- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 2023.
- Maxim Raginsky and Alexander Rakhlin. Information-based complexity, feedback and dynamics in convex programming. *IEEE Transactions on Information Theory*, 57(10):7036–7056, 2011.
- Mamshad Nayeem Rizve, Kevin Duarte, Yogesh S Rawat, and Mubarak Shah. In defense of pseudo-labeling: An uncertainty-aware pseudo-label selection framework for semi-supervised learning. *arXiv preprint arXiv:2101.06329*, 2021.
- Daniel Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of optimistic exploration. In *Advances in Neural Information Processing Systems*, pp. 2256–2264, 2013.
- Abulhair Saparov and He He. Language models are greedy reasoners: A systematic formal analysis of chain-of-thought. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=qFVVBzXxR2V>.
- Igal Sason and Sergio Verdú. f -divergence inequalities. *IEEE Transactions on Information Theory*, 62(11):5973–6006, 2016.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 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.
- Max Simchowitz, Kevin Jamieson, and Benjamin Recht. The simulator: Understanding adaptive sampling in the moderate-confidence regime. In *Conference on Learning Theory*, pp. 1794–1834. PMLR, 2017.
- Mohit Singh and Nisheeth K Vishnoi. Entropy, optimization and counting. In *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, pp. 50–59, 2014.
- 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, J Andrew Bagnell, and Wen Sun. Understanding preference fine-tuning through the lens of coverage. *arXiv:2406.01462*, 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.
- Kevin Swersky, Yulia Rubanova, David Dohan, and Kevin Murphy. Amortized bayesian optimization over discrete spaces. In *Conference on Uncertainty in Artificial Intelligence*, pp. 769–778. PMLR, 2020.
- Kunal Talwar. Computational separations between sampling and optimization. *Advances in neural information processing systems*, 32, 2019.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. *arXiv:2307.09288*, 2023.

- Joseph F Traub, Grzegorz W Wasilkowski, and Henryk Woźniakowski. Information-based complexity. 1988.
- S. A. van de Geer. *Empirical Processes in M-Estimation*. Cambridge University Press, 2000.
- Ziyu Wan, Xidong Feng, Muning Wen, Stephen Marcus McAleer, Ying Wen, Weinan Zhang, and Jun Wang. Alphazero-like tree-search can guide large language model decoding and training. *Forty-first International Conference on Machine Learning*, 2024.
- Tianlu Wang, Ilia Kulikov, Olga Golovneva, Ping Yu, Weizhe Yuan, Jane Dwivedi-Yu, Richard Yuanzhe Pang, Maryam Fazel-Zarandi, Jason Weston, and Xian Li. Self-taught evaluators. *arXiv preprint arXiv:2408.02666*, 2024.
- Xuezhi Wang and Denny Zhou. Chain-of-thought reasoning without prompting. *arXiv preprint arXiv:2402.10200*, 2024.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560*, 2022.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Wing Hung Wong and Xiaotong Shen. Probability inequalities for likelihood ratios and convergence rates of sieve mles. *The Annals of Statistics*, 1995.
- Tianhao Wu, Weizhe Yuan, Olga Golovneva, Jing Xu, Yuandong Tian, Jiantao Jiao, Jason Weston, and Sainbayar Sukhbaatar. Meta-rewarding language models: Self-improving alignment with llm-as-a-meta-judge. *arXiv preprint arXiv:2407.19594*, 2024a.
- Yangzhen Wu, Zhiqing Sun, Shanda Li, Sean Welleck, and Yiming Yang. An empirical analysis of compute-optimal inference for problem-solving with language models. *arXiv preprint arXiv:2408.00724*, 2024b.
- Yue Wu, Zhiqing Sun, Huizhuo Yuan, Kaixuan Ji, Yiming Yang, and Quanquan Gu. Self-play preference optimization for language model alignment. *arXiv preprint arXiv:2405.00675*, 2024c.
- Tengyang Xie and Nan Jiang. Q* approximation schemes for batch reinforcement learning: A theoretical comparison. In *Conference on Uncertainty in Artificial Intelligence*, 2020.
- Tengyang Xie, Dylan J Foster, Yu Bai, Nan Jiang, and Sham M Kakade. The role of coverage in online reinforcement learning. In *The Eleventh International Conference on Learning Representations*, 2023.
- Tengyang Xie, Dylan J Foster, Akshay Krishnamurthy, Corby Rosset, Ahmed Awadallah, and Alexander Rakhlin. Exploratory preference optimization: Harnessing implicit Q*-approximation for sample-efficient rlhf. *arXiv:2405.21046*, 2024.
- Wei Xiong, Hanze Dong, Chenlu Ye, Han Zhong, Nan Jiang, and Tong Zhang. Gibbs sampling from human feedback: A provable KL-constrained framework for RLHF. *arXiv:2312.11456*, 2023.
- Joy Qiping Yang, Salman Salamatian, Ziteng Sun, Ananda Theertha Suresh, and Ahmad Beirami. Asymptotics of language model alignment. *arXiv preprint arXiv:2404.01730*, 2024.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- 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:2402.07314*, 2024.
- Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. Self-rewarding language models. *arXiv preprint arXiv:2401.10020*, 2024.

- Andrea Zanette, Martin J Wainwright, and Emma Brunskill. Provable benefits of actor-critic methods for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 2021.
- Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.
- Tong Zhang. From ϵ -entropy to KL-entropy: Analysis of minimum information complexity density estimation. *The Annals of Statistics*, 34(5):2180–2210, 2006.
- Stephen Zhao, Rob Brekelmans, Alireza Makhzani, and Roger Baker Grosse. Probabilistic inference in language models via twisted sequential monte carlo. *International Conference on Machine Learning*, pp. 60704–60748, 2024.
- 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, 2024.
- Banghua Zhu, Michael Jordan, and Jiantao Jiao. Principled reinforcement learning with human feedback from pairwise or k-wise comparisons. In *International Conference on Machine Learning*, pp. 43037–43067. PMLR, 2023.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

Contents of Appendix

I	Additional Discussion and Results	19
A	Detailed Discussion of Related Work	19
B	Guarantees for Inference-Time Sharpening	20
C	Guarantees for SFT-Sharpening with Adaptive Sampling	21
D	Computational and Representational Challenges in Sharpening	22
	D.1 Computational Challenges	23
	D.2 Representational Challenges	23
E	Additional Experiments and Details	25
	E.1 Inference-time validation experiments	28
	E.2 Experiments with other self reward functions	29
	E.3 Effect of SFT-Sharpening	29
II	Proofs	34
F	Preliminaries	34
	F.1 Guarantees for Approximate Maximizers	34
	F.2 Technical Tools	34
G	Proofs from Section 3.1	36
H	Proofs from Section 3.3	36
I	Proofs from Section 4.1 and Appendix C	40
J	Proofs from Section 4.2	42
	J.1 Proof of Theorem 4.2	43
	J.2 Proof of Theorem 4.3 and [UNDEFINED]	47

Part I

Additional Discussion and Results

A DETAILED DISCUSSION OF RELATED WORK

In this section, we discuss related work in greater detail, including relevant works not already covered.

Self-improvement and self-training. Our work is most directly related to a growing body of empirical research that studies self-improvement/self-training for language models in a supervision-free setting in which there is no external feedback (Huang et al., 2022; Wang et al., 2022; Bai et al., 2022b; Pang et al., 2023), and takes a first step toward providing a theoretical understanding for these methods. This line of work is closely related to a body of research on “LLM-as-a-Judge” techniques and related work, which investigates approaches to designing self-reward functions r_{self} , often based on specific prompting techniques (Zheng et al., 2024; Yuan et al., 2024; Wu et al., 2024a; Wang et al., 2024).

There is a somewhat complementary line of research that develops algorithms based on self-training and self-play (Zelikman et al., 2022; Chen et al., 2024; Wu et al., 2024c; Qu et al., 2024), but leverages various forms of external feedback (e.g., positive examples for SFT or explicit reward signal). These methods typically outperform self-improvement methods, which do not use any external feedback (Zelikman et al., 2022). However, in many scenarios, obtaining external feedback can be costly or laborious; it may require collecting high-quality labeled/annotated data, rewriting examples in a formal language, etc. Thus, these methods are not directly comparable to methods based on self-improvement.

Lastly, we mention in passing that the self-improvement problem we study is related to a more classical line of research on *self-distillation* (Buciluă et al., 2006; Hinton et al., 2015; Devlin, 2018; Pham et al., 2021; Rizve et al., 2021), but this specific form of self-training has received limited investigation in the context of language modeling.

Alignment and RLHF. The specific algorithms for self-improvement/sharpening we study can be viewed as special cases of standard alignment algorithms, including classical RLHF methods (Christiano et al., 2017; Bai et al., 2022a; Ouyang et al., 2022), direct alignment (Rafailov et al., 2023), and (inference-time or training-time) best-of- N methods (Amini et al., 2024; Sessa et al., 2024; Gui et al., 2024; Pace et al., 2024). However, the maximum likelihood sharpening objective (2) used for our theoretical results has been relatively unexplored within the alignment literature.

Inference-time decoding. Many inference-time decoding strategies such as greedy/low-temperature decoding, beam-search (Meister et al., 2020), and chain-of-thought decoding (Wang & Zhou, 2024) can be viewed as instances of inference-time sharpening for specific choices of the self-reward function r_{self} . More sophisticated inference-time search strategies such as tree search and MCTS (Yao et al., 2024; Wan et al., 2024; Mudgal et al., 2023; Zhao et al., 2024) are also related, though this line of working frequently makes use of external reward signals or verification, which is somewhat complementary to our work.

Theoretical guarantees for self-training. On the theoretical side, current understanding of self-training is limited. One line of work, focusing on the *self-distillation* objective (Hinton et al., 2015) for binary classification and regression, aims to provide convergence guarantees for self-training in stylized setups such as linear models (Mobahi et al., 2020; Das & Sanghavi, 2023; Das et al., 2024; Pareek et al., 2024), with Allen-Zhu & Li (2020) giving guarantees for feedforward neural networks. Perhaps most closely related to our work is Frei et al. (2022), who show that self-training on a model’s pseudo-labels can amplify the margin for linear logistic regression. However, to the best of our knowledge, our work is the first to study self-training in a general framework that subsumes language modeling.

Our theoretical results for RLHF-Sharpener are also related to a recent body of work that provides sample complexity guarantees for alignment methods (Zhu et al., 2023; Xiong et al., 2023; Ye et al., 2024; Huang et al., 2024; Liu et al., 2024; Song et al., 2024; Xie et al., 2024), but our results leverage the unique structure of the maximum-likelihood sharpening self-reward function

$r_{\text{self}}(y | x) = \log \pi_{\text{base}}(y | x)$, and provide guarantees for the sharpening objective in [Definition 3.1](#) instead of the usual notion of reward suboptimality used in reinforcement learning theory.

Lastly, we mention that our results—particularly our *amortization* perspective on self-improvement—are related to recent work that studies fundamental representational advantages of allowing additional inference time ([Malach, 2023](#); [Li et al., 2024](#)). These work focus on truly sequential tasks, while our work focuses on the complementary question of amortizing *parallel* computation. Thus the representational implications are quite different.

Optimization versus sampling. The maximum-likelihood sharpening we introduce in [Section 3](#) connects the study of *self-improvement* to a large body of research in theoretical computer science on computational tradeoffs (e.g., separations and equivalences) for optimization and sampling ([Barahona, 1982](#); [Kirkpatrick et al., 1983](#); [Lovász & Vempala, 2006](#); [Singh & Vishnoi, 2014](#); [Ma et al., 2019](#); [Talwar, 2019](#); [Eldan et al., 2022](#)). On the one hand, this line of research highlights that there exist natural classes of distributions for which sampling is tractable, yet maximum likelihood optimization is intractable, and vice-versa. On the other hand, various works in this line of research also demonstrate *computational reductions* between optimization and sampling, whereby optimization can be reduced to sampling and vice-versa.

Our setting indeed includes natural model classes where one should not expect there to be a computational reduction from optimization ($\arg \max_{y \in \mathcal{Y}} \pi_{\text{base}}(y | x)$) to sampling ($y \sim \pi_{\text{base}}(\cdot | x)$), and hence inference-time sharpening is computationally intractable ([Proposition D.1](#)). Of course, coverage assumptions eliminate this intractability. For training-time sharpening (where the goal is to *amortize* across prompts by training a sharpened model, as formulated in [Section 3](#)) the obstacle in natural, concrete model classes is not just computational but in fact *representational* ([Proposition D.2](#)). Regarding the latter point, we note that while amortized Bayesian inference has received extensive investigation empirically ([Beal, 2003](#); [Gershman & Goodman, 2014](#); [Swersky et al., 2020](#); [Bengio et al., 2021](#); [Hu et al., 2023](#)), we are unaware of theoretical guarantees outside of this work.

B GUARANTEES FOR INFERENCE-TIME SHARPENING

In this section, we give theoretical guarantees for the inference-time best-of- N sampling algorithm for sharpening described in [Section 3.1](#), under the maximum-likelihood sharpening self-reward function $r_{\text{self}}(y | x; \pi_{\text{base}}) = \log \pi_{\text{base}}(y | x)$.

Recall that given a prompt $x \in \mathcal{X}$, the inference-time best-of- N sampling algorithm draws N responses $y_1, \dots, y_n \sim \pi_{\text{base}}(\cdot | x)$, then return the response $\hat{y} = \arg \max_{y_i} \log \pi_{\text{base}}(y_i | x)$. We show that this algorithm returns an approximate maximizer for the maximum-likelihood sharpening objective whenever the base policy π_{base} has sufficient coverage. Recall that for a parameter $\gamma \in [0, 1)$ we define

$$\mathbf{y}_{\gamma}^*(x) := \left\{ y \mid \pi_{\text{base}}(y | x) \geq (1 - \gamma) \cdot \max_{y \in \mathcal{Y}} \pi_{\text{base}}(y | x) \right\}$$

as the set of $(1 - \gamma)$ -approximate maximizers for $\log \pi_{\text{base}}(y | x)$.

Proposition B.1. *Let a prompt $x \in \mathcal{X}$ be given. For any $\rho \in (0, 1)$ and $\gamma \in [0, 1)$, as long as*

$$N \geq \frac{\log(\rho^{-1})}{\pi_{\text{base}}(\mathbf{y}_{\gamma}^*(x) | x)},$$

inference-time best-of- N sampling produces a response $\hat{y} \in \mathbf{y}_{\gamma}^(x)$ with probability at least $1 - \rho$.*

Proof of Proposition B.1. Fix a prompt $x \in \mathcal{X}$, failure probability $\rho \in (0, 1)$, and parameter $\gamma \in (0, 1)$.

By definition of the set $\mathbf{y}_{\gamma}^*(x)$, $\hat{y} \in \mathbf{y}_{\gamma}^*(x)$ if and only if there exists $i \in [N]$ such that $y_i \in \mathbf{y}_{\gamma}^*(x)$. The complement of this event, i.e., that $y_i \notin \mathbf{y}_{\gamma}^*(x)$ for all $i \in [N]$, has probability

$$\mathbb{P}(y_i \notin \mathbf{y}_{\gamma}^*(x), \forall i \in [N]) = (1 - \pi_{\text{base}}(\mathbf{y}_{\gamma}^*(x) | x))^N.$$

Rearranging the right-hand-side, we have

$$(1 - \pi_{\text{base}}(\mathbf{y}_{\gamma}^*(x) | x))^N = \exp\left(-N \log\left(\frac{1}{1 - \pi_{\text{base}}(\mathbf{y}_{\gamma}^*(x) | x)}\right)\right) \leq \exp(-N \cdot \pi_{\text{base}}(\mathbf{y}_{\gamma}^*(x) | x)),$$

since $\log(x) \geq 1 - \frac{1}{x}$ for $x > 0$, which implies that $\log\left(\frac{1}{1 - \pi_{\text{base}}(\mathbf{y}_\gamma^* | x)}\right) \geq \pi_{\text{base}}(\mathbf{y}_\gamma^* | x)$. Thus, as long as $N \geq \frac{\log(\rho^{-1})}{\pi_{\text{base}}(\mathbf{y}_\gamma^* | x)}$, we have

$$\mathbb{P}(y_i \notin \mathbf{y}_\gamma^*(x), \forall i \in [N]) \leq \exp(-N \cdot \pi_{\text{base}}(\mathbf{y}_\gamma^* | x)) \leq \exp(-\log(\rho^{-1})) = \rho.$$

We conclude that with probability at least $1 - \rho$, there exists $i \in [N]$ such that $y_i \in \mathbf{y}_\gamma^*(x)$, and $\hat{y} \in \mathbf{y}_\gamma^*(x)$ as a result. \square

C GUARANTEES FOR SFT-SHARPENING WITH ADAPTIVE SAMPLING

SFT-Sharpener is a simple and natural self-training scheme, and converges to a sharpened policy as $n, N \rightarrow \infty$. However, using a fixed response sample size N may be wasteful for prompts where the model is confident. To this end, in this section we introduce and analyze, a variant of SFT-Sharpener based on *adaptive sampling*, which adjusts the number of sampled responses adaptively.

Algorithm. We present the adaptive SFT-Sharpener algorithm only for the special case of the maximum-likelihood sharpening self-reward. Let a *stopping parameter* $\mu > 0$ be given. For $x_i \in \mathcal{X}$, and $y_{i,1}, y_{i,2} \dots \sim \pi_{\text{base}}(\cdot | x_i)$, define a stopping time (e.g., [Benjamini & Hochberg \(1995\)](#)) via:

$$N_\mu(x_i) := \inf \left\{ k : \frac{1}{\max_{1 \leq j \leq k} \pi_{\text{base}}(y_{i,j} | x_i)} \leq \frac{k}{\mu} \right\}. \quad (10)$$

The adaptive SFT-Sharpener algorithm computes adaptively sampled responses y_i^{AdaBoN} via

$$y_i^{\text{AdaBoN}} \sim \arg \max \{ \log \pi_{\text{base}}(y_{i,j} | x_i) \mid y_{i,1}, \dots, y_{i,N_\mu(x_i)} \},$$

then trains the sharpened model through SFT:

$$\hat{\pi}^{\text{AdaBoN}} = \arg \max_{\pi \in \Pi} \sum_{i=1}^n \log \pi(y_i^{\text{AdaBoN}} | x_i).$$

Critically, by using scheme in [Eq. \(10\)](#), this algorithm can stop sampling responses for the prompt x_i if it becomes clear that the confidence is large.

Theoretical guarantee. We now show that adaptive SFT-Sharpener enjoys provable benefits over its non-adaptive counterpart through the dependence on the accuracy parameter $\epsilon > 0$.

Given $x \in \mathcal{X}$, and $y_1, y_2 \dots \sim \pi_{\text{base}}(x)$, let $N_\mu(x) := \inf \{ k : \frac{1}{\max_{1 \leq i \leq k} \pi_{\text{base}}(y_i | x)} \leq k/\mu \}$, and define a random variable $y^{\text{AdaBoN}}(x) \sim \arg \max \{ \log \pi_{\text{base}}(y_i | x) \mid y_1, \dots, y_{N_\mu(x)} \}$. Let $\pi_\mu^{\text{AdaBoN}}(x)$ denote the distribution over $y^{\text{AdaBoN}}(x)$. We make the following realizability assumption.

Assumption C.1. The model class Π satisfies $\pi_\mu^{\text{AdaBoN}} \in \Pi$.

Compared to SFT-Sharpener, we require a somewhat stronger coverage coefficient given by

$$\bar{C}_{\text{cov}} = \mathbb{E}_{x \sim \mu} \left[\frac{1}{\max_{y \in \mathcal{Y}} \pi_{\text{base}}(y | x)} \right].$$

This definition coincides with [Eq. \(6\)](#) when the arg-max response is unique, but is larger in general.

Our main theoretical guarantee for adaptive SFT-Sharpener is as follows.

Theorem C.1. Let $\delta, \rho \in (0, 1)$ be given. Set $\mu = \ln(2\delta^{-1})$, and assume [Assumption C.1](#) holds. Then with probability at least $1 - \rho$, the adaptive SFT-Sharpener algorithm has

$$\mathbb{P}_{x \sim \mu} [\hat{\pi}(\mathbf{y}^*(x) | x) \leq 1 - \delta] \lesssim \frac{\log(|\Pi|\rho^{-1})}{\delta n},$$

and has sample complexity $\mathbb{E}[n] = n \cdot \bar{C}_{\text{cov}} \log(\delta^{-1})$. Taking $n \gtrsim \frac{\log(|\Pi|\rho^{-1})}{\delta \epsilon}$ ensures that with probability at least $1 - \rho$,

$$\mathbb{P}_{x \sim \mu} [\hat{\pi}(\mathbf{y}^*(x) | x) \leq 1 - \delta] \leq \epsilon,$$

and gives total sample complexity

$$\mathbb{E}[m] = O\left(\frac{\bar{C}_{\text{cov}} \log(|\Pi| \rho^{-1}) \log(\delta^{-1})}{\delta \epsilon}\right).$$

Compared to the result for SFT-Sharpener in [Theorem 4.1](#), this shows that adaptive SFT-Sharpener achieves sample complexity scaling with $\frac{1}{\epsilon}$ instead of $\frac{1}{\epsilon^2}$. We believe the dependence on \bar{C}_{cov} for this algorithm is tight, as the adaptive stopping rule used in the algorithm can be overly conservative when $|\mathbf{y}^*(x)|$ is large.

A matching lower bound. We now prove a complementary lower bound, which shows that the ϵ -dependence in [Theorem C.1](#) is tight. To do so, we consider the following adaptive variant of the sample-and-evaluate framework.

Definition C.1 (Adaptive sample-and-evaluate framework). *In the Adaptive Sample-and-Evaluate framework, the learner is allowed to sample n prompts $x \sim \mu$, and sample an arbitrary, adaptively chosen number of samples $y_1, y_2, \dots \sim \pi_{\text{base}}(\cdot | x)$ before sampling a new prompt $x' \sim \mu$. In this framework we define sample complexity m as the total number of pairs (x, y) sampled by the algorithm, which is a random variable.*

Our main lower bound is as follows.

