

Annealed Softmax Greedy in Many-Armed Bayesian Bandits

William Overman
Stanford University

Mohsen Bayati
Stanford University

Abstract

Reinforcement learning with verifiable rewards (RLVR) and group-based policy optimization methods such as GRPO update a stochastic policy by sampling multiple completions per prompt and increasing the policy’s probability on those with higher reward, regularized by a KL penalty toward a reference policy. These updates do not include explicit mechanisms that track epistemic uncertainty. This paper studies a stylized explanation for why such uncertainty-agnostic updates can nevertheless be effective. We analyze an annealed softmax (Boltzmann) policy that selects actions according to a softmax of empirical mean rewards in a many-armed Bayesian Bernoulli bandit. Under a linear upper-tail condition on the prior (the $\beta = 1$ case of β -regularity), which implies an abundance of near-optimal arms, we prove that annealed softmax greedy achieves Bayes regret $\tilde{O}(m + T/m)$, and in particular $\tilde{O}(\sqrt{T})$ when the number of arms scales as $m = \Theta(\sqrt{T})$. This is the near-optimal Bayes regret rate in this regime, attained also by empirical-mean greedy. Under β -regularity, many arms maintain empirical means close to the optimum throughout learning, so when softmax samples an arm other than the empirically best, that arm tends to be another near-optimal one rather than a clearly inferior one. By contrast, with a small number of arms, the same kind of softmax policy can suffer linear regret (Cesa-Bianchi et al., 2017). The result also provides a structural analogy to RLVR, where a base policy with a non-negligible probability of producing a correct completion plays the role of β -regularity.

1 Introduction

Reinforcement learning with verifiable rewards (RLVR) is now a common component of language-model post-training pipelines. In RLVR, a model is treated as a stochastic policy over candidate solutions and is trained using rewards that can be checked automatically, such as exact-match accuracy in mathematics or unit-test success in code. This makes it possible to optimize behavior at scale without relying on human preference labels. Group-based policy optimization methods such as Group Relative Policy Optimization (GRPO) are prominent examples of this paradigm (Shao et al., 2024).

RLVR sits somewhat outside the classical exploration picture in reinforcement learning and bandit theory. In those settings, performance is often tied to *epistemic uncertainty*: the learner must gather information to distinguish promising actions from poor ones. Group-based RLVR pipelines, by contrast, sample multiple completions per prompt, raise the policy probability of those with higher reward, and apply a KL penalty toward a reference policy; the randomness comes from repeated sampling rather than from an uncertainty-aware exploration rule. Under verifiable rewards, such updates can be viewed as iteratively amplifying the probability of already successful outputs, with the gain coming largely from redistribution within the model’s existing support (Liu

Emails: wpo@stanford.edu, bayati@stanford.edu.

et al., 2025b). This raises a basic question: when can uncertainty-agnostic reweighting work well in problems that seem to require exploration?

A related discussion concerns the base model’s *coverage* of good solutions. Here $\text{pass}@k$ denotes the probability that at least one of k independent samples from the model is correct, with $\text{pass}@1$ being the per-sample success rate. One line of work reports that RLVR often improves $\text{pass}@1$ while making little progress on $\text{pass}@k$ for large k , suggesting that rewarded reasoning paths may already be present in the base model’s output distribution (Yue et al., 2025). Other recent work reports settings where RLVR *does* extend reasoning performance, indicating that the empirical picture depends on training details and evaluation choices (Cui et al., 2025). We do not try to resolve that broader debate. We instead study a simplified model that isolates one mechanism by which reweighting alone can be effective.

1.1 A stylized lens: many-armed bandits and annealed softmax

We consider a stochastic many-armed bandit model with Bernoulli rewards and Beta priors, where each arm represents a “completion mode” and reward corresponds to a verifiable success event. Within this model, we analyze an annealed softmax policy that, at each time step t , selects an arm with probability proportional to $\exp(\eta_t m_{i,t})$, where $m_{i,t}$ is the empirical mean reward of arm i at time t and η_t is an inverse-temperature schedule that increases over time, a setup referred to as *annealing*. The policy does not use an optimism term, a posterior sample, or per-arm confidence intervals.

A classical result of Cesa-Bianchi et al. (2017) shows that this type of exploration (also referred to as Boltzmann exploration) can suffer linear regret with a small number of arms under any monotone temperature schedule, and that strong guarantees in that regime typically require schedules that adapt to arm-specific information. The standard failure mode is “cooling too fast”: early noise gets amplified, causing premature concentration on suboptimal actions. At first glance, such results argue against exactly the kind of uncertainty-agnostic annealing that appears in RLVR practice.

The many-armed Bayesian setting allows a different possibility. If the prior places enough mass near the optimum, the action space contains many arms whose rewards are already close to optimal. Even if the algorithm spreads probability across several arms, many of them are good enough that the regret cost remains small. This idea is studied in the many-armed bandit literature, where greedy-style procedures achieve Bayesian regret guarantees under upper-tail regularity conditions on the prior (Bayati et al., 2020).

1.2 Main results and contributions

The question we study is therefore: can annealed softmax, despite ignoring epistemic uncertainty, inherit the regret behavior of greedy methods in the many-armed Bayesian regime? We answer this positively. The main contributions are:

1. *Bayes regret guarantee.* Under the linear upper-tail condition on the prior (the $\beta = 1$ case of β -regularity), annealed softmax achieves Bayes regret $\tilde{O}(m + T/m)$, yielding $\tilde{O}(\sqrt{T})$ when $m = \Theta(\sqrt{T})$. This is the near-optimal Bayes regret rate in this regime (also attained by empirical-mean greedy (Bayati et al., 2020)), achieved using only empirical-mean scores, with no optimism term, no posterior sample, and no per-arm confidence interval. The analysis extends to general $\beta > 0$, with optimized rate $\tilde{O}(m + T(\log T/m)^{1/\beta})$ (Remark 5.3).

2. *Mechanism.* Annealed softmax differs from greedy only by placing some probability mass on arms that are not currently empirically best. Under β -regularity, many arms maintain empirical means close to the optimum throughout learning, so when softmax samples an arm other than the empirically best, that arm tends to be another near-optimal one rather than a clearly inferior one. This contrasts with the failure mode studied by Cesa-Bianchi et al. (2017), where the same kind of policy can suffer linear regret with a small number of arms.
3. *Structural analogue in RLVR.* The β -regular prior condition has a direct analogue in RLVR: a base policy with non-negligible probability of producing a correct completion plays the role of β -regularity, with pass@ k probing the upper tail of the policy’s completion distribution. In that setting, repeated sampling surfaces correct completions reliably, and subsequent reweighting can improve pass@1 without an explicit uncertainty-aware exploration rule. We note that this is only a structural analogy and not a formal result about GRPO; the theorem lives in the non-contextual bandit model, and extending the analysis to the contextual or sequential settings where GRPO operates is an open problem (Shao et al., 2024; Liu et al., 2025b; Yue et al., 2025; Cui et al., 2025).

1.3 Paper organization

Section 2 reviews related work. Section 3 introduces the many-armed Bayesian Bernoulli bandit model, the β -regular prior assumption, and notation. Section 4 presents the greedy baseline and the annealed softmax greedy (ASG) algorithm. Section 5 states the main regret theorem and its extensions. Section 6 complements the theoretical analysis with simulations on Bernoulli bandits. Section 7 concludes with discussion. Appendix A contains the proofs.

2 Related Work

Our paper connects three strands of literature: the theory of Boltzmann exploration in bandits, the many-armed bandit framework with “free exploration,” and the empirical RLVR/GRPO pipeline for LLM post-training. We give a concise overview here and defer a comprehensive survey to Appendix B.

Boltzmann exploration and its limitations. Boltzmann (softmax) action selection is a standard randomized alternative to greedy or optimistic policies in bandits and RL. Cesa-Bianchi et al. (2017) prove that, for fixed- K stochastic bandits, *any monotone* temperature schedule can be forced into suboptimal behavior—either exploring too long or committing too early—and propose per-arm learning rates that explicitly track uncertainty as a remedy. This negative result is the starting point of our analysis: we identify a structural regime (many arms, thick-tailed prior) that circumvents the impossibility without requiring uncertainty-aware schedules.