Theorem C.2 (Lower bound for sharpening under adaptive sampling). *Fix an integer $d \geq 1$ and parameters $\epsilon \in (0, 1)$ and $C \geq 1$. There exists a class of models Π such that (i) $\log |\Pi| \asymp d(1 + \log(C\epsilon^{-1}))$, (ii) $\sup_{\pi \in \Pi} C_{\text{cov}}(\pi) \lesssim C$, and (iii) $\mathbf{y}^\pi(x)$ is a singleton for all $\pi \in \Pi$, for which any sharpening algorithm $\hat{\pi}$ in the adaptive sample-and-evaluate framework that achieves $\mathbb{E}[\mathbb{P}_{x \sim \mu}[\hat{\pi}(\mathbf{y}^{\pi_{\text{base}}}(x) | x) > 1/2]] \geq 1 - \epsilon$ for all $\pi_{\text{base}} \in \Pi$ must collect a total number of samples $m = n \cdot N$ at least*

$$\mathbb{E}[m] \gtrsim \frac{C \log |\Pi|}{\epsilon \cdot (1 + \log(C\epsilon^{-1}))}.$$

[Theorem C.2](#) is a special case of a more general theorem, [Theorem 3.1'](#), which is stated and proven in [Appendix H](#).

D COMPUTATIONAL AND REPRESENTATIONAL CHALLENGES IN SHARPENING

In this section, we make several basic observations about the inherent computational and representational challenges of maximum-likelihood sharpening. First, in [Appendix D.1](#), we focus on computational challenges, and show that computing a sharpened response for a given prompt x can be computationally intractable in general, even when sampling $y \sim \pi_{\text{base}}(\cdot | x)$ can be performed efficiently. Then, in [Appendix D.2](#), we shift our focus to representational challenges, and show that even if π_{base} is an autoregressive model, the “sharpened” version of π_{base} may not be representable as an autoregressive model with the same architecture. These results motivate the statistical assumptions (coverage and realizability) made in our analysis of SFT-Sharpener and RLHF-Sharpener in [Section 4](#).

To make the results in this section precise, we work in perhaps the simplest special case of autoregressive language modelling, where the model class consists of *multi-layer linear softmax models*. Formally, let \mathcal{X} be the space of prompts, and let $\mathcal{Y} := \mathcal{V}^H$ be the space of responses, where \mathcal{V} is the vocabulary space and H is the horizon. For a collection of fixed/known d -dimensional feature mappings $\phi_h : \mathcal{X} \times \mathcal{V}^h \rightarrow \mathbb{R}^d$ and a norm parameter B , we define the model class $\Pi_{\phi, B, H}$ as the set of models

$$\pi_\theta(y_{1:H} | x) = \prod_{h=1}^H \pi_{\theta_h}(y_h | x, y_{1:h-1}) \quad (11)$$

where

$$\pi_\theta(y_h | x, y_{1:h-1}) \propto \exp(\langle \phi(x, y_{1:h}), \theta_h \rangle)$$

and $\theta = (\theta_1, \dots, \theta_H) \in (\mathbb{R}^d)^H$ is any tuple with $\|\theta_h\|_2 \leq B$ for all $h \in [H]$.

D.1 COMPUTATIONAL CHALLENGES

Given query access to ϕ , for any given parameter vector θ and prompt x , *sampling* from a linear softmax model π_θ (Eq. (11)) is computationally tractable, since it only requires time $\text{poly}(H, |\mathcal{V}|, d)$. Similarly, *evaluating* $\pi_\theta(y_{1:H} \mid x)$ for given prompt x and response $y_{1:H}$ is computationally tractable. However, the following proposition shows that computing the sharpened response $\arg \max_{y_{1:H} \in \mathcal{V}^H} \pi_\theta(y_{1:H} \mid x)$ for a given parameter θ and response x is NP-hard. Hence, even inference-time sharpening is computationally intractable in the worst case.

Proposition D.1. *Set $\mathcal{X} = \{\perp\}$ and $\mathcal{V} = \{-1, 1\}$. Set $d = d(H) := H + H^2 + H^3$. Identifying $[d]$ with $[H] \sqcup [H]^2 \sqcup [H]^3$, we define $\phi_h : \mathcal{X} \times \mathcal{V}^h \rightarrow \mathbb{R}^d$ by $\phi_h(\perp, y_{1:h})_i = y_i$ and $\phi_h(\perp, y_{1:h})_{(i,j)} = y_i y_j$ and $\phi_h(\perp, y_{1:h})_{(i,j,k)} = y_i y_j y_k$. There is a function $B(H) \leq \text{poly}(H)$ such that the following problem is NP-hard: given $\theta = (\theta_1, \dots, \theta_H)$ with $\max_{h \in [H]} \|\theta_h\|_2 \leq B(H)$, compute any element of $\arg \max_{y_{1:H} \in \mathcal{V}^H} \pi_\theta(y_{1:H} \mid x)$.*

Note that our results in Section 4 and Appendix B bypass this hardness through the assumption that the coverage parameter C_{cov} is bounded.

Proof of Proposition D.1. Fix H and recall that $d(H) = H + H^2 + H^3$. We define three collection of basis vectors: $\{e_h\}_{h \in [H]}$ cover the first H coordinates, $\{e_{(h,h')}\}_{h,h' \in [H]^2}$ cover the next H^2 coordinates, and $\{e_{(h,h',h'')}\}_{h,h',h'' \in [H]^3}$ cover the last H^3 coordinates. Suppose we define $\theta_1, \dots, \theta_{H-2} = 0$, so that $\pi_\theta(y_h \mid x, y_{1:h-1}) = 1/2$ for all $1 \leq h \leq H-2$. Define $\theta_{H-1} = \sum_{1 \leq i,j \leq H-2} J_{ij} e_{(i,j,H-1)}$ for a matrix $J \in \mathbb{R}^{(H-2) \times (H-2)}$ to be specified later, and define $\theta_H = \frac{B}{2}(e_{(H-1,H)} + e_H)$. Then $2^{H-2} \cdot \pi_\theta(y_{1:H} \mid \perp) \leq 1/2$ for any $y_{1:H}$ with $y_{H-1} = -1$ or $y_H = -1$, since this implies that $\pi_{\theta_H}(y_H \mid \perp, y_{1:H-1}) \leq 1/2$. Meanwhile, for any $y_{1:H}$ with $y_{H-1} = y_H = 1$, we have

$$2^{H-2} \cdot \pi_\theta(y_{1:H} \mid \perp) = \frac{\exp\left(\sum_{i,j \leq H-2} J_{ij} y_i y_j\right)}{\exp\left(\sum_{i,j \leq H-2} J_{ij} y_i y_j\right) + \exp\left(-\sum_{i,j \leq H-2} J_{ij} y_i y_j\right)} \cdot \frac{\exp(B)}{\exp(B) + \exp(-B)}.$$

Let G be any graph on vertex set $[H-2]$ and let $J = -A(G)$ where $A(G)$ is the adjacency matrix of G . Then among $y_{1:H}$ with $y_{H-1} = y_H = 1$, $2^{H-2} \cdot \pi_\theta(y_{1:H} \mid \perp)$ is maximized when $y_{1:H-2}$ corresponds to a max-cut in G . If G has an odd number of edges, then some max-cut removes strictly more than half of the edges, and for the corresponding sequence $y_{1:H}$ we have $2^{H-2} \cdot \pi_\theta(y_{1:H} \mid \perp) \geq (1/2 + \Omega(1)) \cdot (1 - \exp(-\Omega(B)))$, which is greater than $1/2$ when we take $B := H$ and H is sufficiently large. Thus, computing $\arg \max_{y_{1:H} \in \mathcal{V}^H} \pi_\theta(y_{1:H} \mid \perp)$ yields a max-cut of G . It is well-known that computing a max-cut in a graph is NP-hard, and the assumption that G has an odd number of edges is without loss of generality. \square

D.2 REPRESENTATIONAL CHALLENGES

To give provable guarantees for our sharpening algorithms, we required certain *realizability* assumptions, which in particular posited that the model class actually contains a “sharpened” version of π_{base} (Assumptions 4.1 and 4.3). In the simple example of a *single-layer* linear softmax model classes (corresponding to $H = 1$ in the above definition), Assumption 4.3 is in fact satisfied, and the sharpened model can be obtained by increasing the temperature of π_{base} . However, multi-layer linear softmax models with $H \gg 1$ better capture autoregressive language models. The following proposition shows that as soon as $H \geq 2$, multi-layer linear softmax model classes may not be closed under sharpening. This illustrates a potential drawback of training-time sharpening compared to inference-time sharpening, which requires no realizability assumptions. It also provides a simple example where greedy decoding does not yield a sequence-level arg-max response (since increasing temperature in a multi-layer softmax model class exactly converges to the greedy decoding).

Proposition D.2. *Let $\mathcal{X} = \{\perp\}$, $\mathcal{V} = [n]$, and $H = d = 2$. For any n sufficiently large, there is a multi-layer linear softmax policy class $\Pi_{\phi,B,H}$ and a policy $\pi_{\text{base}} \in \Pi_{\phi,B,H}$ such that $y_{1:H}^* := \arg \max_{y_{1:H} \in \mathcal{V}^H} \pi_\theta(y_{1:H} \mid \perp)$ is unique but for all $B' > B$ and $\pi \in \Pi_{\phi,B',H}$, it holds that $\pi(y_{1:H}^* \mid \perp) \leq 1/2$.*

Proof of Proposition D.2. Throughout, we omit the dependence on the prompt \perp for notational clarity. Since $H = 2$, the model class consists of models π_θ of the form

$$\pi_\theta(a) = \pi_{\theta_1}(y_1)\pi_{\theta_2}(y_2 | y_1) = \frac{\exp(\langle \phi_1(y_1), \theta_1 \rangle)}{Z_{\theta_1}} \frac{\exp(\langle \phi_2(y_{1:2}), \theta_2 \rangle)}{Z_{\theta_2}(y_1)} \quad (12)$$

for $Z_{\theta_1} := \sum_{y_1 \in \mathcal{V}} \exp(\langle \phi_1(y_1), \theta_1 \rangle)$ and $Z_{\theta_2}(y_1) := \sum_{y_2 \in \mathcal{V}} \exp(\langle \phi_2(y_{1:2}), \theta_2 \rangle)$.

Define ϕ_1 by:

$$\phi_1(i) = \begin{cases} e_1 & \text{if } i = 1 \\ e_1 & \text{if } i = 2 \\ e_2 & \text{if } i \geq 3 \end{cases}.$$

Define ϕ_2 by:

$$\phi_2(i, j) = \begin{cases} e_1 & \text{if } i = 2, j = 1 \\ e_2 & \text{if } i = 2, j \neq 1 \\ 0 & \text{if } i \neq 2 \end{cases}.$$

Define $\pi_{\text{base}} := \pi_{\theta^*}$ where $\theta_1^* := \theta_2^* := B \cdot e_1$ for a parameter $B \geq \log(n)$. Then $\pi_{\text{base}}(1) = \pi_{\text{base}}(2)$ and $\pi_{\text{base}}(i) \leq e^{-B} \pi_{\text{base}}(2)$ for all $i \in \{3, \dots, n\}$. Moreover, $\pi_{\text{base}}(\cdot | i) = \text{Unif}([n])$ for all $i \neq 2$, and $\pi_{\text{base}}(j | 2) \leq e^{-B} \pi_{\text{base}}(1 | 2)$ for all $j \neq 1$. Thus,

$$\pi_{\text{base}}(2, 1) = \pi_{\text{base}}(2)\pi_{\text{base}}(1 | 2) \geq \frac{1}{2 + (n-2)e^{-B}} \cdot \frac{1}{1 + (n-1)e^{-B}} \geq \Omega(1)$$

whereas $\pi_{\text{base}}(i, j) = O(1/n)$ for all $(i, j) \neq (2, 1)$. Thus, $(2, 1)$ is the sequence-level argmax for sufficiently large n . However, for any π_θ of the form described in Eq. (12), we have

$$\pi_\theta(2, 1) \leq \pi_\theta(2) \leq \frac{\pi_\theta(2)}{\pi_\theta(1) + \pi_\theta(2)} = \frac{1}{2}$$

since $\phi(1) = \phi(2)$. This means that there is no B' for which $\Pi_{\phi, B', H}$ contains an (ϵ, δ) -sharpened policy for π_{base} for any $\delta > 1/2$. \square

BoN-Norm: % Lift over Greedy							Majority: % Lift over Greedy						
Phi3 (Mini)	4.4	-0.4	7.0	-1.2	7.2	5.2	Phi3 (Mini)	19.7	5.1	0.3	7.1	12.5	8.7
Phi3.5 (Mini)	0.8	5.0	1.5	6.0	4.9	0.7	Phi3.5 (Mini)	19.4	8.7	1.2	11.1	6.9	1.5
Phi3 (Small)	2.1	7.3	-4.6	3.3	2.4	-11.8	Phi3 (Small)	17.1	11.4	-14.4	2.5	7.9	-9.4
Phi3 (Medium)	0.1	-1.7	0.5	0.2	6.5	-1.3	Phi3 (Medium)	16.0	6.1	5.0	3.6	13.2	1.3
Mistral-7B	28.6	17.1	2.8	7.6	25.9	4.0	Mistral-7B	72.2	48.8	38.5	19.8	19.5	10.5
Llama3.2-3B	4.3	1.3	3.6	-12.7	79.3	20.3	Llama3.2-3B	17.5	11.6	9.6	11.5	99.3	59.7
GPT-3.5	3.4	1.2	-5.9	-1.8	7.1	-7.9	GPT-3.5	35.5	18.2	4.6	7.0	11.0	7.4
	MATH	GSM8K	ProntoQA	Bio	Phys	Chem		MATH	GSM8K	ProntoQA	Bio	Phys	Chem

(a)

Pass@50: Accuracy (%)							Greedy: Accuracy (%)						
Phi3 (Mini)	96.6	97.3	98.8	99.2	99.9	98.1	Phi3 (Mini)	66.0	87.1	50.4	80.6	65.7	52.0
Phi3.5 (Mini)	97.0	96.8	90.9	97.5	98.8	96.7	Phi3.5 (Mini)	67.6	84.4	50.8	77.1	68.6	55.0
Phi3 (Small)	97.9	97.3	90.3	98.8	97.2	96.3	Phi3 (Small)	72.3	79.3	59.4	85.4	74.5	66.0
Phi3 (Medium)	98.4	98.3	83.3	99.0	99.0	94.7	Phi3 (Medium)	73.4	86.7	47.3	88.2	70.6	60.0
Mistral-7B	81.9	94.0	99.7	98.8	98.2	98.0	Mistral-7B	23.0	46.9	50.0	61.8	36.3	42.0
Llama3.2-3B	93.8	95.8	100.0	99.6	99.8	99.2	Llama3.2-3B	58.2	76.6	47.7	61.8	14.7	30.0
GPT-3.5	96.0	96.6	95.1	98.5	99.9	99.4	GPT-3.5	55.9	70.3	49.6	68.8	51.0	53.0
	MATH	GSM8K	ProntoQA	Bio	Phys	Chem		MATH	GSM8K	ProntoQA	Bio	Phys	Chem

(c)

(b)

(d)

Figure 2: Performance of alternative decoding schemes beyond BoN. Percent improvement of accuracy over greedy decoding for self-improvement with length-normalized log probability (a) and majority voting (b), with both demonstrating efficacy on a range of model-task pairs. (c) Measure of coverage of correct answer, demonstrating that most model-task pairs produce the correct answer most of the time with at least one completion out of 50. (d) Accuracy of greedy decoding baseline on each model-task pair.

E ADDITIONAL EXPERIMENTS AND DETAILS

In this section we detail the precise setup required to replicate our empirical results. All of our experiments were run either on 40G NVIDIA A100 GPUs, 192G AMD MI300X GPUs, or through the OpenAI API. We considered the following models. All models, except for gpt-3.5-turbo-instruct, are available on <https://huggingface.co> and we provide HuggingFace model identifiers below.

1. Phi models: We experiment with several models from the Phi family of models (Abdin et al., 2024), specifically Phi3-Mini (“microsoft/Phi-3-mini-4k-instruct”), Phi3-Small (“microsoft/Phi-3-small-8k-instruct”), Phi3-Medium (“microsoft/Phi-3-medium-4k-instruct”), and Phi3.5-Mini (“microsoft/Phi-3.5-mini-instruct”).
2. Llama3.2-3B-Instruct (“meta-llama/Llama-3.2-3B-Instruct”) (Dubey et al., 2024)
3. Mistral-7B-Instruct-v0.3 (“mistralai/Mistral-7B-Instruct-v0.3”) (Jiang et al., 2023)
4. gpt-3.5-turbo-instruct (Brown et al., 2020): We access this model via the OpenAI API.
5. llama2-7b-game24-policy-hf (“OhCherryFire/llama2-7b-game24-policy-hf”): We use the model of Wan et al. (2024), which is a Llama-2 model finetuned on the GameOf24 task (Yao et al., 2024). We use this model only the GameOf24 task.

We consider the following tasks:

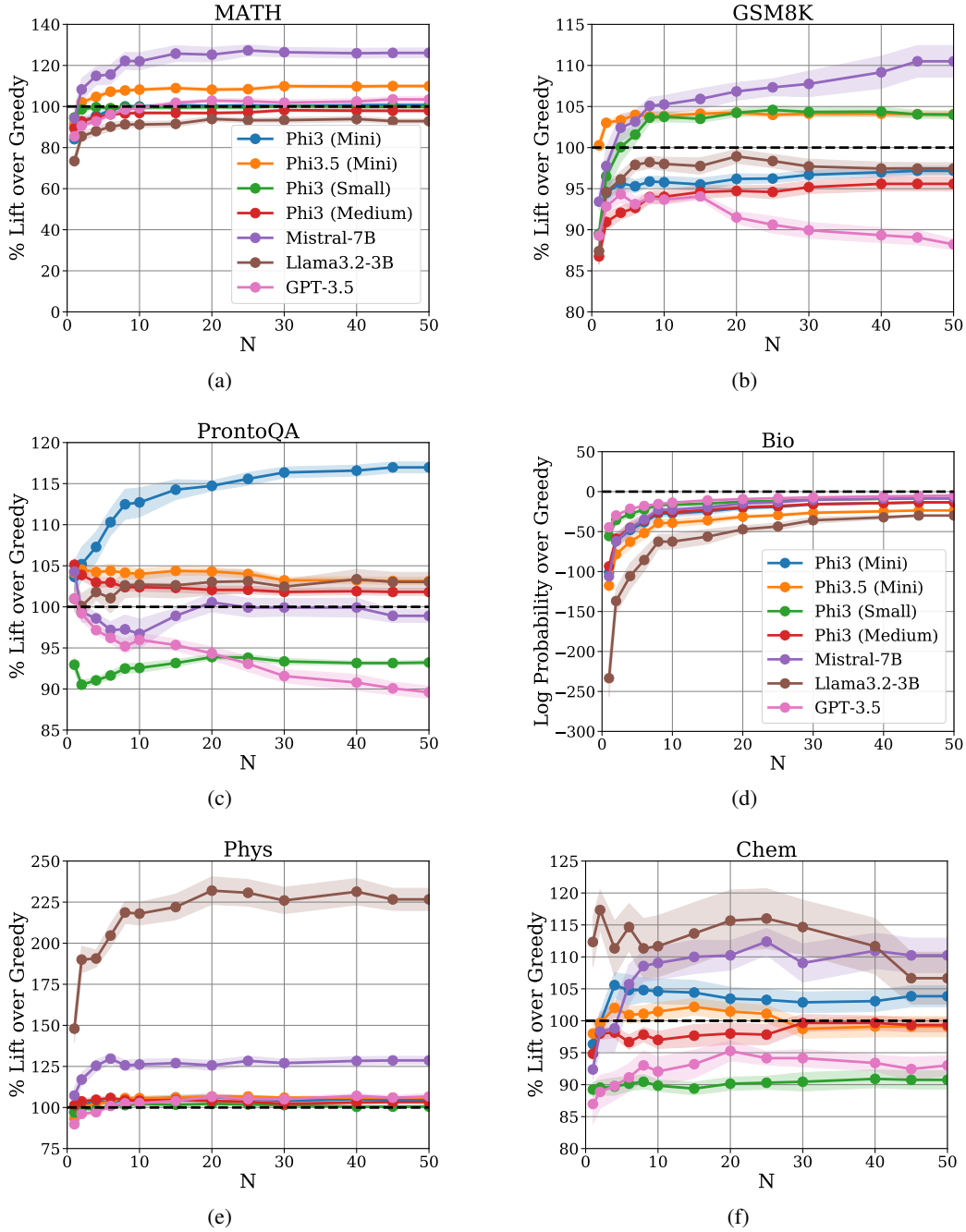


Figure 3: Percent lift of BoN-sharpening over greedy decoding accuracy as N is varied for each task. For many task-model pairs, the accuracy improves as N increases, demonstrating the effect of sequence-level log probability sharpening.

1. MATH: We use the above models to generate responses to prompts from the MATH (Hendrycks et al., 2021), which consists of more difficult math questions. We consider “all” subsets and take the first 256 examples of the test set where the solution matches the regular expression (`\d*`).⁹

⁹<https://huggingface.co/datasets/lighteval/MATH>.

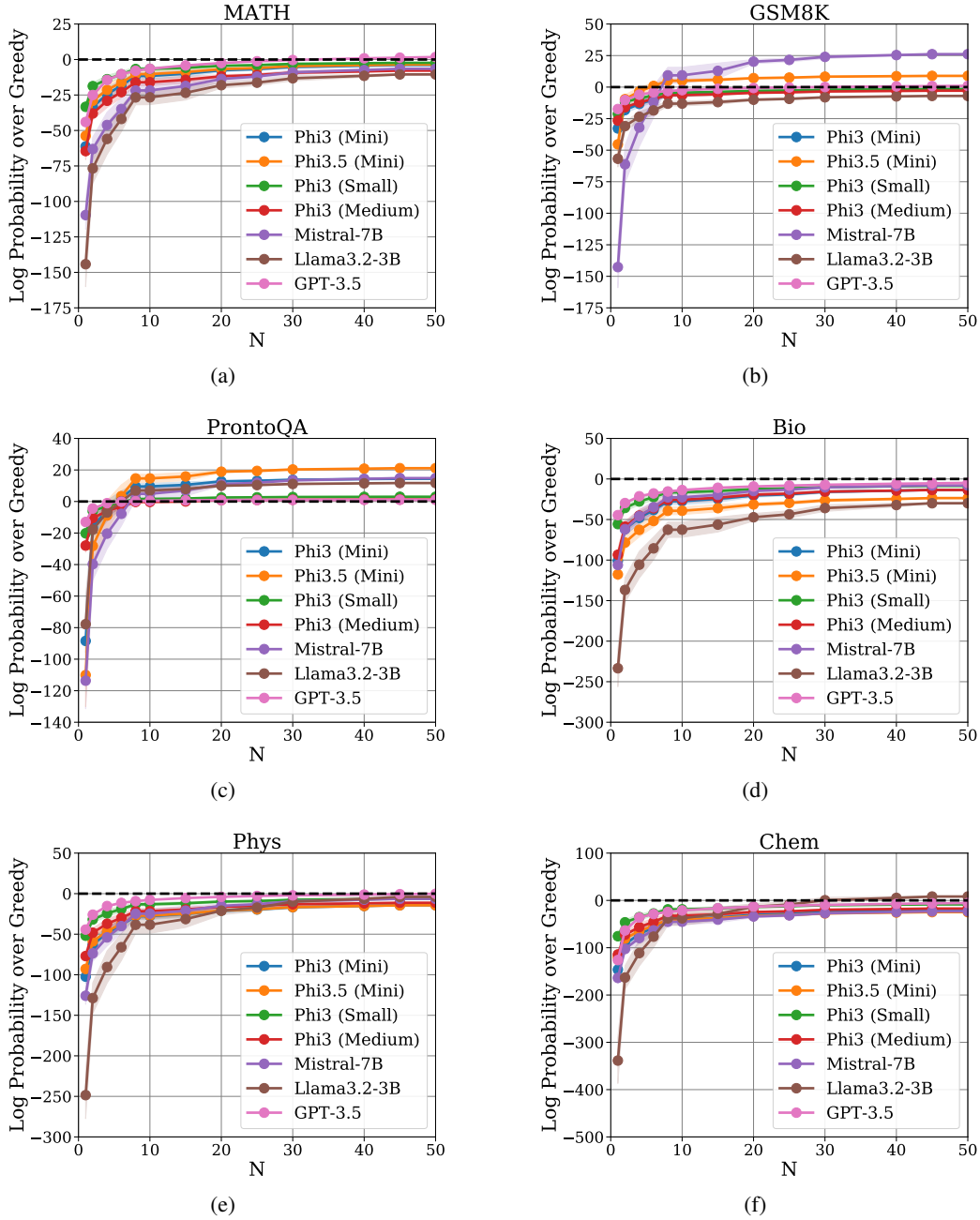


Figure 4: Effect of N on the difference in average sequence level log probabilities between inference time BoN-sharpening and greedy decoding on a variety of model-task pairs. As predicted by theory, as N increases, the likelihood of the resulting sequence increases.

2. GSM8k: We use the above models to generate responses to prompts from the GSM-8k dataset (Cobbe et al., 2021) where the goal is to generate a correct answer to an elementary school math question. We take the first 256 examples from the test set in the main subset.¹⁰

¹⁰<https://huggingface.co/datasets/openai/gsm8k>.

3. ProntoQA: We use the above models to generate responses to prompts from the ProntoQA dataset (Saparov & He, 2023), which consists of chain-of-thought-style reasoning questions with boolean answers. We take the first 256 examples from the training set.¹¹
4. MMLU: We use the above models to generate responses to prompts from three subsets of the MMLU dataset (Hendrycks et al., 2020), specifically college_biology (Bio), college_physics (Phys), and college_chemistry (Chem) all of which consist of multiple choice questions¹². We take the first 256 examples of the test set.
5. GameOf24: We use only the model of Wan et al. (2024) (i.e., llama2-7b-game24-policy-hf), on the GameOf24 task (Yao et al., 2024). The prompts are four numbers and the goal is to combine the numbers with standard arithmetic operations to reach the number ‘24.’ Here we use both the train and test splits of the dataset.¹³

E.1 INFERENCE-TIME VALIDATION EXPERIMENTS

To form the plots in Figure 1 and in Figures 3 and 4, for each (model, task) pair, we sampled N generations per prompt with temperature 1 and returned the best of the N generations according to the maximum-likelihood sharpening self-reward function $r_{\text{self}}(y | x) = \log \pi_{\text{base}}(y | x)$; we compare against greedy decoding as a baseline, whose accuracy is displayed in Figure 2(d).

Implementation details. For all models and datasets except for GameOf24, we used 1-shot prompting to ensure that models conform to the desired output format and to elicit chain of thought reasoning (for GameOf24 we do not provide a demonstration in the prompt). We set the maximum length of decoding to be 512 tokens. We used 10 seeds for all (model, task) pairs with a maximum value of $N = 50$ in Best-of- N sampling. We simulated N responses for $N < 50$ by subsampling the 50 generated samples. For Best-of- N sampling, we always use temperature 1.0. Since greedy decoding is a deterministic strategy, we only use 1 seed for each (model, task) pair. In all experiments, we collect both the responses and their log-likelihoods under the *reference model* (i.e., the original model from which samples were generated).

Results. Results for most datasets are presented in Figures 3 and 4. Because we only consider a single model for GameOf24, we separate this task into Figure 5. For all datasets, we visualize both performance—measured as normalized improvement in accuracy over greedy decoding—and log-likelihoods—under π_{base} —of the selected responses.

In all cases, Best-of- N sampling (using $r_{\text{self}}(y | x) = \log \pi_{\text{base}}(y | x)$) improves over the naïve sampling strategy, wherein we simply sample a single generation with temperature 1.0. In all datasets, we also see improvements over the standard *greedy decoding* strategy, at least for some models. Analogously, for every model, there is at least one dataset for which Best-of- N sampling improves over greedy decoding.

We further explore the relationship between sequence level log probabilities and generation quality in Figure 6, where we plot the empirical distributions of responses sampled with temperature 1 from the base model for a variety of model-dataset pairs, conditioned on whether or not the response is correct. It is clear from the figures that the distribution of log probabilities conditioned on correctness stochastically dominates that conditioned on incorrectness in each case, which provides yet more evidence that log likelihoods represent a reasonable self-improvement target.

We mention several other observations from the experiments. First, in most cases, performance and log-likelihood saturate at relatively small values of N , typically around 10 or 20. This suggests that significant improvements can be obtained with relatively low computational overhead. Second, in some cases, performance can degrade as N increases. We found that this happens for two reasons: (1) the performance of the reference model is quite low and so r_{self} provides a poor signal (e.g., with Llama3.2-3B-Instruct) and (2) the Best-of- N criteria selects for short responses, which have higher log-likelihood but cannot leverage the computational/representational benefits of chain-of-thought, and thus yield worse performance (e.g., with gpt-3.5-turbo-instruct on GSM8k).

¹¹<https://huggingface.co/datasets/longface/prontoqa-train>.

¹²<https://huggingface.co/datasets/cais/mmlu>.

¹³<https://github.com/princeton-nlp/tree-of-thought-llm/tree/master/src/tot/data/24>

Model	Dataset	% Lift over Greedy (Accuracy)	Lift over Greedy (Likelihood)
Phi3.5-Mini	MATH	19.24 ± 2.41	48.33 ± 0.17
Phi3.5-Mini	GSM8k	1.82 ± 0.64	1.49 ± 0.55
Phi3.5-Mini	ProntoQA	12.46 ± 1.08	5.64 ± 0.01
Mistral-7B	MATH	8.88 ± 5.55	5.71 ± 3.00

Table 1: Empirical Performance of SFT-Sharpening

E.2 EXPERIMENTS WITH OTHER SELF REWARD FUNCTIONS

Although we focus on $r_{\text{self}}(y | x) = \log \pi_{\text{base}}(y | x)$ throughout the paper, the sharpening framework is significantly more general. As such, we also ran experiments with other choices for r_{self} , specifically:

1. Length-normalized log-likelihood: $r_{\text{self}}(y | x) = \log \pi_{\text{base}}(y | x) / |y|$ where $|y|$ is the length, in tokens, of the response.
2. Majority (self-consistency): All datasets except GameOf24 have multiple-choice, boolean, or numerical answers. Although we allow responses to contain chain-of-thought tokens, we can extract the answer from each response and use the most-frequently-occurring answer. This can be seen as a sample-based approximation to the following self-reward function: $r_{\text{self}}(y | x) = \sum_{y': y'_{\text{ans}} = y_{\text{ans}}} \pi_{\text{base}}(y' | x)$, where y_{ans} are the “answer” tokens in the full response y .

Finally, as a skyline we consider the *coverage* criterion (Brown et al., 2024), where we simply check if any of the sampled responses corresponds to the correct answer. This criterion is a skyline and does not fit into the self-improvement framework due to the fact that it uses knowledge of the ground truth (external) task reward function.

Results are displayed in Figure 2. For length-normalized log-likelihood and majority, we see qualitatively similar behavior to (unnormalized) log-likelihood in the sense that inference-time sharpening via these self-reward functions offers improvements over both vanilla (temperature 1.0) sampling and greedy decoding. In both cases, the improvements are generally much larger than those obtained with log-likelihood. Finally, examining the coverage criteria, we see that with $N = 50$ samples, these models almost always produce a correct answer on these tasks, raising the possibility of other self-reward functions that further improve performance.

E.3 EFFECT OF SFT-Sharpening

In addition to inference-time experiments demonstrating the validity of the amortization objective considered in our theory, we also demonstrate empirically that amortization can be effected with SFT-Sharpening. Due to the realities of limited computational resources, we choose a strict subset of the model-task pairs considered in Appendix E.1 that have particularly promising inference-time BoN performance and apply SFT-Sharpening to amortize the inference time cost of multiple generations.

For each of the chosen model-dataset pairs (cf. Table 1), we sample $N = 50$ responses with temperature 1 for each prompt in the dataset and select the most likely (according to the relevant reference model). We then combine these likely responses with the prompts in order to form a training corpus and train a Low Rank Adaptation (Hu et al., 2021) to the model, sweeping over LoRA rank, learning rate scheduler, and weight decay in order to return the best optimized model.¹⁴ We report the specific hyperparameters chosen in Table 2. On all models, we used a learning rate of 3×10^{-4} with linear decay to zero and gradient clamping at 0.1.

Results. In Table 1 we report our results for the best model during training of each model-dataset pair, averaged across 3 random seeds, where responses are sampled with temperature 1 from the fine-tuned model. We report both the percent lift in accuracy on the dataset with respect to the greedy generation of the reference model and the increase in average sequence level log likelihood

¹⁴In all experiments involving Phi3.5-Mini we use a batch size of 4; unfortunately, due to a known numerical issue with LoRA on Mistral-7B-Instruct-v0.3 involving batch size > 1 , we use a batch of 1 in this case. Because of this choice, instead of the 30 epochs we use to train our other models, for Mistral-7B-Instruct-v0.3, we run only 10 epochs.

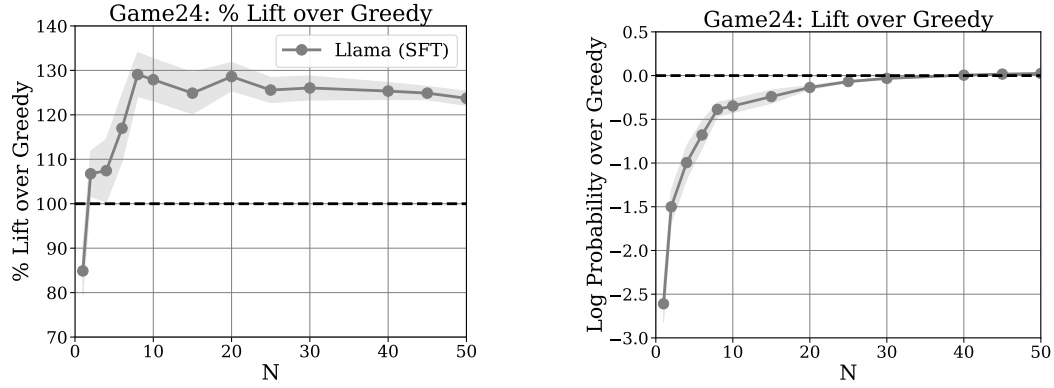


Figure 5: Effect of inference-time BoN-sharpening on GameOf24 with the finetuned llama2-7b-game24-policy-hf from Wan et al. (2024).

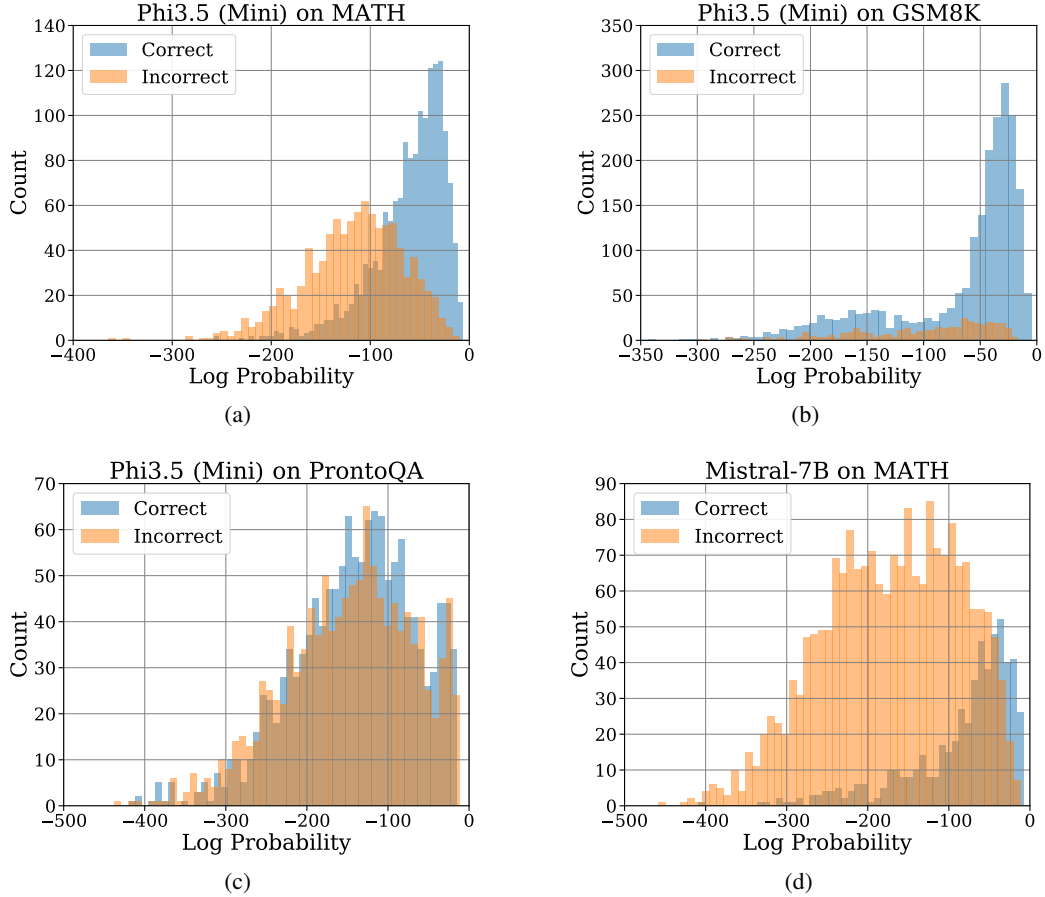


Figure 6: Distribution of sequence-level log probabilities of sampled responses with temperature 1, conditioned on whether or not the response is correct for 4 model-dataset pairs: (a) (Phi3.5-Mini, MATH); (b) (Phi3.5-Mini, GSM8k); (c) (Phi3.5-Mini, ProntoQA); (d) (Mistral-7B-Instruct-v0.3, MATH). In all cases, conditioning on the response being correct leads to a noticeable increase in log probabilities, further justifying the use of sequence-level log probabilities as a valid self-improvement score.

Model	Dataset	Weight Decay	LoRA Rank
Phi3.5-Mini	MATH	0.1	16
Phi3.5-Mini	GSM8k	0.5	16
Phi3.5-Mini	ProntoQA	0.0	16
Mistral-7B-Instruct-v0.3	MATH	1.0	8

Table 2: Empirical Performance of SFT-Sharpener

with respect to the same. In all cases, we see improvement on both metrics, demonstrating that some amortization is possible with SFT-Sharpener. In Figures 7 and 8, we display the evolution throughout training of these same metrics for each of the model-dataset pairs. While Phi3.5-Mini is quite well-behaved on MATH and ProntoQA, there appears to be a fair amount of noise in the training on GSM8k, with the log probability being a significantly less useful proxy for accuracy on this dataset than the others. In the case of Mistral-7B-Instruct-v0.3 on MATH, while we do see some improvement after sufficient training, the optimization suffers an initial substantial drop and then spends $\sim 90\%$ of the gradient steps recovering; we speculate that this is a function of insufficient hyper-parameter tuning of the optimization itself, rather than a fundamental barrier.

Finally, in Figure 9, we investigate the effect that the choice of N has on SFT-Sharpener for Phi3.5-Mini on MATH. In particular, in forming our training set, we choose $N \in \{10, 25, 50\}$ and repeat the procedure described above, averaging our results over three seeds. We find that increasing N leads to a modest increase in the sequence-level log-likelihood and a consequent increment in the accuracy of the fine-tuned model, in accordance with our theory.

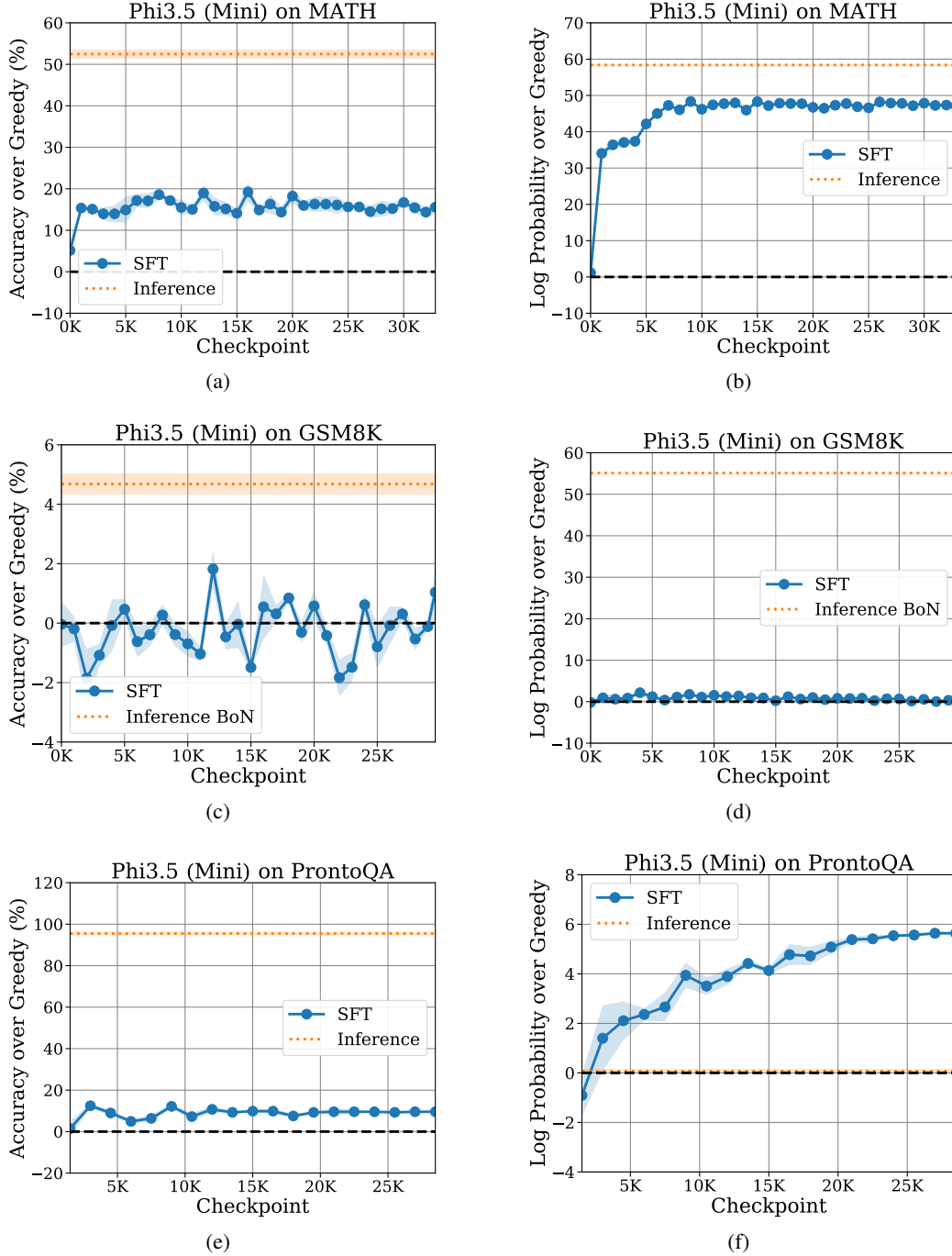


Figure 7: Evolution of Phi3.5-Mini under SFT-Sharpener on different datasets for $N = 50$ as measured by % lift over Greedy in accuracy and difference in average sequence-level log probability under the reference model of generated responses. The fine-tuned model produces generations with high probability under the reference model and a consequent increase in accuracy; the model still is not able to match the performance of the inference-time BoN approach.

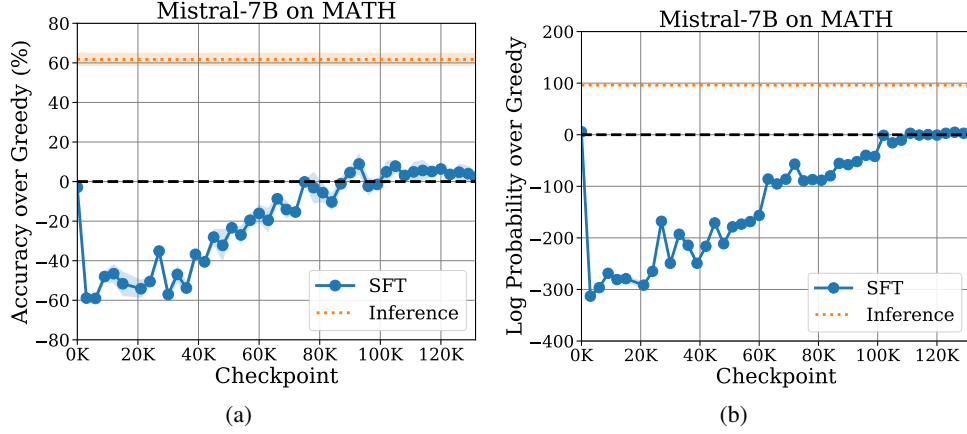


Figure 8: Evolution of Mistral-7B-Instruct-v0.3 under SFT-Sharpener on MATH for $N = 50$ as measured both by % lift over Greedy in accuracy and difference in average sequence-level log probability under the reference model of generated responses.