Many-armed bandits and free exploration. When the number of arms is large relative to the horizon, the distribution of arm qualities (rather than per-arm estimation) governs achievable regret. In the Bayesian setting, Bayati et al. (2020) show that under an upper-tail regularity condition on the prior (“ β -regularity”), a subsampled greedy policy achieves Bayes regret $\tilde{O}(\max\{m, T/m\})$, yielding $\tilde{O}(\sqrt{T})$ with $m = \Theta(\sqrt{T})$ arms. The key mechanism is *free exploration*: discarding a poorly performing arm still leaves many near-optimal alternatives. Our work extends this viewpoint from greedy to annealed softmax policies, showing that the same prior-tail structure suppresses the

“softmax leakage” that causes failures in the fixed- K regime. Related infinite-armed formulations include [Berry et al. \(1997\)](#); [Wang et al. \(2008\)](#); [Carpentier and Valko \(2015\)](#).

RLVR, GRPO, and the “bounded-by-base” phenomenon. Reinforcement learning with verifiable rewards has become a core post-training technique for reasoning-oriented LLMs. GRPO [Shao et al. \(2024\)](#) performs group-sampled policy updates under KL regularization, and recent analyses relate its dynamics to iterative soft reweighting of completions [Mroueh \(2025\)](#); [Liu et al. \(2025b\)](#). A central empirical observation is that RLVR often improves pass@1 while failing to improve, or even degrading, large- k pass@ k , suggesting that optimization primarily redistributes probability mass within the base model’s existing support rather than discovering new reasoning paths [Yue et al. \(2025\)](#); though subsequent work argues this picture is nuanced and depends on training details and evaluation choices [Cui et al. \(2025\)](#); [Wen et al. \(2025\)](#); [Liu et al. \(2025a\)](#). Our bandit model provides a stylized formalization of this “bounded-by-base” effect: the β -regular tail condition translates “good pass@ k ” into “many near-optimal arms,” under which softmax reweighting suffices without explicit exploration.

3 Problem Setup

This section sets up the model and the notation. More specifically, it formalizes the bandit model, defines the prior regularity condition and the Bayes regret criterion, and introduces the empirical-mean notation used in the algorithm. Throughout the paper we assume $T \geq m$.

3.1 Many-armed Bayesian Bernoulli bandit

Fix a time horizon $T \in \mathbb{N}$ and consider $m \in \mathbb{N}$ arms indexed by $i \in [m] := \{1, \dots, m\}$. Each arm has an unknown mean $\mu_i \in [0, 1]$ drawn i.i.d. from a prior Γ on $[0, 1]$. Conditional on μ_i , each arm i has an i.i.d. sequence of Bernoulli rewards

$$X_{i,s} \mid \mu_i \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(\mu_i) \quad (s = 1, 2, \dots),$$

where $X_{i,s}$ is the reward from the s -th pull of arm i (s is arm-specific, not the global time index). At each round $t \in [T]$, a policy selects an arm $A_t \in [m]$ (possibly randomized and history-dependent) and observes $X_{A_t, N_{A_t}(t)}$, where

$$N_i(t) := \sum_{s=1}^t \mathbf{1}\{A_s = i\}$$

is the number of pulls of arm i up to time t .

3.2 Prior regularity: β -regular upper tail

We assume Γ has a polynomially thick upper tail near 1.

Definition 3.1 (β -regular prior near 1). A distribution Γ on $[0, 1]$ is β -regular if there exist constants $0 < c_0 \leq C_0 < \infty$ and $\varepsilon_0 \in (0, 1)$ such that, for all $\varepsilon \in (0, \varepsilon_0]$,

$$c_0 \varepsilon^\beta \leq \Gamma([1 - \varepsilon, 1]) \leq C_0 \varepsilon^\beta.$$

We will present the main theorem first for the *linear tail* case $\beta = 1$ (which is the regime treated in [Bayati et al. \(2020\)](#) for simplicity), and then state a clean general- β extension as a corollary/remark.

3.3 Bayes regret and a convenient surrogate

Let $\mu_* := \max_{i \in [m]} \mu_i$. The (frequentist) regret of a policy π given $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m)$ is

$$R_T(\pi \mid \boldsymbol{\mu}) := \sum_{t=1}^T (\mu_* - \mu_{A_t}).$$

The *Bayes regret* is

$$\text{BR}_{T,m}(\pi) := \mathbb{E}[R_T(\pi \mid \boldsymbol{\mu})],$$

where the expectation is over $\boldsymbol{\mu} \sim \Gamma^{\otimes m}$ (i.e., μ_1, \dots, μ_m are i.i.d. draws from Γ , so $\Gamma^{\otimes m}$ is the m -fold product measure on $[0, 1]^m$), reward randomness, and any internal randomness of π .