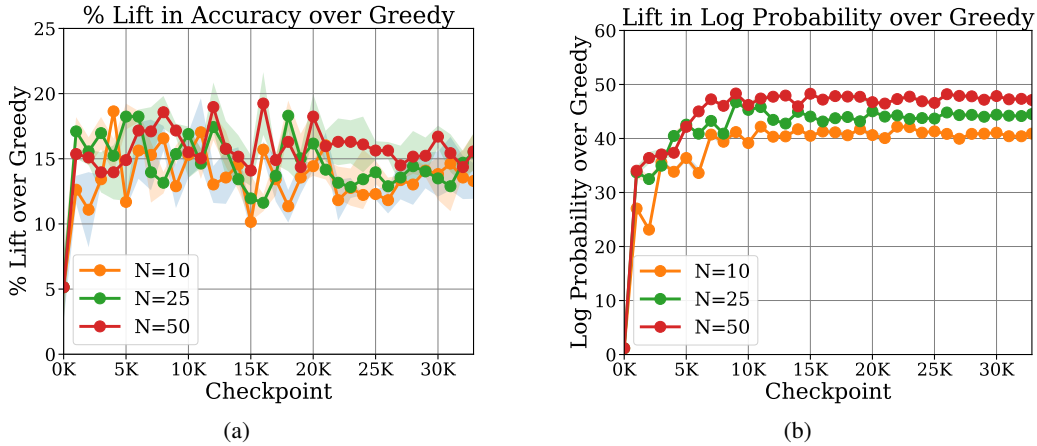


Figure 9: Effect that choice of N has on SFT-Sharpener for Phi3.5-Mini on MATH. We report both (a) % lift over greedy in accuracy and (b) lift in sequence-level log likelihood averaged over the dataset. In both cases, we see that increasing N leads to more lift, in accordance with theory.

Part II

Proofs

F PRELIMINARIES

F.1 GUARANTEES FOR APPROXIMATE MAXIMIZERS

Recall that the theoretical guarantees for sharpening algorithms in [Section 4](#) provide convergence to the set $\mathbf{y}^*(x) := \arg \max_{y \in \mathcal{Y}} \pi_{\text{base}}(y \mid x)$ of (potentially non-unique) maximizers for the maximum-likelihood sharpening self-reward function $\log \pi_{\text{base}}(y \mid x)$. These guarantees require that the base model π_{base} places sufficient provability mass on $\mathbf{y}^*(x)$, which may be unrealistic. To address this, throughout this appendix we state and prove more general versions of our theoretical results that allow for approximate maximizers, and consequently enjoy weaker coverage assumptions

For a parameter $\gamma \in [0, 1)$ we define

$$\mathbf{y}_\gamma^*(x) := \left\{ y \mid \pi_{\text{base}}(y \mid x) \geq (1 - \gamma) \cdot \max_{y \in \mathcal{Y}} \pi_{\text{base}}(y \mid x) \right\}$$

as the set of $(1 - \gamma)$ -approximate maximizers for $\log \pi_{\text{base}}(y \mid x)$. We quantify the quality of a sharpened model as follows.

Definition F.1 (Sharpened model). *We say that a model $\hat{\pi}$ is $(\epsilon, \delta, \gamma)$ -sharpened relative to π_{base} if*

$$\mathbb{P}_{x \sim \mu} [\hat{\pi}(\mathbf{y}_\gamma^*(x) \mid x) \geq 1 - \delta] \geq 1 - \epsilon.$$

That is, an $(\epsilon, \delta, \gamma)$ -sharpened policy places at least $1 - \delta$ mass on $(1 - \gamma)$ -approximate arg-max responses on all but an ϵ -fraction of prompts under μ .

Lastly, we will make use of the following generalized coverage coefficient

$$C_{\text{cov}, \gamma} = \mathbb{E}_{x \sim \mu} \left[\frac{1}{\pi_{\text{base}}(\mathbf{y}_\gamma^*(x) \mid x)} \right],$$

which has $C_{\text{cov}, \gamma} \leq C_{\text{cov}}$.

F.2 TECHNICAL TOOLS

For a pair of probability measures \mathbb{P} and \mathbb{Q} with a common dominating measure ω , Hellinger distance is defined via

$$D_{\text{H}}^2(\mathbb{P}, \mathbb{Q}) = \int \left(\sqrt{\frac{d\mathbb{P}}{d\omega}} - \sqrt{\frac{d\mathbb{Q}}{d\omega}} \right)^2 d\omega.$$

Lemma F.1 (MLE for conditional density estimation (e.g., [Wong & Shen \(1995\)](#); [van de Geer \(2000\)](#); [Zhang \(2006\)](#))). *Consider a conditional density $\pi^* : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$. Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ be a dataset in which (x_i, y_i) are drawn i.i.d. as $x_i \sim \mu \in \Delta(\mathcal{X})$ and $y_i \sim \pi^*(\cdot \mid x)$. Suppose we have a finite function class $\Pi \subset (\mathcal{X} \rightarrow \Delta(\mathcal{Y}))$ such that $\pi^* \in \Pi$. Define the maximum likelihood estimator*

$$\hat{\pi} := \arg \max_{\pi \in \Pi} \sum_{(x, y) \in \mathcal{D}} \log \pi(y \mid x).$$

Then with probability at least $1 - \rho$,

$$\mathbb{E}_{x \sim \mu} [D_{\text{H}}^2(\hat{\pi}(\cdot \mid x), \pi^*(\cdot \mid x))] \leq \frac{2 \log(|\Pi| \rho^{-1})}{n}.$$

Lemma F.2 (Elliptic potential lemma). *Let $\lambda, K > 0$, and let $A_1, \dots, A_T \in \mathbb{R}^{d \times d}$ be positive semi-definite matrices with $\text{Tr}(A_t) \leq K$ for all $t \in [T]$. Fix $\Gamma_0 = \lambda I_d$ and $\Gamma_t = \lambda I_d + \sum_{i=1}^t A_i$ for $t \in [T]$. Then*

$$\sum_{t=1}^T \text{Tr}(\Gamma_{t-1}^{-1} A_t) \leq \frac{dK \log \frac{(T+1)K}{\lambda}}{\lambda \log(1 + K/\lambda)}.$$

Proof of Lemma F.2. Fix $t \in [T]$. Since $\text{Tr}(A_t) \leq 1$, there is some $p_t \in \Delta(\mathbb{R}^d)$ such that $A_t = \mathbb{E}_{a \sim p_t} aa^\top$ and $\mathbb{P}[\|a\|_2 \leq 1] = 1$. Now observe that

$$\begin{aligned} \log \det(\Gamma_t) &= \log \det(\Gamma_{t-1} + A_t) \\ &= \log \det(\Gamma_{t-1}) + \log \det(I_d + \Gamma_{t-1}^{-1/2} A_t \Gamma_{t-1}^{-1/2}) \\ &= \log \det(\Gamma_{t-1}) + \log \det \left(\mathbb{E}_{a \sim p_t} \left[I_d + \Gamma_{t-1}^{-1/2} aa^\top \Gamma_{t-1}^{-1/2} \right] \right) \\ &\geq \log \det(\Gamma_{t-1}) + \mathbb{E}_{a \sim p_t} \log \det(I_d + \Gamma_{t-1}^{-1/2} aa^\top \Gamma_{t-1}^{-1/2}) \\ &= \log \det(\Gamma_{t-1}) + \mathbb{E}_{a \sim p_t} \log(1 + a^\top \Gamma_{t-1}^{-1} a). \end{aligned}$$

Now $a^\top \Gamma_{t-1}^{-1} a \leq 1/\lambda$ with probability 1, where $\lambda = \lambda_{\min}(\Gamma_0)$. We know that $\lambda x \log(1 + 1/\lambda) \leq \log(1 + x)$ for all $x \in [0, 1/\lambda]$. Thus,

$$\log \det(\Gamma_t) \geq \log \det(\Gamma_{t-1}) + \lambda \log(1 + 1/\lambda) \mathbb{E}_{a \sim p_t} a^\top \Gamma_{t-1}^{-1} a.$$

Summing over $t \in [T]$, we get

$$\log \det(\Gamma_T) \geq \log \det(\Gamma_0) + \lambda \log(1 + 1/\lambda) \sum_{t=1}^T \text{Tr}(\Gamma_{t-1}^{-1} A_t).$$

Finally note that $\lambda_{\max}(\Gamma_T) \leq T + 1$ so $\log \det(\Gamma_T) \leq d \log T$, whereas $\log \det(\Gamma_0) \geq d \log \lambda$. Thus,

$$\sum_{t=1}^T \text{Tr}(\Gamma_{t-1}^{-1} A_t) \leq \frac{d \log \frac{T+1}{\lambda}}{\lambda \log(1 + 1/\lambda)}$$

as claimed. \square

Lemma F.3 (Freedman's inequality, e.g. Agarwal et al. (2014)). Let $(Z_t)_{t=1}^T$ be a martingale difference sequence adapted to filtration $(\mathcal{F}_t)_{t=0}^{T-1}$. Suppose that $|Z_t| \leq R$ holds almost surely for all t . For any $\delta \in (0, 1)$ and $\eta \in (0, 1/R)$, it holds with probability at least $1 - \delta$ that

$$\sum_{t=1}^T Z_t \leq \eta \sum_{t=1}^T \mathbb{E}[Z_t^2 | \mathcal{F}_{t-1}] + \frac{\log(1/\delta)}{\eta}.$$

Corollary F.1. Let $(Z_t)_{t=1}^T$ be a sequence of random variables adapted to filtration $(\mathcal{F}_t)_{t=0}^{T-1}$. Suppose that $Z_t \in [0, R]$ holds almost surely for all t . For any $\delta \in (0, 1)$, it holds with probability at least $1 - \delta$ that

$$\sum_{t=1}^T \mathbb{E}[Z_t | \mathcal{F}_{t-1}] \leq 2 \sum_{t=1}^T Z_t + 4R \log(1/\delta).$$

Proof of Corollary F.1. Observe that for any $t \in [T]$,

$$\begin{aligned} \mathbb{E}[(Z_t - \mathbb{E}[Z_t | \mathcal{F}_{t-1}])^2 | \mathcal{F}_{t-1}] &\leq \mathbb{E}[Z_t^2 | \mathcal{F}_{t-1}] \\ &\leq R \cdot \mathbb{E}[Z_t | \mathcal{F}_{t-1}]. \end{aligned}$$

Applying Lemma F.3 to the sequence $(\mathbb{E}[Z_t | \mathcal{F}_{t-1}] - Z_t)_{t=1}^T$, which is a martingale difference sequence with elements supported almost surely on $[-R, R]$, we get for any $\eta \in (0, 1/R)$ that with probability at least $1 - \delta$,

$$\begin{aligned} \sum_{t=1}^T (\mathbb{E}[Z_t | \mathcal{F}_{t-1}] - Z_t) &\leq \eta \sum_{t=1}^T \mathbb{E}[(Z_t - \mathbb{E}[Z_t | \mathcal{F}_{t-1}])^2 | \mathcal{F}_{t-1}] + \frac{\log(1/\delta)}{\eta} \\ &\leq \eta R \sum_{t=1}^T \mathbb{E}[Z_t | \mathcal{F}_{t-1}] + \frac{\log(1/\delta)}{\eta}. \end{aligned}$$

Set $\eta = 1/(2R)$. Simplifying gives

$$\sum_{t=1}^T \mathbb{E}[Z_t | \mathcal{F}_{t-1}] \leq 2 \sum_{t=1}^T Z_t + 4R \log(1/\delta).$$

as claimed. \square

G PROOFS FROM SECTION 3.1

Proof of Proposition 3.1. We prove the result by induction. Fix $x \in \mathcal{X}$, and let $y_1^*, \dots, y_H^* := y^*(x)$. Fix $h \in [H]$, and assume by induction that $\hat{y}_{h'} = y_{h'}^*$ for all $h' < h$. We claim that in this case,

$$\pi_h(y_h^* \mid \hat{y}_1, \dots, \hat{y}_{h-1}, x) = \pi_h(y_h^* \mid y_1^*, \dots, y_{h-1}^*, x) > 1/2,$$

which implies that $\hat{y}_h = y_h^*$. To see this, we observe that by Bayes' rule,

$$\begin{aligned} \pi(y_1^*, \dots, y_H^* \mid x) &\leq \pi(y_1^*, \dots, y_h^* \mid x) \\ &= \prod_{h'=1}^h \pi_{h'}(y_{h'}^* \mid y_1^*, \dots, y_{h'-1}^*, x) \leq \pi_h(y_h^* \mid y_1^*, \dots, y_{h-1}^*, x). \end{aligned}$$

If we were to have $\pi_h(y_h^* \mid \hat{y}_1, \dots, \hat{y}_{h-1}, x) = \pi_h(y_h^* \mid y_1^*, \dots, y_{h-1}^*, x) \leq 1/2$, it would contradict the assumption that $\pi(y_1^*, \dots, y_H^* \mid x) > 1/2$. This proves the result. \square

H PROOFS FROM SECTION 3.3

Below, we state and prove a generalization of Theorems 3.1 and C.2 which allows for approximate maximizers in the sense of Definition F.1, as well as a more general coverage coefficient.

To state the result, for a model π , we define

$$\mathbf{y}_\gamma^\pi(x) = \left\{ y \mid \pi(y \mid x) \geq (1 - \gamma) \cdot \max_{y \in \mathcal{Y}} \pi(y \mid x) \right\}.$$

Next, for any integer $p \in \mathbb{N}$, we define

$$C_{\text{cov}, \gamma, p}(\pi) = \left(\mathbb{E} \left[\frac{1}{(\pi(\mathbf{y}_\gamma^\pi(x) \mid x))^p} \right] \right)^{1/p},$$

with the convention that $C_{\text{cov}, \gamma, p} = C_{\text{cov}, \gamma, p}(\pi_{\text{base}})$. For our negative results, we select $\gamma = 1/2$. Thus, our lower bounds which we are about to state and prove hold *in a regime where the best y has bounded margin away from suboptimal responses*.

Theorem 3.1' (Lower bound for sharpening). *Fix integers $d \geq 1$ and $p \geq 1$ and parameters $\epsilon \in (0, 1)$ and $C \geq 1$, and set $\gamma = 1/2$. There exists a class of models Π such that i) $\log |\Pi| \asymp d(1 + \log(C\epsilon^{-1/p}))$, ii) $\sup_{\pi \in \Pi} C_{\text{cov}, \gamma, p}(\pi) \lesssim C$, and iii) $\mathbf{y}_\gamma^\pi(x)$ is a singleton for all $\pi \in \Pi$, for which any sharpening algorithm $\hat{\pi}$ that attains $\mathbb{E}[\mathbb{P}_{x \sim \mu}[\hat{\pi}(\mathbf{y}_\gamma^{\pi_{\text{base}}}(x)) > 1/2]] \geq 1 - \epsilon$ for all $\pi_{\text{base}} \in \Pi$ must collect a total number of samples $m = n \cdot N$ at least*

$$m \gtrsim \begin{cases} \frac{C \log |\Pi|}{\epsilon^{1+1/p}(1+\log(C\epsilon^{-1/p}))} & \text{sample-and-evaluate oracle,} \\ \frac{C \log |\Pi|}{\epsilon^{1/p}(1+\log(C\epsilon^{-1/p}))} & \text{adaptive sample-and-evaluate oracle.} \end{cases}$$

Proof of Theorem 3.1'. Let parameter $d, p \in \mathbb{N}$ and $\epsilon > 0$ be given, and set $\gamma = 1/2$. Let $M \in \mathbb{N}$ and $\Delta > 0$ be parameter to be chosen later. Let $\mathcal{X} = \{x_0, x_1, \dots, x_d\}$ and $\mathcal{Y} = \{y_0, y_1, \dots, y_M\}$ be arbitrary discrete sets (with $|\mathcal{X}| = d + 1$ and $|\mathcal{Y}| = M + 1$).

Construction of prompt distribution and model class. We use the same construction for the non-adaptive and adaptive lower bounds in the theorem statement. We define the prompt distribution μ via

$$\mu := (1 - \Delta)\delta_{x_0} + \frac{\Delta}{d} \sum_{i=1}^d \delta_{x_i},$$

where δ_x denotes the Dirac delta distribution on element x .

As the first step toward constructing the model class Π , we introduce a family of distributions (P_0, P_1, \dots, P_M) on \mathcal{Y} as follows

$$P_0 = \delta_{y_0}, \quad \forall i \geq 1, \quad P_i = \frac{1}{(1 - \gamma)M} \delta_{y_i} + \sum_{j \in [M] \setminus \{i\}} \frac{1}{M} \left(1 - \frac{\gamma}{(M - 1)(1 - \gamma)} \right) \delta_{y_j}.$$

Next, for or any index $\mathcal{I} = (j_1, j_2, \dots, j_d) \in [M]^d$, define a model

$$\pi^{\mathcal{I}}(x_i) = \begin{cases} P_0 & i = 0 \\ P_{j_i} & i > 0 \end{cases}.$$

We define the model class as

$$\Pi := \{\pi^{\mathcal{I}} : \mathcal{I} \in [M]^d\},$$

which we note has

$$\log |\Pi| = d \log M.$$

Preliminary technical results. Define

$$\mathbf{y}_{\gamma}^{\mathcal{I}}(x) := \{y : \pi^{\mathcal{I}}(y | x) \geq (1 - \gamma) \max_{y \in \mathcal{Y}} \pi^{\mathcal{I}}(y | x)\}.$$

The following property is immediate.

Lemma H.1. *Let $\mathcal{I} = (j_1, \dots, j_d) \in [d]^M$. Then $\mathbf{y}_{\gamma}^{\mathcal{I}}(x_i) = \{y_{j_i}\}$ if $i > 0$, and $\mathbf{y}_{\gamma}^{\mathcal{I}}(x_0) = \{y_0\}$.*

In view of this result, we define $y^{\mathcal{I}}(x) = \arg \max_y \pi^{\mathcal{I}}(y | x)$ as the unique arg-max response for x .

Going forward, let us fix the algorithm under consideration. Let $\mathbb{P}^{\mathcal{I}}[\cdot]$ denote the law over the dataset used by the algorithm when the true instance is $\pi^{\mathcal{I}}$ (including possible randomness and adaptivity from the algorithm itself), and let $\mathbb{E}^{\mathcal{I}}[\cdot]$ denote the corresponding expectation. The following lemma is a basic technical result.

Lemma H.2 (Reduction to classification). *Let $\hat{\pi}$ be the model produced by an algorithm with access to a sample-and-evaluate oracle for $\pi^{\mathcal{I}}$. Suppose that for some $\epsilon \geq 0$,*

$$\mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} \mathbb{P}_{x \sim \mu} [\hat{\pi}(\mathbf{y}_{\gamma}^{\mathcal{I}}(x) | x) > 1/2] \geq 1 - \epsilon.$$

Define $\hat{\mathcal{I}} = (\hat{j}_1, \dots, \hat{j}_d)$ via $\hat{j}_i = \arg \max_j \hat{\pi}(y_j | x_i)$, and write $\mathcal{I} = (j_1^, \dots, j_d^*)$. Then,*

$$\frac{1}{d} \sum_{i=1}^d \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{\hat{j}_i \neq j_i^*\}] \leq \epsilon / \Delta.$$

Proof of Lemma H.2. As established in Lemma H.1, under instance \mathcal{I} , $\mathbf{y}_{\gamma}^{\mathcal{I}}(x_i) = \{y_{j_i^*}\}$ for any $i \in [d]$. Thus, whenever $\hat{\pi}(\mathbf{y}_{\gamma}^{\mathcal{I}}(x_i)) > 1/2$, $j_i^* = \arg \max_j \hat{\pi}(y_j | x_i) =: \hat{j}_i$. The result follows by noting that the event $\{\exists i \in [d] : x = x_i\}$ occurs with probability at least Δ under $x \sim \mu$. \square

Lower bound under sample-and-evaluate oracle. Recall that in the non-adaptive framework, the sample complexity m is fixed. In light of Lemma H.2, it suffices to establish the following claim.

Lemma H.3. *There exists a universal constant $c > 0$ such that for all $M \geq 8$, if $m \leq cdM/\Delta$, then $\mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{\hat{j}_i \neq j_i^*\}] \geq 1/8$ for all i .*

With this, the result follows by selecting $\Delta = 16\epsilon$, with which Lemma H.2 implies that any algorithm with $\mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} \mathbb{P}_{x \sim \mu} [\hat{\pi}(\mathbf{y}_{\gamma}^{\mathcal{I}}(x) | x) > 1/2] \geq 1 - \epsilon$ must have $m \gtrsim dM/\Delta$, then. To conclude, we choose $M \approx 1 + C\epsilon^{-1/p}$, which gives $m \approx dM/\Delta \approx dC\epsilon^{-(1+1/p)} \approx \epsilon^{-(1+1/p)} \log \Pi / \log(1 + C\epsilon^{1/p})$. Finally, we check that with this choice, all $\pi \in \Pi$ satisfy

$$\begin{aligned} C_{\text{cov}, \gamma, p}(\pi) &= (\mathbb{P}_{x \sim \mu}[x = x_0] + (M(1 - \gamma))^p \mathbb{P}_{x \sim \mu}[x \neq x_0])^{1/p} \\ &= ((1 - \Delta) + (M(1 - \gamma))^p \Delta)^{1/p} \\ &\lesssim ((1 - \Delta) + (8C(1 - \gamma))^p)^{1/p} \lesssim C. \end{aligned}$$

Proof of Lemma H.3. Let $i \in [d]$ be fixed. Of the $m = n \cdot N$ tuples $(x, y, \log \pi_{\text{base}}(y | x))$ that are observed by the algorithm, let m_i denote (random) the number of such examples for which $x = x_i$. From Markov's inequality, we have

$$\mathbb{P}[m_i \leq 2\Delta m/d] \geq \frac{1}{2} \tag{13}$$

Going forward, let $\mathcal{D} = \{(x, y, \log \pi_{\text{base}}(y | x))\}$ denote the dataset collected by the algorithm, which has $|\mathcal{D}| = m$. Let \mathcal{E}_i denote the event that, for prompt $x = x_i$, (i) there are at least two distinct responses y_j for which $(x_i, y_j) \notin \mathcal{D}$; and (ii) there are no pairs $(x_i, y) \in \mathcal{D}$ for which $\pi_{\text{base}}(y | x_i) > \frac{1}{M}$. Since \mathcal{E}_i is a measurable function of \mathcal{D} , we can write

$$\begin{aligned} \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{\hat{j}_i \neq j_i^*\}] &\geq \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{\hat{j}_i \neq j_i^*\} \cdot \mathbb{I}\{\mathcal{E}_i\}] \\ &= \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{\mathcal{E}_i\} \mathbb{E}_{\mathcal{I} \sim \mathbb{P}[\mathcal{I} = \cdot | \mathcal{D}]} [\mathbb{I}\{\hat{j}_i \neq j_i^*\}]], \end{aligned} \quad (14)$$

where $\mathcal{I} \sim \mathbb{P}[\mathcal{I} = \cdot | \mathcal{D}]$ is sampled from the posterior distribution over \mathcal{I} conditioned on the dataset \mathcal{D} . Observe that conditioned on \mathcal{E}_i , the posterior distribution over j_i^* under $\mathcal{I} \sim \mathbb{P}[\mathcal{I} = \cdot | \mathcal{D}]$ is uniform over the set of indices $j \in [M]$ for which $(x_i, y_j) \notin \mathcal{D}$, and this set has size at least 2. Hence, $\mathbb{I}\{\mathcal{E}_i\} \mathbb{E}_{\mathcal{I} \sim \mathbb{P}[\mathcal{I} = \cdot | \mathcal{D}]} [\mathbb{I}\{\hat{j}_i \neq j_i^*\}] \geq \frac{1}{2}$, and resuming from Eq. (16), we have

$$\begin{aligned} \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{\hat{j}_i \neq j_i^*\}] &\geq \frac{1}{2} \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{\mathcal{E}_i\}] \geq \frac{1}{2} \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{P}^{\mathcal{I}} [\mathcal{E}_i \cap \{m_i \leq 2\Delta m/d\}] \\ &\geq \frac{1}{4} \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{P}^{\mathcal{I}} [\mathcal{E}_i | m_i \leq 2\Delta m/d], \end{aligned}$$

where the last inequality is from Eq. (13). Finally, we can check that, under the law $\mathbb{P}^{\mathcal{I}}$, the probability of the event \mathcal{E}_i —conditioned on the value m_i —is at least the probability that $(x_i, y_{j_i^*}), (x_i, y_{j'}) \notin \mathcal{D}$ for an arbitrary fixed index $j' \neq j_i^*$, which on the event $\{m_i \leq 2\Delta m/d\}$ is at least

$$\left(1 - \frac{3}{M}\right)^{m_i} \geq \left(1 - \frac{3}{M}\right)^{2\Delta m/d},$$

where we have used that $\gamma = 1/2$. The value above is at least $\frac{1}{4}$ whenever $m \leq c \cdot dM/\Delta$ for a sufficiently small absolute constant $c > 0$. For this value of m , we conclude that $\mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{\hat{j}_i \neq j_i^*\}] \geq \frac{1}{4} \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{P}^{\mathcal{I}} [\mathcal{E}_i | \{m_i \leq 2\Delta m/d\}] \geq \frac{1}{8}$. \square

Lower bound under adaptive sample-and-evaluate oracle. In the adaptive framework, we let m_i denote the (potentially random) number of tuples $(x, y, \log \pi_{\text{base}}(y | x))$ observed by the algorithm in which $x = x_i$. Note that unlike the non-adaptive framework, the distribution over m_i depends on the underlying instance \mathcal{I} with which the algorithm interacts.

To begin, from Lemma H.2 and Markov’s inequality, if $\hat{\pi}$ satisfies the guarantee $\mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{P}_{x \sim \mu} [\hat{\pi}(\mathbf{y}_{\gamma}^{\mathcal{I}}(x)) > 1/2] \geq 1 - \epsilon$, then there exists a set of indices $S_{\text{good}} \subset [d]$ such that¹⁵

$$|S_{\text{good}}| \geq \lfloor d/2 \rfloor, \quad \forall i \in S_{\text{good}}, \quad \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{\hat{j}_i \neq j_i^*\}] \leq \frac{2\epsilon}{\Delta}. \quad (15)$$

We now appeal to the following lemma.

Lemma H.4. *As long as $M \geq 6$, it holds that for all $i \in [d]$,*

$$\mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{\hat{j}_i \neq j_i^*\}] \geq \frac{1}{4e} \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{m_i \leq M/3\}].$$

Combining Lemma H.4 with Eq. (15), it follows that there exist absolute constant $c_1, c_2, c_3 > 0$ such that if $\Delta = c_1 \cdot \epsilon$, then for all $i \in S_{\text{good}}$,

$$\mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{P}^{\mathcal{I}} [m_i \geq c_2 M] \geq c_3.$$

Thus, with this choice for Δ , we have that $i \in S_{\text{good}}$,

$$\mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [m_i] \gtrsim M,$$

¹⁵We emphasize that the set S_{good} is not a random variable, and depends only on the algorithm itself.

and we can lower bound the algorithm's expected sample complexity by summing over $i \in S_{\text{good}}$:

$$\mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}}[m] \geq \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} \left[\sum_{i \in S_{\text{good}}} m_i \right] \gtrsim |S_{\text{good}}| M \gtrsim dM.$$

The result now follows by tuning $M \approx 1 + C\epsilon^{-1/p}$ as in the proof of the lower bound for non-adaptive sampling, which gives $\mathbb{E}[m] \gtrsim dM \approx dC\epsilon^{-1/p} \approx \epsilon^{-1/p} \log \Pi / \log(1 + C\epsilon^{1/p})$ and $C_{\text{cov}, \gamma, p}(\pi) \lesssim C$ for all $\pi \in \Pi$.

Proof of Lemma H.4. Let $i \in [d]$ be fixed. Let $\mathcal{D} = \{(x, y, \log \pi_{\text{base}}(y | x))\}$ denote the dataset collected by the algorithm at termination, which has $|\mathcal{D}| = m$. Let \mathcal{E}_i denote the event that, for prompt $x = x_i$, (i) there are at least two distinct responses y_j for which $(x_i, y_j) \notin \mathcal{D}$; and (ii) there are no pairs $(x_i, y) \in \mathcal{D}$ for which $\pi_{\text{base}}(y | x_i) > \frac{1}{M}$. Since \mathcal{E}_i is a measurable function of \mathcal{D} , we can write

$$\begin{aligned} \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{\hat{j}_i \neq j_i^*\}] &\geq \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{\hat{j}_i \neq j_i^*\} \cdot \mathbb{I}\{\mathcal{E}_i\}] \\ &= \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{\mathcal{E}_i\} \mathbb{E}_{\mathcal{I} \sim \mathbb{P}[\mathcal{I} = \cdot | \mathcal{D}]} [\mathbb{I}\{\hat{j}_i \neq j_i^*\}]], \end{aligned} \quad (16)$$

where $\mathcal{I} \sim \mathbb{P}[\mathcal{I} = \cdot | \mathcal{D}]$ is sampled from the posterior distribution over \mathcal{I} conditioned on the dataset \mathcal{D} . Observe that conditioned on \mathcal{E}_i , the posterior distribution over j_i^* under $\mathcal{I} \sim \mathbb{P}[\mathcal{I} = \cdot | \mathcal{D}]$ is uniform over the set of indices $j \in [M]$ for which $(x_i, y_j) \notin \mathcal{D}$, and this set has size at least 2. Hence, $\mathbb{I}\{\mathcal{E}_i\} \mathbb{E}_{\mathcal{I} \sim \mathbb{P}[\mathcal{I} = \cdot | \mathcal{D}]} [\mathbb{I}\{\hat{j}_i \neq j_i^*\}] \geq \frac{1}{2}$, and resumming from Eq. (16), we have

$$\begin{aligned} \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{\hat{j}_i \neq j_i^*\}] &\geq \frac{1}{2} \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{\mathcal{E}_i\}] \\ &\geq \frac{1}{2} \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{P}^{\mathcal{I}} [\mathcal{E}_i \cap \{m_i \leq M/3\}] \\ &= \frac{1}{2} \mathbb{E}_{\mathcal{I} \sim \text{Unif}} [\mathbb{P}^{\mathcal{I}} [\mathcal{E}_i | m_i \leq M/3] \cdot \mathbb{P}^{\mathcal{I}} [m_i \leq M/3]]. \end{aligned}$$

The event \mathcal{E}_i is a superset of the event $\mathcal{E}_{i,j'}$ that $(x_i, y_{j_i^*}), (x_i, y_{j'}) \notin \mathcal{D}$ for an arbitrary fixed index $j' \neq j_i^*$. Thus,

$$\mathbb{P}^{\mathcal{I}} [\mathcal{E}_i | m_i \leq M/3] \geq \mathbb{P}^{\mathcal{I}} [\mathcal{E}_{i,j'} | m_i \leq M/3]$$

Moreover, we can realize the law of $\mathbb{P}^{\mathcal{I}}$ considering an infinite tape, associated to index i , of i.i.d. samples $y \sim \pi_{\text{base}}(\cdot | x_i)$, and letting values of y form the samples $(x, y, \log \pi_{\text{base}}(y | x)) \in \mathcal{D}$ with $x = x_i$ corresponding to the first m_i elements on this tape (see, e.g. (Simchowitz et al., 2017) for an argument of this form). On the event $\{m_i \leq M/3\}$, then, m_i samples in $(x, y, \log \pi_{\text{base}}(y | x)) \in \mathcal{D}$ with $x = x_i$ are a subset of the first $M/3$ samples from the index- i tape. Viewed in this way, we can lower bound the probability of $\mathcal{E}_{i,j'}$ of by the probability of the event $\tilde{\mathcal{E}}_{i,j'}$ that the first $M/3$ y 's on the index- i tape contain neither j_i^* , nor the designated index j' . As these first $M/3$ y 's are not chosen adaptively, the probability of $\tilde{\mathcal{E}}_{i,j'}$ is at least

$$\left(1 - \frac{3}{M}\right)^{m_i} \geq \left(1 - \frac{3}{M}\right)^{M/3} \geq \frac{1}{2e},$$

as long as $M \geq 6$ and $\gamma = 1/2$. We conclude that

$$\mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{\hat{j}_i \neq j_i^*\}] \geq \frac{1}{4e} \mathbb{E}_{\mathcal{I} \sim \text{Unif}} \mathbb{E}^{\mathcal{I}} [\mathbb{I}\{m_i \leq M/3\}].$$

□

□

I PROOFS FROM SECTION 4.1 AND APPENDIX C

The following theorem is a generalization of [Theorem 4.1'](#) which allows for approximate maximizers in the sense of [Definition F.1](#).

Theorem 4.1'. *Let $\rho, \delta \in (0, 1)$ be given, and suppose we set $N = N^* \log(2\delta^{-1})$ for a parameter $N^* \in \mathbb{N}$. Then for any $n \in \mathbb{N}$, SFT-Sharpener ensures that with probability at least $1 - \rho$, for any $\gamma \in (0, 1)$, the output model $\hat{\pi}$ satisfies*

$$\mathbb{P}_{x \sim \mu} [\hat{\pi}(\mathbf{y}_\gamma^*(x) | x) \leq 1 - 2\delta] \lesssim \frac{1}{\delta} \cdot \frac{\log(|\Pi|\rho^{-1})}{n} + \frac{C_{\text{cov}, \gamma}}{N^*}.$$

In particular, given $(\epsilon, \delta, \gamma)$, by setting $n = C_{4.1} \frac{\log|\Pi|}{\delta\epsilon}$ and $N^ = C_{4.1} \frac{C_{\text{cov}, \gamma}}{\epsilon}$ for a sufficiently large absolute constant $C_{4.1} > 0$, we are guaranteed that*

$$\mathbb{P}_{x \sim \mu} [\hat{\pi}(\mathbf{y}_\gamma^*(x) | x) \leq 1 - \delta] \leq \epsilon$$

The total sample complexity is

$$m = O\left(\frac{C_{\text{cov}, \gamma} \log(|\Pi|\rho^{-1}) \log(\delta^{-1})}{\delta\epsilon^2}\right).$$

Proof of Theorem 4.1'. Under realizability of π_N^{BoN} ([Assumption 4.1](#)), [Lemma F.1](#) implies that the output of SFT-Sharpener satisfies, with probability at least $1 - \rho$,

$$\mathbb{E}_{x \sim \mu} [D_H^2(\hat{\pi}(\cdot | x), \pi_N^{\text{BoN}}(\cdot | x))] \leq \varepsilon_{\text{stat}}^2 := \frac{2 \log(|\Pi|/\rho)}{n}. \quad (17)$$

Henceforth we condition on the event that [Eq. \(17\)](#) holds. Let

$$\mathcal{X}_{\text{good}} := \left\{ x \in \mathcal{X} \mid N^* \geq \frac{1}{\pi_{\text{base}}(\mathbf{y}_\gamma^*(x) | x)} \right\}$$

denote the set of prompts for which π_{base} places sufficiently high mass on $\mathbf{y}_\gamma^*(x)$. We can bound

$$\begin{aligned} \mathbb{P}_{x \sim \mu} [\hat{\pi}(\mathbf{y}_\gamma^*(x) | x) \leq 1 - \delta] \\ \leq \mathbb{P}_{x \sim \mu} [\hat{\pi}(\mathbf{y}_\gamma^*(x) | x) \leq 1 - \delta, x \in \mathcal{X}_{\text{good}}] + \mathbb{P}_{x \sim \mu} [x \notin \mathcal{X}_{\text{good}}]. \end{aligned} \quad (18)$$

To bound the first term in [Eq. \(18\)](#), note that if $x \in \mathcal{X}_{\text{good}}$, then $\pi_N^{\text{BoN}}(\mathbf{y}_\gamma^*(x) | x) \geq 1 - \delta/2$. Indeed, observe that $y \sim \pi_N^{\text{BoN}}(\cdot | x) \notin \mathbf{y}_\gamma^*(x)$ if and only if $y_1, \dots, y_N \sim \pi_{\text{base}}(x)$ have $y_i \notin \mathbf{y}_\gamma^*(x)$ for all i , which happens with probability $(1 - \pi_{\text{base}}(\mathbf{y}_\gamma^*(x) | x))^N \leq (1 - 1/N^*)^N \leq \delta/2$ since $x \in \mathcal{X}_{\text{good}}$. It follows that for any such x , we can lower bound (using the data processing inequality)

$$\begin{aligned} D_H^2(\hat{\pi}(\cdot | x), \pi_N^{\text{BoN}}(\cdot | x)) &\geq \left(\sqrt{1 - \hat{\pi}(\mathbf{y}_\gamma^*(x) | x)} - \sqrt{1 - \pi_N^{\text{BoN}}(\mathbf{y}_\gamma^*(x) | x)} \right)^2 \\ &\gtrsim \delta \cdot \mathbb{I}\{\hat{\pi}(\mathbf{y}_\gamma^*(x) | x) \leq 1 - \delta\}. \end{aligned} \quad (19)$$

By [Eqs. \(17\) and \(19\)](#), it follows that

$$\mathbb{P}_{x \sim \mu} [\hat{\pi}(\mathbf{y}_\gamma^*(x) | x) \leq 1 - 2\delta, x \in \mathcal{X}_{\text{good}}] \lesssim \frac{\varepsilon_{\text{stat}}^2}{\delta}.$$

For the second term in [Eq. \(18\)](#), we bound

$$\begin{aligned} \mathbb{P}_{x \sim \mu} [x \notin \mathcal{X}_{\text{good}}] &= \mathbb{P}_{x \sim \mu} \left[N^* < \frac{1}{\pi_{\text{base}}(\mathbf{y}_\gamma^*(x) | x)} \right] \\ &= \mathbb{P}_{x \sim \mu} \left[\frac{1}{N^* \pi_{\text{base}}(\mathbf{y}_\gamma^*(x) | x)} > 1 \right] \\ &\leq \frac{1}{N^*} \mathbb{E}_{x \sim \mu} \left[\frac{1}{\pi_{\text{base}}(\mathbf{y}_\gamma^*(x) | x)} \right] \\ &\leq \frac{C_{\text{cov}, \gamma}}{N^*} \end{aligned}$$

via Markov's inequality and the definition of $C_{\text{cov}, \gamma}$. Substituting both bounds into Eq. (18) completes the proof. \square

Proof of Theorem C.1. The proof begins similarly to Theorem 4.1. By realizability of π_{N_μ} , Lemma F.1 implies that the output of SFT-Sharpener satisfies, with probability at least $1 - \rho$,

$$\mathbb{E}_{x \sim \mu} [D_{\text{H}}^2(\hat{\pi}(\cdot | x), \pi_{N_\mu}(\cdot | x))] \leq \varepsilon_{\text{stat}}^2 := \frac{2 \log(|\Pi|/\rho)}{n}.$$

Condition on the event that this guarantee holds. We invoke the following lemma, proven in the sequel.

Lemma I.1. *Let P be a distribution on a discrete space \mathcal{Y} . Let $\mathbf{y}^* = \arg \max_{y \in \mathcal{Y}} P(y)$ and let $P^* := \max_{y \in \mathcal{Y}} P(y)$. Let $y_1, y_2, \dots \sim P$, and for any stopping time τ , define*

$$\hat{y}_\tau \in \arg \max \{P(y) : y \in \{y_1, \dots, y_\tau\}\}.$$

Next, for a parameter $\mu > 0$, define the stopping time

$$N_\mu := \inf \left\{ k : \frac{1}{\max_{1 \leq i \leq k} P(y_i)} \leq k/\mu \right\}.$$

Then

$$\mathbb{E}[N_\mu] \leq \frac{\mu + (1/|\mathbf{y}^*|)}{P^*}.$$

In addition, for any stopping time $\tau \geq N_\mu$ (including $\tau = N_\mu$ itself), we have $\mathbb{P}[\hat{y}_\tau \notin \mathbf{y}^] \leq e^{-|\mathbf{y}^*|\mu}$.*

This lemma, with our choice of μ , ensures that for all $x \in \mathcal{X}$,

$$\pi_{N_\mu}(\mathbf{y}^*(x) | x) \geq 1 - e^{-\mu} = 1 - \delta/2.$$

Following the reasoning in Eq. (19), this implies that

$$D_{\text{H}}^2(\hat{\pi}(\cdot | x), \pi_{N_\mu}(\cdot | x)) \gtrsim \delta \cdot \mathbb{I}\{\hat{\pi}(\mathbf{y}^*(x) | x) \leq 1 - \delta\},$$

so that

$$\mathbb{P}_{x \sim \mu}[\hat{\pi}(\mathbf{y}^*(x) | x) \leq 1 - \delta] \lesssim \frac{\varepsilon_{\text{stat}}^2}{\delta}$$

as desired.

To bound the expected sample complexity, we observe that

$$\mathbb{E}[m] = n \cdot \mathbb{E}[N_\mu(x)] \stackrel{(i)}{\leq} \mathbb{E} \left[\frac{1 + \mu}{\pi_{\text{base}}(\mathbf{y}^*(x) | x)} \right] = (1 + \mu) \bar{C}_{\text{cov}},$$

where inequality (i) invokes Lemma I.1 once more. \square

Proof of Lemma I.1. Define $N^* := \mu/P^*$. To bound the tails of N_μ , define

$$\tau = \inf \{k \mid k \geq N^* \text{ and } \mathbf{y}^* \cap \{y_1, \dots, y_k\} \neq \emptyset\}.$$

It follows from the definition that $N_\mu \leq \tau$, since for any $k \geq N^*$, if there exists $i \leq k$ such that $y_i \in \mathbf{y}^*$, then

$$\frac{1}{P(y_i)} = \frac{1}{P^*} = \frac{N^*}{\mu} \leq \frac{k}{\mu}.$$

Thus, for $k \geq N^*$, we can bound

$$\mathbb{P}[N_\mu > k] \leq \mathbb{P}[\tau > k] = \mathbb{P}[\mathcal{Y}^* \cap \{y_1, \dots, y_k\} = \emptyset] \leq (1 - |\mathbf{y}^*|P^*)^k,$$

and consequently

$$\begin{aligned}\mathbb{E}[N_\mu] &\leq \mathbb{E}[\tau] \leq \mathbb{E}[\tau \mathbb{I}\{\tau \leq N^*\}] + \mathbb{E}[\tau \mathbb{I}\{\tau > N^*\}] \\ &\leq N^* + \sum_{k > N^*} (1 - |\mathbf{y}^*| P^*)^k \\ &\leq N^* + \frac{1}{|\mathbf{y}^*| P(y^*)} = \frac{\mu + 1/|\mathbf{y}^*|}{P(y^*)}.\end{aligned}$$

To check correctness, observe that $N_\mu \geq N^*$, because for all $y \in \mathcal{Y}$, $\frac{1}{P(y)} \geq N^*/\mu$. Hence, any stopping time $\tau \geq N_\mu$ also satisfies $\tau \geq N^*$, and moreover has $\hat{y}_\tau \in \mathbf{y}^*$ whenever $\mathbf{y}^* \cap \{y_1, y_2, \dots, y_\tau\} \neq \emptyset$. This fails to occur with probability no more than

$$\left(1 - \frac{|\mathbf{y}^*|}{P^*}\right)^{N^*} = \left(1 - \frac{|\mathbf{y}^*|}{P^*}\right)^{\mu/P^*} \leq e^{-|\mathbf{y}^*|/\mu}.$$

□

J PROOFS FROM SECTION 4.2

The following result is a generalization of [Lemma 4.1](#).

Lemma 4.1'. *For all $\gamma \in (0, 1)$, the model π_β^* satisfies $\mathcal{C}_{\pi_\beta^*} \leq (1 - \gamma)^{-1} C_{\text{cov}, \gamma}$ and $\mathcal{C}_{\pi_{\text{base}}/\pi_\beta^*; \beta} \leq |\mathcal{Y}|$.*

Proof of Lemma 4.1'. For any fixed $x \in \mathcal{X}$, we have

$$\begin{aligned}\mathbb{E}_{y \sim \pi_\beta^*(\cdot | x)} \left[\frac{\pi_\beta^*(y | x)}{\pi_{\text{base}}(y | x)} \right] &= \mathbb{E}_{y \sim \pi_\beta^*(\cdot | x)} \left[\frac{\pi_{\text{base}}^{1+\beta^{-1}}(y | x)}{\pi_{\text{base}}(y | x)} \right] \cdot \left(\sum_{y' \in \mathcal{Y}} \pi_{\text{base}}^{1+\beta^{-1}}(y' | x) \right)^{-1} \\ &\leq \max_{y \in \mathcal{Y}} \pi_{\text{base}}^{\beta^{-1}}(y | x) \cdot \left(\sum_{y' \in \mathcal{Y}} \pi_{\text{base}}^{1+\beta^{-1}}(y' | x) \right)^{-1} \\ &\leq (1 - \gamma)^{-1} \pi_{\text{base}}^{\beta^{-1}}(\mathbf{y}_\gamma^*(x) | x) \cdot \left(\sum_{y' \in \mathcal{Y}} \pi_{\text{base}}^{1+\beta^{-1}}(y' | x) \right)^{-1} \\ &= (1 - \gamma)^{-1} \frac{\pi_{\text{base}}^{1+\beta^{-1}}(\mathbf{y}_\gamma^*(x) | x)}{\pi_{\text{base}}(\mathbf{y}_\gamma^*(x) | x)} \cdot \left(\sum_{y' \in \mathcal{Y}} \pi_{\text{base}}^{1+\beta^{-1}}(y' | x) \right)^{-1} \\ &= (1 - \gamma)^{-1} \frac{\sum_{y \in \mathbf{y}_\gamma^*(x)} \pi_{\text{base}}^{1+\beta^{-1}}(y | x)}{\pi_{\text{base}}(\mathbf{y}_\gamma^*(x) | x)} \cdot \left(\sum_{y' \in \mathcal{Y}} \pi_{\text{base}}^{1+\beta^{-1}}(y' | x) \right)^{-1} \\ &\leq (1 - \gamma)^{-1} \frac{1}{\pi_{\text{base}}(\mathbf{y}_\gamma^*(x) | x)}.\end{aligned}$$

It follows that $\mathcal{C}_{\pi_\beta^*} \leq (1 - \gamma)^{-1} C_{\text{cov}, \gamma}$ as claimed.

For the second result, we have

$$\mathcal{C}_{\pi_{\text{base}}/\pi_\beta^*; \beta} = \mathbb{E}_{\pi_{\text{base}}} \left[\frac{1}{\pi_{\text{base}}(y | x)} \cdot \left(\sum_{y' \in \mathcal{Y}} \pi_{\text{base}}^{1+\beta^{-1}}(y' | x) \right)^\beta \right] \leq \mathbb{E}_{\pi_{\text{base}}} \left[\frac{1}{\pi_{\text{base}}(y | x)} \right] = |\mathcal{Y}|.$$

□

J.1 PROOF OF THEOREM 4.2

We state and prove a generalized version of [Theorem 4.2](#). In the assumptions below, we fix a parameter $\gamma \in [0, 1)$; the setting $\gamma = 0$ corresponds to [Theorem 4.2](#).

Assumption J.1 (Coverage). *All $\pi \in \Pi$ satisfy $\mathcal{C}_\pi \leq C_{\text{conc}}$ for a parameter $C_{\text{conc}} \geq (1 - \gamma)^{-1} C_{\text{cov}, \gamma}$, and $\mathcal{C}_{\pi_{\text{base}}/\pi; \beta} \leq C_{\text{loss}}$ for a parameter $C_{\text{loss}} \geq |\mathcal{Y}|$.*

By [Lemma 4.1'](#), this assumption is consistent with the assumption that $\pi_\beta^* \in \Pi$.

Assumption J.2 (Margin). *For all $x \in \text{supp}(\mu)$, the initial model π_{base} satisfies*

$$\pi_{\text{base}}(\mathbf{y}_\gamma^*(x) \mid x) \geq (1 + \gamma_{\text{margin}}) \cdot \pi_{\text{base}}(y \mid x) \quad \forall y \notin \mathbf{y}_\gamma^*(x)$$

for a parameter $\gamma_{\text{margin}} > 0$.

Theorem 4.2'. *Assume that $\pi_\beta^* \in \Pi$ ([Assumption 4.3](#)), and that [Assumption 4.4](#) and [Assumption 4.2](#) hold with respect to some $\gamma \in [0, 1)$, with parameters C_{conc} , C_{loss} , and $\gamma_{\text{margin}} > 0$. For any $\delta, \rho \in (0, 1)$, the DPO algorithm in [Eq. \(4\)](#) ensures that with probability at least $1 - \rho$,*

$$\mathbb{P}_{x \sim \mu} [\hat{\pi}(\mathbf{y}_\gamma^*(x) \mid x) \leq 1 - \delta] \lesssim \frac{1}{\gamma_{\text{margin}} \delta} \cdot \tilde{O} \left(\sqrt{\frac{C_{\text{conc}} \log^3(C_{\text{loss}} |\Pi| \rho^{-1})}{n}} + \beta \log(C_{\text{conc}}) + \gamma \right)$$

where $\tilde{O}(\cdot)$ hides factors logarithmic in n and C_{conc} and doubly logarithmic in Π , C_{loss} , and ρ^{-1} .

We first state and prove some supporting technical lemmas, then proceed to the proof of [Theorem 4.2'](#).

J.1.1 TECHNICAL LEMMAS

Lemma J.1. *Suppose $\beta \in [0, 1]$. For any model π , with probability at least $1 - \delta$ over the draw of $x \sim \mu$, $y, y' \sim \pi_{\text{base}}(\cdot \mid x)$, we have that for all $s > 0$,*

$$\mathbb{P} \left[\left| \beta \log \left(\frac{\pi(y \mid x)}{\pi_{\text{base}}(y \mid x)} \right) - \beta \log \left(\frac{\pi(y' \mid x)}{\pi_{\text{base}}(y' \mid x)} \right) \right| > \log(2C_{\pi_{\text{base}}/\pi; \beta}) + s \right] \leq \exp(-s).$$

Proof of Lemma J.1. Define

$$X := \left| \beta \log \left(\frac{\pi(y \mid x)}{\pi_{\text{base}}(y \mid x)} \right) - \beta \log \left(\frac{\pi(y' \mid x)}{\pi_{\text{base}}(y' \mid x)} \right) \right|.$$

By the Chernoff method, we have that with probability at least $1 - \delta$,

$$\begin{aligned} X &\leq \log(\mathbb{E}[\exp(X)]) + \log(\delta^{-1}) \\ &= \log \left(\mathbb{E}_{x \sim \mu, y, y' \sim \pi_{\text{base}}(x)} \left[\exp \left(\left| \beta \log \left(\frac{\pi(y \mid x)}{\pi_{\text{base}}(y \mid x)} \right) - \beta \log \left(\frac{\pi(y' \mid x)}{\pi_{\text{base}}(y' \mid x)} \right) \right| \right) \right] \right) + \log(\delta^{-1}) \\ &\leq \log \left(\mathbb{E}_{x \sim \mu, y, y' \sim \pi_{\text{base}}(x)} \left[\exp \left(\beta \log \left(\frac{\pi(y \mid x)}{\pi_{\text{base}}(y \mid x)} \right) - \beta \log \left(\frac{\pi(y' \mid x)}{\pi_{\text{base}}(y' \mid x)} \right) \right) \right] \right. \\ &\quad \left. + \mathbb{E}_{x \sim \mu, y, y' \sim \pi_{\text{base}}(x)} \left[\exp \left(\beta \log \left(\frac{\pi(y' \mid x)}{\pi_{\text{base}}(y' \mid x)} \right) - \beta \log \left(\frac{\pi(y \mid x)}{\pi_{\text{base}}(y \mid x)} \right) \right) \right] \right) + \log(\delta^{-1}) \\ &= \log \left(2 \mathbb{E}_{x \sim \mu, y, y' \sim \pi_{\text{base}}(x)} \left[\exp \left(\beta \log \left(\frac{\pi(y \mid x)}{\pi_{\text{base}}(y \mid x)} \right) - \beta \log \left(\frac{\pi(y' \mid x)}{\pi_{\text{base}}(y' \mid x)} \right) \right) \right] \right) + \log(\delta^{-1}) \\ &= \log \left(\mathbb{E}_{x \sim \mu, y, y' \sim \pi_{\text{base}}(x)} \left[\left(\frac{\pi(y \mid x)}{\pi_{\text{base}}(y \mid x)} \cdot \frac{\pi_{\text{base}}(y' \mid x)}{\pi(y' \mid x)} \right)^\beta \right] \right) + \log(2\delta^{-1}). \end{aligned}$$

As long as $\beta \leq 1$, by Jensen's inequality, we can bound

$$\begin{aligned}
& \mathbb{E}_{x \sim \mu, y, y' \sim \pi_{\text{base}}(x)} \left[\left(\frac{\pi(y | x)}{\pi_{\text{base}}(y | x)} \cdot \frac{\pi_{\text{base}}(y' | x)}{\pi(y' | x)} \right)^\beta \right] \\
& \leq \mathbb{E}_{x \sim \mu, y, y' \sim \pi_{\text{base}}(x)} \left[\left(\mathbb{E}_{y \sim \pi_{\text{base}}(x)} \left[\frac{\pi(y | x)}{\pi_{\text{base}}(y | x)} \right] \cdot \frac{\pi_{\text{base}}(y' | x)}{\pi(y' | x)} \right)^\beta \right] \\
& = \mathbb{E}_{x \sim \mu, y' \sim \pi_{\text{base}}(x)} \left[\left(\frac{\pi_{\text{base}}(y' | x)}{\pi(y' | x)} \right)^\beta \right] \\
& = \mathcal{C}_{\pi_{\text{base}}/\pi; \beta},
\end{aligned}$$

which proves the result. \square

Lemma J.2. Let $\beta \in [0, 1]$. For all models π , we have

$$\mathbb{E}_{x \sim \mu, y, y' \sim \pi_{\text{base}}(\cdot | x)} \left[\left| \beta \log \left(\frac{\pi(y | x)}{\pi_{\text{base}}(y | x)} \right) - \beta \log \left(\frac{\pi(y' | x)}{\pi_{\text{base}}(y' | x)} \right) \right|^4 \right] \leq O(\log^4(\mathcal{C}_{\pi_{\text{base}}/\pi; \beta}) + 1).$$

Proof of Lemma J.2. Define

$$X := \left| \beta \log \left(\frac{\pi(y | x)}{\pi_{\text{base}}(y | x)} \right) - \beta \log \left(\frac{\pi(y' | x)}{\pi_{\text{base}}(y' | x)} \right) \right|.$$

Set $k = \log(2\mathcal{C}_{\pi_{\text{base}}/\pi; \beta})$. We can bound

$$\begin{aligned}
\mathbb{E}[X^4] &= \mathbb{E} \left[\int_0^\infty \mathbb{I}\{X^4 > t\} dt \right] \\
&= 4 \mathbb{E} \left[\int_0^\infty \mathbb{I}\{X > t\} t^3 dt \right] \\
&= 4 \int_0^\infty \mathbb{P}[X > t] t^3 dt \\
&\leq k^4 + 4 \int_k^\infty \mathbb{P}[X > t] t^3 dt \\
&\leq k^4 + 4 \int_k^\infty e^{k-t} t^3 dt \\
&= k^4 + 4(k^3 + 3k^2 + 6k + 6) \\
&= O(k^4 + 1),
\end{aligned}$$

where the third-to-last line uses Lemma J.1. \square

J.1.2 PROOF OF THEOREM 4.2'

Proof of Theorem 4.2'. For any model $\pi \in \Pi$, define $J(\pi) := \mathbb{E}_\pi[\log \pi_{\text{base}}(y | x)]$. Let $\hat{\pi} \in \Pi$ denote the model returned by the DPO algorithm in Eq. (8). Let $\mathbb{E}_{\pi, \pi'}[\cdot]$ denote shorthand for $\mathbb{E}_{x \sim \mu, y \sim \pi(x), y' \sim \pi'(x)}[\cdot]$, and for any $r : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ define $\Delta^r(x, y, y') := r(x, y) - r(x, y')$. Define

$$r^*(x, y) := \log \pi_{\text{base}}(y | x) = \beta \log \left(\frac{\pi_\beta^*(y | x)}{\pi_{\text{base}}(y | x)} \right) + Z(x),$$

and let $\hat{r}(x, y) := \beta \log \left(\frac{\hat{\pi}(y | x)}{\pi_{\text{base}}(y | x)} \right)$. By a standard argument (Huang et al., 2024), we have

$$\hat{\pi} \in \arg \max_{\pi : \mathcal{X} \rightarrow \Delta(\mathcal{Y})} \mathbb{E}_\pi[\hat{r}(x, y)] - \beta D_{\text{KL}}(\pi \| \pi_{\text{base}}). \quad (20)$$

Therefore for any comparator model $\pi^* : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$ (not necessarily in the model class Π), we have

$$\begin{aligned}
J(\pi^*) - J(\hat{\pi}) &= \mathbb{E}_{\pi^*}[r^*(x, y)] - \mathbb{E}_{\hat{\pi}}[r^*(x, y)] \\
&= \mathbb{E}_{\pi^*}[\hat{r}(x, y)] - \beta D_{\text{KL}}(\pi^* \parallel \pi_{\text{base}}) - \mathbb{E}_{\hat{\pi}}[\hat{r}(x, y)] + \beta D_{\text{KL}}(\hat{\pi} \parallel \pi_{\text{base}}) \\
&\quad + \mathbb{E}_{\pi^*}[r^*(x, y) - \hat{r}(x, y)] + \beta D_{\text{KL}}(\pi^* \parallel \pi_{\text{base}}) + \mathbb{E}_{\hat{\pi}}[\hat{r}(x, y) - r^*(x, y)] - \beta D_{\text{KL}}(\hat{\pi} \parallel \pi_{\text{base}}) \\
&\leq \mathbb{E}_{\pi^*}[r^*(x, y) - \hat{r}(x, y)] + \beta D_{\text{KL}}(\pi^* \parallel \pi_{\text{base}}) + \mathbb{E}_{\hat{\pi}}[\hat{r}(x, y) - r^*(x, y)] - \beta D_{\text{KL}}(\hat{\pi} \parallel \pi_{\text{base}}) \\
&= \mathbb{E}_{\pi^*, \pi_{\text{base}}}[\Delta^{r^*}(x, y, y') - \Delta^{\hat{r}}(x, y, y')] + \mathbb{E}_{\hat{\pi}, \pi_{\text{base}}}[\Delta^{\hat{r}}(x, y, y') - \Delta^{r^*}(x, y, y')] \\
&\quad + \beta D_{\text{KL}}(\pi^* \parallel \pi_{\text{base}}) - \beta D_{\text{KL}}(\hat{\pi} \parallel \pi_{\text{base}}) \tag{21}
\end{aligned}$$

where the inequality uses [Eq. \(20\)](#). To bound the right-hand-side above, we will use the following lemma, which is proven in the sequel.

Lemma J.3. *For any model π and any $\eta > 0$, we have that*

$$\begin{aligned}
&\mathbb{E}_{\pi, \pi_{\text{base}}} \left[\left| \Delta^{r^*}(x, y, y') - \Delta^{\hat{r}}(x, y, y') \right| \right] \\
&\lesssim \mathcal{C}_{\pi}^{1/2} \cdot \left(\mathbb{E}_{\pi_{\text{base}}, \pi_{\text{base}}} \left[\left| \Delta^{r^*}(x, y, y') - \Delta^{\hat{r}}(x, y, y') \right|^2 \mathbb{I}\{|\Delta^{r^*}| \leq \eta, |\Delta^{\hat{r}}| \leq \eta\} \right] \right)^{1/2} \\
&\quad + \mathcal{C}_{\pi}^{1/2} (\log(\mathcal{C}_{\pi_{\text{base}}/\hat{\pi}; \beta}) + \log(\mathcal{C}_{\pi_{\text{base}}/\pi_{\hat{\beta}}; \beta})) \cdot \left(\mathbb{P}_{\pi_{\text{base}}, \pi_{\text{base}}} [|\Delta^{r^*}| > \eta] + \mathbb{P}_{\pi_{\text{base}}, \pi_{\text{base}}} [|\Delta^{\hat{r}}| > \eta] \right)^{1/4}.
\end{aligned}$$

Using [Lemma J.3](#) to bound the first two terms of [Eq. \(21\)](#), and using the fact that all $\pi \in \Pi$ have $\mathcal{C}_{\pi} \leq C_{\text{conc}}$ and $\mathcal{C}_{\pi_{\text{base}}/\pi; \beta} \leq C_{\text{loss}}$, we have that

$$\begin{aligned}
J(\pi^*) - J(\hat{\pi}) &\lesssim (\mathcal{C}_{\pi^*} + C_{\text{conc}})^{1/2} \cdot \left(\mathbb{E}_{\pi_{\text{base}}, \pi_{\text{base}}} \left[\left| \Delta^{r^*}(x, y, y') - \Delta^{\hat{r}}(x, y, y') \right|^2 \mathbb{I}\{|\Delta^{r^*}| \leq \eta, |\Delta^{\hat{r}}| \leq \eta\} \right] \right)^{1/2} \\
&\quad + (\mathcal{C}_{\pi^*} + C_{\text{conc}})^{1/2} \log(C_{\text{loss}}) \cdot \left(\mathbb{P}_{\pi_{\text{base}}, \pi_{\text{base}}} [|\Delta^{r^*}| > \eta] + \mathbb{P}_{\pi_{\text{base}}, \pi_{\text{base}}} [|\Delta^{\hat{r}}| > \eta] \right)^{1/4} + \beta D_{\text{KL}}(\pi^* \parallel \pi_{\text{base}}). \tag{22}
\end{aligned}$$

Let us overload notation and write $\Delta^{\pi}(x, y, y') = \beta \log\left(\frac{\pi(y|x)}{\pi_{\text{base}}(y|x)}\right) - \beta \log\left(\frac{\pi(y'|x)}{\pi_{\text{base}}(y'|x)}\right)$, so that $\Delta^{\hat{\pi}} = \Delta^{\hat{r}}$ and $\Delta^{\pi_{\hat{\beta}}} = \Delta^{r^*}$. Since $\pi_{\hat{\beta}} \in \Pi$, the definition of $\hat{\pi}$ in [Eq. \(4\)](#) implies that

$$\begin{aligned}
\sum_{(x, y, y') \in \mathcal{D}_{\text{pref}}} \left(\Delta^{\hat{\pi}}(x, y, y') - \Delta^{\pi_{\hat{\beta}}}(x, y, y') \right)^2 &\leq \min_{\pi \in \Pi} \sum_{(x, y, y') \in \mathcal{D}_{\text{pref}}} \left(\Delta^{\pi}(x, y, y') - \Delta^{\pi_{\hat{\beta}}}(x, y, y') \right)^2 \\
&\leq \sum_{(x, y, y') \in \mathcal{D}_{\text{pref}}} \left(\Delta^{\pi_{\hat{\beta}}}(x, y, y') - \Delta^{\pi_{\hat{\beta}}}(x, y, y') \right)^2 \\
&= 0.
\end{aligned}$$

Define $B_{n, \rho} := \log(2nC_{\text{loss}}|\Pi|\rho^{-1})$. It is immediate that

$$\sum_{(x, y, y') \in \mathcal{D}_{\text{pref}}} \left(\Delta^{\hat{\pi}}(x, y, y') - \Delta^{\pi_{\hat{\beta}}}(x, y, y') \right)^2 \mathbb{I}\{|\Delta^{\hat{\pi}}| \leq B_{n, \rho}, |\Delta^{\pi_{\hat{\beta}}}| \leq B_{n, \rho}\} \leq 0.$$

From here, Bernstein's inequality and a union bound implies that with probability at least $1 - \rho$,

$$\begin{aligned}
&\mathbb{E}_{\pi_{\text{base}}, \pi_{\text{base}}} \left[\left| \Delta^{\hat{\pi}}(x, y, y') - \Delta^{\pi_{\hat{\beta}}}(x, y, y') \right|^2 \mathbb{I}\{|\Delta^{\hat{\pi}}| \leq B_{n, \rho}, |\Delta^{\pi_{\hat{\beta}}}| \leq B_{n, \rho}\} \right] \\
&\lesssim \frac{B_{n, \rho}^2 \log(|\Pi|\rho^{-1})}{n} =: \varepsilon_{\text{stat}}^2.
\end{aligned}$$

In particular, if we combine this with [Eq. \(22\)](#) and set $\eta = B_{n, \rho}$, then [Lemma J.1](#) implies that

$$J(\pi^*) - J(\hat{\pi}) \lesssim (\mathcal{C}_{\pi^*} + C_{\text{conc}})^{1/2} \cdot \varepsilon_{\text{stat}} + (\mathcal{C}_{\pi^*} + C_{\text{conc}})^{1/2} \log(C_{\text{loss}}) \cdot \rho^{1/4} + \beta D_{\text{KL}}(\pi^* \parallel \pi_{\text{base}}).$$

Note that the above bound holds for any $\pi^* : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$. We define π^* by

$$\pi^*(y | x) := \frac{\pi_{\text{base}}(y | x) \mathbb{I}[y \in \mathbf{y}_\gamma^*(x)]}{\pi_{\text{base}}(\mathbf{y}_\gamma^*(x) | x)},$$

which can be seen to satisfy $\mathcal{C}_{\pi^*} \leq C_{\text{cov}, \gamma} \leq C_{\text{conc}}$ and $D_{\text{KL}}(\pi^* \| \pi_{\text{base}}) \leq \log(\mathcal{C}_{\pi^*}) \leq \log(C_{\text{conc}})$. With this choice, we can further bound the expression above by

$$J(\pi^*) - J(\hat{\pi}) \lesssim (C_{\text{conc}})^{1/2} \cdot \varepsilon_{\text{stat}} + (C_{\text{conc}})^{1/2} \log(C_{\text{loss}}) \cdot \rho^{1/4} + \beta \log(C_{\text{conc}})$$

Given a desired failure probability ρ , applying the bound above with $\rho' := \rho \wedge (\varepsilon_{\text{stat}} / \log(C_{\text{loss}}))^4$ then gives

$$J(\pi^*) - J(\hat{\pi}) \lesssim (C_{\text{conc}})^{1/2} \cdot \varepsilon_{\text{stat}} + \beta \log(C_{\text{conc}}).$$

Finally, we observe that for our choice of π^* , under the margin condition with parameter γ , we have

$$\begin{aligned} J(\pi^*) - J(\hat{\pi}) &= \mathbb{E}_{x \sim \mu} \mathbb{E}_{y, y' \sim \pi^*, \hat{\pi}} \left[\log \left(\frac{\pi_{\text{base}}(y | x)}{\pi_{\text{base}}(y' | x)} \right) \right] \\ &\gtrsim \gamma_{\text{margin}} \cdot \mathbb{E}_{x \sim \mu} \mathbb{E}_{y' \sim \hat{\pi}} [\mathbb{I}\{y' \notin \mathbf{y}_\gamma^*(x)\}] - \gamma \\ &\gtrsim \gamma_{\text{margin}} \delta \cdot \mathbb{E}_{x \sim \mu} [\mathbb{I}\{\hat{\pi}(\mathbf{y}_\gamma^*(x) | x) \leq 1 - \delta\}] - \gamma \end{aligned}$$

where the first inequality uses [Assumption J.2](#) together with the fact that $y \in \mathbf{y}_\gamma^*(x)$ with probability 1 over $x \sim \mu$ and $y \sim \pi^*(\cdot | x)$. This proves the result. \square

Proof of Lemma J.3. For any $\eta > 0$, we can bound

$$\begin{aligned} \mathbb{E}_{\pi, \pi_{\text{base}}} \left[\left| \Delta^{r^*}(x, y, y') - \Delta^{\hat{r}}(x, y, y') \right| \right] &\leq \mathbb{E}_{\pi, \pi_{\text{base}}} \left[\left| \Delta^{r^*}(x, y, y') - \Delta^{\hat{r}}(x, y, y') \right| \mathbb{I}\left\{ |\Delta^{r^*}| \leq \eta, |\Delta^{\hat{r}}| \leq \eta \right\} \right] \\ &\quad + \mathbb{E}_{\pi, \pi_{\text{base}}} \left[\left| \Delta^{r^*}(x, y, y') - \Delta^{\hat{r}}(x, y, y') \right| \mathbb{I}\left\{ |\Delta^{r^*}| > \eta \vee |\Delta^{\hat{r}}| > \eta \right\} \right]. \end{aligned}$$

For the second term above, we can use Cauchy-Schwarz to bound

$$\begin{aligned} &\mathbb{E}_{\pi, \pi_{\text{base}}} \left[\left| \Delta^{r^*}(x, y, y') - \Delta^{\hat{r}}(x, y, y') \right| \mathbb{I}\left\{ |\Delta^{r^*}| > \eta \vee |\Delta^{\hat{r}}| > \eta \right\} \right] \\ &\leq \mathcal{C}_\pi^{1/2} \cdot \left(\mathbb{E}_{\pi_{\text{base}}, \pi_{\text{base}}} \left[\left| \Delta^{r^*}(x, y, y') - \Delta^{\hat{r}}(x, y, y') \right|^2 \mathbb{I}\left\{ |\Delta^{r^*}| > \eta \vee |\Delta^{\hat{r}}| > \eta \right\} \right] \right)^{1/2} \\ &\lesssim \mathcal{C}_\pi^{1/2} \cdot \left(\mathbb{P}_{\pi_{\text{base}}, \pi_{\text{base}}} \left[|\Delta^{r^*}| > \eta \right] + \mathbb{P}_{\pi_{\text{base}}, \pi_{\text{base}}} \left[|\Delta^{\hat{r}}| > \eta \right] \right)^{1/4} \\ &\quad \cdot \left(\mathbb{E}_{\pi_{\text{base}}, \pi_{\text{base}}} \left[\left| \Delta^{r^*}(x, y, y') \right|^4 \right] + \mathbb{E}_{\pi_{\text{base}}, \pi_{\text{base}}} \left[\left| \Delta^{\hat{r}}(x, y, y') \right|^4 \right] \right)^{1/4} \\ &\lesssim \mathcal{C}_\pi^{1/2} \cdot \left(\mathbb{P}_{\pi_{\text{base}}, \pi_{\text{base}}} \left[|\Delta^{r^*}| > \eta \right] + \mathbb{P}_{\pi_{\text{base}}, \pi_{\text{base}}} \left[|\Delta^{\hat{r}}| > \eta \right] \right)^{1/4} \cdot (\log(\mathcal{C}_{\pi_{\text{base}}/\pi; \beta}) + \log(\mathcal{C}_{\pi_{\text{base}}/\pi^*; \beta})), \end{aligned}$$

where the last inequality follows from [Lemma J.2](#).

Meanwhile, for the first term, for any $\lambda > 0$ we can bound

$$\begin{aligned} &\mathbb{E}_{\pi, \pi_{\text{base}}} \left[\left| \Delta^{r^*}(x, y, y') - \Delta^{\hat{r}}(x, y, y') \right| \mathbb{I}\left\{ |\Delta^{r^*}| \leq \eta, |\Delta^{\hat{r}}| \leq \eta \right\} \right] \\ &\leq \mathcal{C}_\pi^{1/2} \left(\mathbb{E}_{\pi_{\text{base}}, \pi_{\text{base}}} \left[\left| \Delta^{r^*}(x, y, y') - \Delta^{\hat{r}}(x, y, y') \right|^2 \mathbb{I}\left\{ |\Delta^{r^*}| \leq \eta, |\Delta^{\hat{r}}| \leq \eta \right\} \right] \right)^{1/2}. \end{aligned}$$

\square

J.2 PROOF OF THEOREM 4.3 AND [UNDEFINED]

In this section we prove [Theorem 4.3](#) as well as ??, the application to linear softmax models. For the formal theorem statements, see [Theorem J.2](#) and [Theorem J.3](#) respectively. The section is organized as follows.

- In [Appendix J.2.1](#), we give necessary background on KL-regularized policy optimization, as well as the Sequential Extrapolation Coefficient.
- [Appendix J.2.2](#) presents a generic guarantee for XPO under a general choice of reward function.
- [Appendix J.2.3](#) instantiates the result above with the self-reward function $r(x, y) := \log \pi_{\text{base}}(y | x)$ to prove [Theorem 4.3](#).
- Finally, [Appendix J.2.4](#) applies the preceding results to prove ??.

J.2.1 BACKGROUND

To begin, we give background on KL-regularized policy optimization and the Sequential Extrapolation Coefficient.

KL-regularized policy optimization. Let $\beta > 0$ be given, and let $r : \mathcal{X} \times \mathcal{Y} \rightarrow [-R_{\max}, R_{\max}]$ be an unknown reward function on prompt/action pairs. Define a value function J_β over model class Π by:

$$J_\beta(\pi) := \mathbb{E}_\pi[r(x, y)] - \beta \cdot D_{\text{KL}}(\mathbb{P}^\pi \parallel \mathbb{P}^{\pi_{\text{base}}}).$$

We refer to this as a *KL-regularized policy optimization* objective (we use the term “policy” following the reinforcement learning literature; for our setting, policies correspond to models). Given query access to r , the goal is to find $\hat{\pi} \in \Pi$ such that

$$J_\beta(\pi_\beta^*) - J_\beta(\hat{\pi}) \leq \epsilon$$

where $\pi_\beta^*(y | x) \propto \pi_{\text{base}}(y | x) \exp(\beta^{-1}r(x, y))$ is the model that maximizes J_β over all models $\pi : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$.

We make use of the following assumptions, as in [Xie et al. \(2024\)](#).

Assumption J.3 (Realizability). *It holds that $\pi_\beta^* \in \Pi$.*

Assumption J.4 (Bounded density ratios). *For all $\pi \in \Pi$, $(x, y) \in \mathcal{X} \times \mathcal{Y}$, $|\beta \log \frac{\pi(y|x)}{\pi_{\text{base}}(y|x)}| \leq V_{\max}$.*

Finally, we require two definitions.

Definition J.1 (Sequential Extrapolation Coefficient for RLHF, [\(Xie et al., 2024\)](#)). *For a model class Π , reward function r , reference model π_{base} , and parameters $T \in \mathbb{N}$ and $\beta, \lambda > 0$, the Sequential Extrapolation Coefficient is defined as*

$$\text{SEC}(\Pi, r, T, \beta, \lambda; \pi_{\text{base}})$$

$$:= \sup_{\pi^{(1)}, \dots, \pi^{(T)} \in \Pi} \left\{ \sum_{t=1}^T \frac{\mathbb{E}^{(t)} \left[\beta \log \frac{\pi^{(t)}(y|x)}{\pi_{\text{base}}(y|x)} - r(x, y) - \beta \log \frac{\pi^{(t)}(y'|x)}{\pi_{\text{base}}(y'|x)} + r(x, y') \right]^2}{\lambda \vee \sum_{i=1}^{t-1} \mathbb{E}^{(i)} \left[\left(\beta \log \frac{\pi^{(i)}(y|x)}{\pi_{\text{base}}(y|x)} - r(x, y) - \beta \log \frac{\pi^{(i)}(y'|x)}{\pi_{\text{base}}(y'|x)} + r(x, y') \right)^2 \right]} \right\}$$

where $\mathbb{E}^{(t)}$ denotes expectation over $x \sim \mu$, $y \sim \pi^{(t)}(\cdot | x)$, and $y' \sim \pi_{\text{base}}(\cdot | x)$.

Definition J.2. *Let $\epsilon > 0$. We say that $\Psi \subseteq \Pi$ is a ϵ -net for model class Π if for every $\pi \in \Pi$ there exists $\pi' \in \Psi$ such that*

$$\max_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} \left| \log \frac{\pi(y | x)}{\pi'(y | x)} \right| \leq \epsilon.$$

We write $\mathcal{N}(\Pi, \epsilon)$ to denote the size of the smallest ϵ -net for Π .

Algorithm 1 Reward-based variant of Exploratory Preference Optimization (Xie et al., 2024)

input: Base model $\pi_{\text{base}} : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$, reward function $r : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, number of iterations $T \in \mathbb{N}$, KL regularization coefficient $\beta > 0$, optimism coefficient $\alpha > 0$.

Initialize: $\pi^{(1)} \leftarrow \pi_{\text{base}}, \mathcal{D}^{(0)} \leftarrow \emptyset$.

for iteration $t = 1, \dots, T$ **do**

Generate sample: $(x^{(t)}, y^{(t)}, \tilde{y}^{(t)})$ via $x^{(t)} \sim \mu, y^{(t)} \sim \pi^{(t)}(\cdot | x^{(t)}), \tilde{y}^{(t)} \sim \pi_{\text{base}}(\cdot | x^{(t)})$.

Update dataset: $\mathcal{D}^{(t)} \leftarrow \mathcal{D}^{(t-1)} \cup \{(x^{(t)}, y^{(t)}, \tilde{y}^{(t)})\}$.

Model optimization with global optimism:

$$\begin{aligned} \pi^{(t+1)} \leftarrow \arg \min_{\pi \in \Pi} & \left\{ \alpha \sum_{(x, y, y') \in \mathcal{D}^{(t)}} \log(\pi(y' | x)) \right. \\ & \left. - \sum_{(x, y, y') \in \mathcal{D}^{(t)}} \left(\beta \log \frac{\pi(y | x)}{\pi_{\text{base}}(y | x)} - \beta \log \frac{\pi(y' | x)}{\pi_{\text{base}}(y' | x)} - (r(x, y) - r(x, y')) \right)^2 \right\}. \end{aligned}$$

return: $\hat{\pi} \leftarrow \arg \max_{t \in [T+1]} J_{\beta}(\pi^{(t)})$. ▷ Can estimate $J_{\beta}(\pi^{(t)})$ using validation data.

J.2.2 GUARANTEES FOR KL-REGULARIZED POLICY OPTIMIZATION WITH XPO

In this section, we give self-contained guarantees for the XPO algorithm (Algorithm 1). XPO was introduced in Xie et al. (2024) for KL-regularized policy optimization in the related setting where the learner only has indirect access to the reward function r through *preference data* (specifically, pairs of actions labeled via a Bradley-Terry model). Standard offline algorithms for this problem, such as DPO, require bounds on concentrability of the model class (see e.g. Eq. (9)). Xie et al. (2024) show that the XPO algorithm avoids this dependence, and instead requires bounded Sequential Extrapolation Coefficient.

Algorithm 1 is a variant of the XPO algorithm which is adapted to reward-based feedback (as opposed to preference-based feedback), and Theorem J.1 shows that this algorithm enjoys guarantees similar to those of Xie et al. (2024) for this setting. Note that this is not an immediate corollary of the results in Xie et al. (2024), since the sample complexity in the preference-based setting scales with $e^{O(R_{\max})}$, and for our application to sharpening it is important to avoid this dependence. However, our algorithm and analysis only diverge from Xie et al. (2024) in a few places.

Theorem J.1 (Variant of Xie et al. (2024, Theorem 3.1)). *Suppose that Assumptions J.3 and J.4 hold. For any $T \in \mathbb{N}$, $\epsilon_{\text{disc}}, \rho \in (0, 1)$, by setting $\alpha := \frac{\beta}{R_{\max} + V_{\max}} \sqrt{\frac{\log(2\mathcal{N}(\Pi, \epsilon_{\text{disc}})T/\rho)}{\text{SEC}(\Pi)T}}$, Algorithm 1 produces a model $\hat{\pi} \in \Pi$ such that with probability at least $1 - \rho$,*

$$\begin{aligned} \beta D_{\text{KL}}(\hat{\pi} \| \pi_{\beta}^*) &= J_{\beta}(\pi_{\beta}^*) - J_{\beta}(\hat{\pi}) \lesssim (R_{\max} + V_{\max}) \sqrt{\frac{\text{SEC}(\Pi) \log(2\mathcal{N}(\Pi, \epsilon_{\text{disc}})T/\rho)}{T}} \\ &\quad + \beta \epsilon_{\text{disc}} \sqrt{\text{SEC}(\Pi)T} \end{aligned}$$

where $\text{SEC}(\Pi) := \text{SEC}(\Pi, r, T, \beta, V_{\max}^2; \pi_{\text{base}})$.

Proof of Theorem J.1. For compactness, we abbreviate $\text{SEC}(\Pi) := \text{SEC}(\Pi, r, T, \beta, V_{\max}^2; \pi_{\text{base}})$. From Equation (37) of Xie et al. (2024), we have

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T J_{\beta}(\pi_{\beta}^*) - J_{\beta}(\pi^{(t)}) \\ & \lesssim \frac{\alpha}{\beta} (R_{\max} + V_{\max})^2 \cdot \text{SEC}(\Pi) + \frac{\beta}{\alpha T} + \frac{V_{\max}}{T} + \frac{1}{T} \sum_{t=2}^T \mathbb{E}_{(x, y) \sim \pi_{\text{base}}} [\beta \log \pi^{(t)}(y | x) - \beta \log \pi_{\beta}^*(y | x)] \\ & \quad + \frac{\beta}{\alpha (R_{\max} + V_{\max})^2 T} \sum_{t=2}^T \mathbb{E}_{\substack{x \sim \mu \\ y, y' \sim \bar{\pi}^{(t)} | x}} \left[\left(\beta \log \frac{\pi^{(t)}(y | x)}{\pi_{\text{base}}(y | x)} - r(x, y) - \beta \log \frac{\pi^{(t)}(y' | x)}{\pi_{\text{base}}(y' | x)} + r(x, y') \right)^2 \right] \end{aligned}$$

where $\bar{\pi}^{(t)} := \frac{1}{t-1} \sum_{i < t} \pi^{(i)} \otimes \pi_{\text{base}}$ denotes the model that, given $x \in \mathcal{X}$, samples $i \sim \text{Unif}([t-1])$ and then samples $y \sim \pi^{(i)}(\cdot | x)$ and $y' \sim \pi_{\text{base}}(\cdot | x)$. For any $2 \leq t \leq T$, define $L^{(t)} : \Pi \rightarrow [0, \infty)$

by

$$L^{(t)}(\pi) := \mathbb{E}_{(x,y) \sim \pi_{\text{base}}} [\beta \log \pi(y | x) - \beta \log \pi_{\beta}^*(y | x)] \\ + \frac{\beta}{\alpha(V_{\max} + R_{\max})^2} \mathbb{E}_{\substack{x \sim \mu \\ y, y' \sim \tilde{\pi}^{(t)} | x}} \left[\left(\beta \log \frac{\pi(y | x)}{\pi_{\text{base}}(y | x)} - r(x, y) - \beta \log \frac{\pi(y' | x)}{\pi_{\text{base}}(y' | x)} + r(x, y') \right)^2 \right].$$

Similarly, define

$$\hat{L}^{(t)}(\pi) := \sum_{(x,y,y') \in \mathcal{D}^{(t)}} [\beta \log \pi(y' | x) - \beta \log \pi_{\beta}^*(y' | x)] \\ + \frac{\beta}{\alpha(V_{\max} + R_{\max})^2} \sum_{(x,y,y') \in \mathcal{D}^{(t)}} \left[\left(\beta \log \frac{\pi(y | x)}{\pi_{\text{base}}(y | x)} - r(x, y) - \beta \log \frac{\pi(y' | x)}{\pi_{\text{base}}(y' | x)} + r(x, y') \right)^2 \right]$$

where $\mathcal{D}^{(t)}$ is the dataset defined in iteration t of [Algorithm 1](#). By [Assumption J.3](#) we have $\pi_{\beta}^* \in \Pi$, so $\inf_{\pi \in \Pi} \hat{L}^{(t)}(\pi) \leq 0$. Moreover by definition, $\pi^{(t)} \in \arg \min_{\pi \in \Pi} \hat{L}^{(t)}$.

Let Ψ be an ϵ_{disc} -net over Π , of size $\mathcal{N}(\Pi, \epsilon_{\text{disc}})$. Fix any $\pi \in \Psi$ and $2 \leq t \leq T$, and define increments $X_i := \hat{L}^{(i)}(\pi) - \hat{L}^{(i-1)}(\pi)$ for $2 \leq i \leq t$, with the notation $\hat{L}^{(1)}(\pi) := 0$ so that $\hat{L}^{(t)}(\pi) = \sum_{i=2}^t X_i$. Let \mathcal{F}_i be the filtration induced by $\mathcal{D}^{(i)}$ and define $\gamma_i := \mathbb{E}[X_i | \mathcal{F}_{i-1}]$. Observe that $(t-1)L^{(t)}(\pi) = \sum_{i=2}^t \gamma_i$. For any i , note that we can write $X_i = Y_i + Z_i$ where $Y_i \in [-V_{\max}, V_{\max}]$ and $Z_i \in [0, \beta/\alpha]$. By [Corollary F.1](#), it holds with probability at least $1 - \rho/(2|\Pi|T)$

$$\sum_{i=2}^t \mathbb{E}[Z_i | \mathcal{F}_{i-1}] \lesssim \frac{\beta}{\alpha} \log(2|\Psi|T/\rho) + \sum_{i=2}^t Z_i.$$

By Azuma-Hoeffding, it holds with probability at least $1 - \rho/(2|\Pi|T)$ that

$$\sum_{i=2}^t \mathbb{E}[Y_i | \mathcal{F}_{i-1}] \lesssim V_{\max} \sqrt{T \log(2|\Psi|T/\rho)} + \sum_{i=2}^t Y_i.$$

Hence, with probability at least $1 - \rho/(|\Psi|T)$ we have

$$(t-1)L^{(t)}(\pi) \lesssim \frac{\beta}{\alpha} \log(2|\Psi|T/\rho) + V_{\max} \sqrt{T \log(2|\Psi|T/\rho)} + \hat{L}^{(t)}(\pi).$$

With probability at least $1 - \rho$ this bound holds for all $\pi \in \Psi$ and $2 \leq t \leq T$. Henceforth condition on this event. Fix any $\pi \in \Pi$ and $2 \leq t \leq T$. Since Ψ is an ϵ -net for Π , we see by definition of $L^{(t)}$ that there is some $\pi' \in \Psi$ such that

$$|L^{(t)}(\pi) - L^{(t)}(\pi')| \lesssim \beta \epsilon_{\text{disc}} + \frac{\beta}{\alpha(V_{\max} + R_{\max})^2} \cdot \beta \epsilon_{\text{disc}} (V_{\max} + R_{\max}) \leq \beta \epsilon_{\text{disc}} \left(1 + \frac{\beta}{\alpha(V_{\max} + R_{\max})} \right)$$

and similarly

$$|\hat{L}^{(t)}(\pi) - \hat{L}^{(t)}(\pi')| \lesssim (t-1)\beta \epsilon_{\text{disc}} \left(1 + \frac{\beta}{\alpha(V_{\max} + R_{\max})} \right).$$

It follows that, for all $2 \leq t \leq T$, since $\hat{L}^{(t)}(\pi^{(t)}) \leq 0$, we get

$$(t-1)L^{(t)}(\pi^{(t)}) \lesssim \frac{\beta}{\alpha} \log(2|\Psi|T/\rho) + V_{\max} \sqrt{T \log(2|\Psi|T/\rho)} + \beta \epsilon_{\text{disc}} T \left(1 + \frac{\beta}{\alpha(V_{\max} + R_{\max})} \right).$$

Hence,

$$\frac{1}{T} \sum_{t=1}^T J_{\beta}(\pi_{\beta}^*) - J_{\beta}(\pi^{(t)}) \\ \lesssim \frac{\alpha}{\beta} (R_{\max} + V_{\max})^2 \cdot \text{SEC}(\Pi) + \frac{\beta}{\alpha T} + \frac{V_{\max}}{T} + \frac{1}{T} \sum_{t=2}^T L^{(t)}(\pi^{(t)}) \\ \lesssim (R_{\max} + V_{\max}) \sqrt{\frac{\text{SEC}(\Pi) \log(2|\Psi|T/\rho)}{T}} + \beta \epsilon_{\text{disc}} \sqrt{\text{SEC}(\Pi) T}$$

by taking

$$\alpha := \frac{\beta}{R_{\max} + V_{\max}} \sqrt{\frac{\log(2|\Psi|T/\rho)}{\text{SEC}(\Pi)T}}.$$

Since the output $\hat{\pi}$ of [Algorithm 1](#) satisfies $\hat{\pi} \in \arg \max_{t \in [T]} J_{\beta}(\pi^{(t)})$, the claimed bound on $J_{\beta}(\pi_{\beta}^*) - J_{\beta}(\hat{\pi})$ is immediate. Finally, observe that by definition of π_{β}^* ,

$$\begin{aligned} J_{\beta}(\pi_{\beta}^*) - J_{\beta}(\hat{\pi}) &= \mathbb{E}_{(x,y) \sim \pi_{\beta}^*} \left[r(x,y) - \beta \log \frac{\pi_{\beta}^*(y | x)}{\pi_{\text{base}}(y | x)} \right] - \mathbb{E}_{(x,y) \sim \hat{\pi}} \left[r(x,y) - \beta \log \frac{\hat{\pi}(y | x)}{\pi_{\text{base}}(y | x)} \right] \\ &= \mathbb{E}_{(x,y) \sim \pi_{\beta}^*} \left[r(x,y) - \beta \log \frac{\pi_{\beta}^*(y | x)}{\pi_{\text{base}}(y | x)} \right] - \mathbb{E}_{(x,y) \sim \hat{\pi}} \left[r(x,y) - \beta \log \frac{\pi_{\beta}^*(y | x)}{\pi_{\text{base}}(y | x)} \right] \\ &\quad + \mathbb{E}_{(x,y) \sim \hat{\pi}} \left[\beta \log \frac{\hat{\pi}(y | x)}{\pi_{\beta}^*(y | x)} \right] \\ &= \beta \log \mathbb{E}_{(x,y) \sim \pi_{\text{base}}} [\exp(r(x,y))] - \beta \log \mathbb{E}_{(x,y) \sim \pi_{\text{base}}} [\exp(r(x,y))] + \beta D_{\text{KL}}(\hat{\pi} \| \pi_{\beta}^*) \\ &= \beta D_{\text{KL}}(\hat{\pi} \| \pi_{\beta}^*). \end{aligned}$$

This completes the proof. \square

J.2.3 APPLYING XPO TO MAXIMUM-LIKELIHOOD SHARPENING

We now prove [Theorem J.2](#), the formal statement of [Theorem 4.3](#), which applies XPO to maximum-likelihood sharpening. This result is a straightforward corollary of [Theorem J.1](#) with the reward function $r_{\text{self}}(x, y) := \log \pi_{\text{base}}(y | x)$, together with the observation that low KL-regularized regret implies sharpness (under [Assumption 4.2](#)).

Theorem J.2 (Sharpening via active exploration). *There are absolute constants $c_{J.2}, C_{J.2} > 0$ so that the following holds. Let $\epsilon, \delta, \gamma_{\text{margin}}, \rho, \beta \in (0, 1)$ and $T \in \mathbb{N}$ be given. For base model π_{base} , define reward function $r(x, y) := \log \pi_{\text{base}}(y | x)$. Let $R_{\max} \geq 1 + \max_{x,y} \log \frac{1}{\pi_{\text{base}}(y | x)}$. Suppose that π_{base} satisfies [Assumption 4.2](#) with parameter γ_{margin} , that $\beta^{-1} \geq 2\gamma_{\text{margin}}^{-1} \log(2|\mathcal{Y}|/\delta)$, and that there is $\epsilon_{\text{disc}} \in (0, 1)$ so that*

$$T \geq C_{J.2} \frac{R_{\max}^2 \text{SEC}(\Pi) \log(2\mathcal{N}(\Pi, \epsilon_{\text{disc}})T/\rho)}{\epsilon^2 \delta^2 \beta^2}$$

and

$$\epsilon_{\text{disc}} \leq c_{J.2} \frac{\epsilon \delta}{\sqrt{\text{SEC}(\Pi)T}}$$

where $\text{SEC}(\Pi) := \text{SEC}(\Pi, r, T, \beta, R_{\max}^2; \pi_{\text{base}})$. Also suppose that $\pi_{\beta}^* \in \Pi$ where $\pi_{\beta}^*(y | x) \propto \pi_{\text{base}}^{1+\beta^{-1}}(y | x)$.

Then applying [Algorithm 1](#) with base model π_{base} , reward function r , iteration count T , regularization β , and optimism parameter $\alpha := \frac{\beta}{R_{\max}} \sqrt{\frac{\log(2\mathcal{N}(\Pi, \epsilon_{\text{disc}})T/\delta)}{\text{SEC}(\Pi)T}}$ yields a model $\hat{\pi} \in \Pi$ such that with probability at least $1 - \rho$,

$$\mathbb{P}_{x \sim \mu}[\hat{\pi}(y^*(x) | x) < 1 - \delta] \leq \epsilon.$$

The total sample complexity is

$$m = \tilde{O} \left(\frac{R_{\max}^2 \text{SEC}(\Pi) \log(\mathcal{N}(\Pi, \epsilon_{\text{disc}})/\rho) \log^2(|\mathcal{Y}|\delta^{-1})}{\gamma_{\text{margin}}^2 \epsilon^2 \delta^2} \right).$$

Proof of Theorem J.2. By definition of r , we have $|r(x, y)| \leq R_{\max}$ for all x, y . By assumption, [Assumption J.3](#) is satisfied, and by definition of R_{\max} , [Assumption 4.5](#) is satisfied with parameter

$V_{\max} := \beta R_{\max} \leq R_{\max}$. It follows from [Theorem J.1](#) that with probability at least $1 - \rho$, the output $\hat{\pi}$ of [Algorithm 1](#) satisfies

$$\beta D_{\text{KL}}(\hat{\pi} \parallel \pi_{\beta}^*) \lesssim (R_{\max} + V_{\max}) \sqrt{\frac{\text{SEC}(\Pi) \log(2\mathcal{N}(\Pi, \epsilon_{\text{disc}})T/\rho)}{T}} + \beta \epsilon_{\text{disc}} \sqrt{\text{SEC}(\Pi)T}.$$

By choice of T and ϵ_{disc} , so long as $C_{\text{J.2}} > 0$ is chosen to be a sufficiently large constant and $c_{\text{J.2}} > 0$ is chosen to be a sufficiently small constant, we have $\beta D_{\text{KL}}(\hat{\pi} \parallel \pi_{\beta}^*) \leq \frac{1}{12} \beta \epsilon \delta$, so by e.g. Equation (16) of [Sason & Verdú \(2016\)](#), $D_{\text{H}}^2(\hat{\pi}, \pi_{\beta}^*) \leq \epsilon \delta / (12)$.

For any $x \in \mathcal{X}$ and $y' \in \mathcal{Y} \setminus \mathbf{y}^*(x)$, by [Assumption 4.2](#) and definition of π_{β}^* we have

$$\begin{aligned} \frac{1}{\pi_{\beta}^*(y' | x)} &\geq \frac{\max_{y \in \mathcal{Y}} \pi_{\beta}^*(y | x)}{\pi_{\beta}^*(y' | x)} = \left(\frac{\max_{y \in \mathcal{Y}} \pi_{\text{base}}(y | x)}{\pi_{\text{base}}(y' | x)} \right)^{1+\beta^{-1}} \\ &\geq (1 + \gamma_{\text{margin}})^{1+\beta^{-1}} \geq e^{\gamma_{\text{margin}}/(2\beta)} \geq \frac{2|\mathcal{Y}|}{\delta} \end{aligned}$$

where the final inequality is by the assumption on β in the theorem statement. Therefore

$$\pi_{\beta}^*(\mathbf{y}^*(x) | x) \geq 1 - \sum_{y' \in \mathcal{Y} \setminus \mathbf{y}^*(x)} \pi_{\beta}^*(y' | x) \geq 1 - \frac{\delta}{2}.$$

Now for any x , we can lower bound

$$\begin{aligned} D_{\text{H}}^2(\hat{\pi}(\cdot | x), \pi_{\beta}^*(\cdot | x)) &\geq \left(\sqrt{1 - \hat{\pi}(\mathbf{y}^*(x) | x)} - \sqrt{1 - \pi_{\beta}^*(\mathbf{y}^*(x) | x)} \right)^2 \\ &\geq \frac{\delta}{12} \cdot \mathbb{I}\{\hat{\pi}(\mathbf{y}^*(x) | x) \leq 1 - \delta\}. \end{aligned}$$

Hence,

$$\begin{aligned} \mathbb{P}_{x \sim \mu}[\hat{\pi}(\mathbf{y}^*(x) | x) < 1 - \delta] &\leq \frac{12}{\delta} \mathbb{E}_{x \sim \mu} D_{\text{H}}^2(\hat{\pi}(\cdot | x), \pi_{\beta}^*(\cdot | x)) \\ &= \frac{12}{\delta} D_{\text{H}}^2(\hat{\pi}, \pi_{\beta}^*) \\ &\leq \epsilon. \end{aligned}$$

as claimed. \square

J.2.4 APPLICATION: LINEAR SOFTMAX MODELS

In this section we apply [Theorem 4.3](#) to the class of linear softmax models, proving [??](#). This demonstrates that [Algorithm 1](#) can achieve an exponential improvement in sample complexity compared to SFT-Sharpener.

Definition J.3 (Linear softmax model). *Let $d \in \mathbb{N}$ be given, and let $\phi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$ be a feature map with $\|\phi(x, y)\|_2 \leq 1$ for all x, y . Let $\pi_{\text{zero}} : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$ be the uniform model $\pi_{\text{zero}}(y | x) := \frac{1}{|\mathcal{Y}|}$, and let $B \geq 1$.¹⁶ We consider the linear softmax model class $\Pi_{\phi, B} := \{\pi_{\theta} : \theta \in \mathbb{R}^d, \|\theta\|_2 \leq B\}$ where $\pi_{\theta} : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$ is defined by*

$$\pi_{\theta}(y | x) \propto \pi_{\text{zero}}(y | x) \exp(\langle \phi(x, y), \theta \rangle).$$

Theorem J.3 (Restatement of [??](#)). *Let $\epsilon, \delta, \gamma_{\text{margin}}, \rho \in (0, 1)$ be given. Suppose that $\pi_{\text{base}} = \pi_{\theta^*} \in \Pi_{\phi, B}$ for some $\theta^* \in \mathbb{R}^d$ with $\|\theta^*\|_2 \leq \frac{\gamma_{\text{margin}} B}{3 \log(2|\mathcal{Y}|/\delta)}$. Also, suppose that π_{base} satisfies [Assumption 4.2](#) with parameter γ_{margin} . Then [Algorithm 1](#) with base model π_{base} , reward function $r(x, y) := \log \pi_{\text{base}}(x, y)$, regularization parameter $\beta := \gamma_{\text{margin}} / (2 \log(2|\mathcal{Y}|/\delta))$, and optimism parameter $\alpha(T) \propto \frac{\beta}{B + \log(|\mathcal{Y}|)} \sqrt{\frac{d \log(BdT/(\epsilon\delta)) + \log(T/\rho)}{dT \log(T)}}$ returns an (ϵ, δ) -sharpened model with probability at least $1 - \rho$, and has sample complexity*

$$m = \text{poly}(\epsilon^{-1}, \delta^{-1}, \gamma_{\text{margin}}^{-1}, d, B, \log(|\mathcal{Y}|/\rho)).$$

¹⁶We use the notation π_{zero} to highlight the fact that $\pi_{\text{zero}} = \pi_{\theta}$ for $\theta = 0$.

Before proving the result, we unpack the conditions. [Theorem J.3](#) requires the base model π_{base} to lie in the model class and also satisfy the margin condition ([Assumption 4.2](#)). For any constant $\epsilon, \delta > 0$, the sharpening algorithm then succeeds with sample complexity $\text{poly}(d, \gamma_{\text{margin}}^{-1}, B, \log(|\mathcal{Y}|))$. These conditions are non-vacuous; in fact, there are fairly natural examples for which non-exploratory algorithm such as SFT-Sharpener require sample complexity $\exp(\Omega(d))$, whereas all of the above parameters are $\text{poly}(d)$. The following is one such example.

Example J.1 (Separation between RLHF-Sharpener and SFT-Sharpener). Set $\mathcal{X} = \{x\}$ and let $\mathcal{Y} \subset \mathbb{R}^d$ be a $1/4$ -packing of the unit sphere in \mathbb{R}^d of cardinality $\exp(\Theta(d))$. Define $\phi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$ by $\phi(x, y) := y$, and let $B = Cd \log d$ for an absolute constant $C > 0$. Fix any $y^* \in \mathcal{Y}$ and define $\pi_{\text{base}} := \pi_{\theta^*} \in \Pi_{\phi, B}$ by $\theta^* := y^*$. Then for any $y \neq y^*$, we have $\langle y, y^* \rangle \leq 1 - \Omega(1)$, so

$$\frac{\pi_{\text{base}}(y^* | x)}{\pi_{\text{base}}(y | x)} = \exp(\langle y^* - y, y^* \rangle) = \exp(\Omega(1)) = 1 + \Omega(1).$$

Thus, π_{base} satisfies [Assumption 4.2](#) with $\gamma_{\text{margin}} = \Omega(1)$. Moreover, $\|\theta^*\|_2 = 1 \leq \frac{\gamma_{\text{margin}} B}{3 \log(2|\mathcal{Y}|/\delta)}$ for any $\delta = 1/\text{poly}(d)$, so long as C is a sufficiently large constant. It follows from ?? that [Algorithm 1](#) computes an (ϵ, δ) -sharpened model with sample complexity $\text{poly}(\epsilon^{-1}, \delta^{-1}, d)$. However, since $\pi_{\text{base}}(y^* | x) \leq \pi_{\text{base}}(y | x) \cdot \exp(2)$ for all $y \in \mathcal{Y}$, it is clear that

$$C_{\text{cov}} = \mathbb{E} \left[\frac{1}{\pi_{\text{base}}(y^*(x) | x)} \right] = \frac{1}{\pi_{\text{base}}(y^* | x)} = \Omega(|\mathcal{Y}|) = \exp(\Omega(d)).$$

Thus, the sample complexity guarantee for SFT-Sharpener in [Theorem 4.1](#) will incur *exponential* dependence on d in the sample complexity. It is straightforward to check that this dependence is real for SFT-Sharpener, and not just an artifact of the analysis, since the model that SFT-Sharpener is trying to learn (via MLE) will itself not be sharp in this example, unless $\exp(\Omega(d))$ samples are drawn per prompt. \triangleleft

We now proceed to the proof of [Theorem J.3](#), which requires the following bounds on the covering number and the Sequential Extrapolation Coefficient of $\Pi_{\phi, B}$.

Lemma J.4. *Let $\epsilon_{\text{disc}} > 0$. Then $\Pi_{\phi, B}$ has an ϵ_{disc} -net of size $(6B/\epsilon_{\text{disc}})^d$.*

Proof of Lemma J.4. By a standard packing argument, there is a set $\{\theta_1, \dots, \theta_N\}$ of size $(6B/\epsilon_{\text{disc}})^d$ such that for every $\theta \in \mathbb{R}^d$ with $\|\theta\|_2 \leq B$ there is some $i \in [N]$ with $\|\theta_i - \theta\|_2 \leq \epsilon_{\text{disc}}/2$. Now for any $x \in \mathcal{X}$ and $y \in \mathcal{Y}$,

$$\begin{aligned} \log \frac{\pi_{\theta}(y | x)}{\pi_{\theta_i}(y | x)} &= \log \frac{\exp(\langle \phi(x, y), \theta \rangle)}{\exp(\langle \phi(x, y), \theta_i \rangle)} + \log \frac{\mathbb{E}_{(x', y') \sim \pi_{\text{zero}}} \exp(\langle \phi(x', y'), \theta_i \rangle)}{\mathbb{E}_{(x', y') \sim \pi_{\text{zero}}} \exp(\langle \phi(x', y'), \theta \rangle)} \\ &= \langle \phi(x, y), \theta - \theta_i \rangle + \log \frac{\mathbb{E}_{(x', y') \sim \pi_{\text{zero}}} [\exp(\langle \phi(x', y'), \theta \rangle) \exp(\langle \phi(x', y'), \theta_i - \theta \rangle)]}{\mathbb{E}_{(x', y') \sim \pi_{\text{zero}}} \exp(\langle \phi(x', y'), \theta \rangle)}. \end{aligned}$$

The first term is bounded by $\epsilon_{\text{disc}}/2$ in magnitude. In the second term, we have $\exp(\langle \phi(x', y'), \theta_i - \theta \rangle) \in [\exp(-\epsilon_{\text{disc}}/2), \exp(\epsilon_{\text{disc}}/2)]$, so the ratio of expectations lies in $[\exp(-\epsilon_{\text{disc}}/2), \exp(\epsilon_{\text{disc}}/2)]$ as well, and so the log-ratio lies in $[-\epsilon_{\text{disc}}/2, \epsilon_{\text{disc}}/2]$. In all, we get $\left| \log \frac{\pi_{\theta}(y | x)}{\pi_{\theta_i}(y | x)} \right| \leq \epsilon_{\text{disc}}$. Thus, $\{\pi_{\theta_1}, \dots, \pi_{\theta_N}\}$ is an ϵ_{disc} -net for Π . \square

Lemma J.5. *Let $r : \mathcal{X} \times \mathcal{Y} \rightarrow [-R_{\text{max}}, R_{\text{max}}]$ be a reward function and let $T \in \mathbb{N}$ and $\beta > 0$. If $\lambda \geq 4\beta^2 B^2 + R_{\text{max}}^2$ then for any $\pi^* \in \Pi_{\phi, B}$,*

$$\text{SEC}(\Pi_{\phi, B}, r, T, \beta, \lambda; \pi^*) \lesssim d \log(T + 1).$$

Proof of Lemma J.5. Fix $\pi^{(1)}, \dots, \pi^{(T)} \in \Pi_{\phi, B}$. By definition, there are some $\theta^{(1)}, \dots, \theta^{(T)} \in \mathbb{R}^d$ with $\|\theta^{(t)}\|_2 \leq B$ and

$$\pi^{(t)}(y | x) \propto \pi_{\text{zero}}(y | x) \exp(\langle \phi(x, y), \theta^{(t)} \rangle)$$

for all $t \in [T]$ and $(x, y) \in \mathcal{X} \times \mathcal{Y}$. Similarly, there is some $\theta^* \in \mathbb{R}^d$ with $\|\theta^*\|_2 \leq B$ and $\pi^*(y | x) \propto \pi_{\text{zero}}(y | x) \exp(\langle \phi(x, y), \theta^* \rangle)$.

Define $\tilde{\phi} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^{d+1}$ by $\tilde{\phi}(x, y) := [\phi(x, y), \frac{r(x, y)}{R_{\max}}]$ and define $\tilde{\theta}^{(t)} := [\beta(\theta^{(t)} - \theta^*), -R_{\max}]$. Then for any $t \in [T]$ we have

$$\begin{aligned} & \frac{\mathbb{E}^{(t)} \left[\beta \log \frac{\pi^{(t)}(y|x)}{\pi^*(y|x)} - r(x, y) - \beta \log \frac{\pi^{(t)}(y'|x)}{\pi^*(y'|x)} + r(x, y') \right]^2}{\lambda \vee \sum_{i=1}^{t-1} \mathbb{E}^{(i)} \left[\left(\beta \log \frac{\pi^{(i)}(y|x)}{\pi^*(y|x)} - r(x, y) - \beta \log \frac{\pi^{(i)}(y'|x)}{\pi^*(y'|x)} + r(x, y') \right)^2 \right]} \\ &= \frac{\mathbb{E}^{(t)} \left[\langle \tilde{\phi}(x, y) - \tilde{\phi}(x, y'), \tilde{\theta}^{(t)} \rangle \right]^2}{\lambda \vee \sum_{i=1}^{t-1} \mathbb{E}^{(i)} \left[\left(\langle \tilde{\phi}(x, y) - \tilde{\phi}(x, y'), \tilde{\theta}^{(i)} \rangle \right)^2 \right]} \\ &\leq \frac{(\tilde{\theta}^{(t)})^\top \Sigma^{(t)} \tilde{\theta}^{(t)}}{\lambda \vee \sum_{i=1}^{t-1} (\tilde{\theta}^{(i)})^\top \Sigma^{(i)} \tilde{\theta}^{(i)}} \end{aligned}$$

where for each $i \in [T]$ we have defined $\Sigma^{(i)} := \mathbb{E}^{(i)} \left[(\tilde{\phi}(x, y) - \tilde{\phi}(x, y'))(\tilde{\phi}(x, y) - \tilde{\phi}(x, y'))^\top \right]$.

Observe that $\|\tilde{\theta}^{(t)}\|_2^2 \leq 4\beta^2 B^2 + R_{\max}^2 \leq \lambda$ by assumption on λ . Therefore,

$$\begin{aligned} \frac{(\tilde{\theta}^{(t)})^\top \Sigma^{(t)} \tilde{\theta}^{(t)}}{\lambda \vee \sum_{i=1}^{t-1} (\tilde{\theta}^{(i)})^\top \Sigma^{(i)} \tilde{\theta}^{(i)}} &\lesssim \frac{(\tilde{\theta}^{(t)})^\top \Sigma^{(t)} \tilde{\theta}^{(t)}}{\lambda + \sum_{i=1}^{t-1} (\tilde{\theta}^{(i)})^\top \Sigma^{(i)} \tilde{\theta}^{(i)}} \\ &\leq \frac{(\tilde{\theta}^{(t)})^\top \Sigma^{(t)} \tilde{\theta}^{(t)}}{(\tilde{\theta}^{(t)})^\top \left(I_d + \sum_{i=1}^{t-1} \Sigma^{(i)} \right) \tilde{\theta}^{(t)}} \\ &\leq \lambda_{\max} \left(\left(I_d + \sum_{i=1}^{t-1} \Sigma^{(i)} \right)^{-1/2} \Sigma^{(t)} \left(I_d + \sum_{i=1}^{t-1} \Sigma^{(i)} \right)^{-1/2} \right) \\ &\leq \text{Tr} \left(\left(I_d + \sum_{i=1}^{t-1} \Sigma^{(i)} \right)^{-1/2} \Sigma^{(t)} \left(I_d + \sum_{i=1}^{t-1} \Sigma^{(i)} \right)^{-1/2} \right) \\ &= \text{Tr} \left(\left(I_d + \sum_{i=1}^{t-1} \Sigma^{(i)} \right)^{-1} \Sigma^{(t)} \right). \end{aligned}$$

Observe that $\text{Tr}(\Sigma^{(t)}) \leq \max_{x, y} \|\tilde{\phi}(x, y)\|_2^2 \lesssim 1$. Hence by [Lemma F.2](#), we have

$$\begin{aligned} & \sum_{t=1}^T \frac{\mathbb{E}^{(t)} \left[\beta \log \frac{\pi^{(t)}(y|x)}{\pi^*(y|x)} - r(x, y) - \beta \log \frac{\pi^{(t)}(y'|x)}{\pi^*(y'|x)} + r(x, y') \right]^2}{\lambda \vee \sum_{i=1}^{t-1} \mathbb{E}^{(i)} \left[\left(\beta \log \frac{\pi^{(i)}(y|x)}{\pi^*(y|x)} - r(x, y) - \beta \log \frac{\pi^{(i)}(y'|x)}{\pi^*(y'|x)} + r(x, y') \right)^2 \right]} \\ &\lesssim \sum_{t=1}^T \text{Tr} \left(\left(I_d + \sum_{i=1}^{t-1} \Sigma^{(i)} \right)^{-1} \Sigma^{(t)} \right) \\ &\lesssim d \log(T+1). \end{aligned}$$

Since $\pi^{(1)}, \dots, \pi^{(T)} \in \Pi$ were arbitrary, this completes the proof. \square

The proof is now immediate from [Theorem J.2](#) and the above lemmas.

Proof of Theorem J.3. By the assumption on θ^* and choice of β , the model π_β^* defined by $\pi_\beta^*(y | x) \propto \pi_{\text{base}}(y | x)^{1+\beta^{-1}}$ satisfies $\pi_\beta^* = \pi_{(1+\beta^{-1})\theta^*} \in \Pi_{\phi, B}$. By [Lemma J.4](#), we have $\mathcal{N}(\Pi_{\phi, B}, \epsilon_{\text{disc}}) \leq (6B/\epsilon_{\text{disc}})^d$. Take $R_{\max} := \sqrt{4\beta^2 B^2 + (2B + \log |\mathcal{Y}|)^2}$. We know that $r(x, y) := \log \pi_{\text{base}}(y | x)$ satisfies $|r(x, y)| \leq 2B + \log |\mathcal{Y}|$ for all x, y . By [Lemma J.5](#), we therefore get that $\text{SEC}(\Pi_{\phi, B}, r, T, \beta, R_{\max}^2; \pi_{\text{base}}) \lesssim d \log(T+1)$. Substituting these bounds into [Theorem J.2](#) yields the claimed result. \square