We will also use the standard surrogate “regret-to-1”

$$\tilde{R}_T(\pi \mid \boldsymbol{\mu}) := \sum_{t=1}^T (1 - \mu_{A_t}) = \sum_{i=1}^m (1 - \mu_i) N_i(T), \quad (1)$$

which satisfies the exact identity

$$R_T(\pi \mid \boldsymbol{\mu}) = \tilde{R}_T(\pi \mid \boldsymbol{\mu}) - T(1 - \mu_*). \quad (2)$$

Thus

$$\text{BR}_{T,m}(\pi) = \mathbb{E}[\tilde{R}_T(\pi \mid \boldsymbol{\mu})] - T \mathbb{E}[1 - \mu_*] \leq \mathbb{E}[\tilde{R}_T(\pi \mid \boldsymbol{\mu})]. \quad (3)$$

Under 1-regularity, $\mathbb{E}[1 - \mu_*]$ is of order $1/m$ (Lemma A.9 below), so the exact identity can sharpen constants, but our main upper bound only needs $\text{BR}_{T,m}(\pi) \leq \mathbb{E}[\tilde{R}_T(\pi \mid \boldsymbol{\mu})]$.

3.4 Empirical means

Let $S_i(t)$ be the number of observed successes of arm i up to time t , so $0 \leq S_i(t) \leq N_i(t)$. Because every arm is pulled once during initialization in Algorithms 1–2, we have $N_i(t) \geq 1$ throughout the greedy/softmax phase. We therefore use the empirical mean

$$\hat{\mu}_{i,t} := \frac{S_i(t)}{N_i(t)} \quad (4)$$

as the score of arm i at time t .

Asymptotic notation. Throughout the paper, $\tilde{O}(\cdot)$ hides polylogarithmic factors in T and m .

4 Algorithms

This section presents the greedy baseline and the annealed softmax greedy (ASG) policy that is the main object of analysis.

4.1 Greedy baseline on empirical means

We begin with the natural greedy benchmark, shown in Algorithm 1. After pulling each arm once so that every arm has an initial estimate, the policy always selects an arm with the largest empirical mean. This is exactly the many-armed greedy rule analyzed by Bayati et al. (2020), and it serves as the baseline throughout the paper.

4.2 Annealed Softmax Greedy on empirical means

Our main object of study is a randomized analogue of greedy, given in Algorithm 2. Instead of deterministically choosing the empirical-best arm at each round, the policy samples from a softmax distribution over empirical means. Arms with higher empirical means are more likely to be selected, but lower-ranked arms still receive some probability mass. The amount of randomness is controlled by a nondecreasing inverse-temperature schedule $\{\eta_t\}_{t \geq 1}$.

Algorithm 1 Greedy (Empirical-Mean Greedy)

Require: Arms $[m]$, horizon T

- 1: **Initialization:**
 - 2: **for** $t = 1, \dots, m$ **do**
 - 3: Pull arm $A_t = t$; observe reward $X_{t,1}$
 - 4: **end for**
 - 5: **Greedy phase:**
 - 6: **for** $t = m + 1, \dots, T$ **do**
 - 7: Compute empirical means $\hat{\mu}_{i,t-1} = \frac{S_i(t-1)}{N_i(t-1)}$ for all $i \in [m]$
 - 8: Pull arm $A_t \in \arg \max_{i \in [m]} \hat{\mu}_{i,t-1}$ (ties broken arbitrarily)
 - 9: Observe reward $X_{A_t, N_{A_t}(t)}$
 - 10: **end for**
-

This policy can be viewed as interpolating between uniform sampling and pure greedy behavior. When η_t is small, the softmax distribution is relatively flat, so the algorithm spreads probability more evenly across all arms. As η_t grows, the distribution becomes more concentrated, and the policy increasingly resembles Greedy. The key question for the rest of the paper is whether this softmax sampling rule, which is agnostic to epistemic uncertainty in arm quality, can achieve near-greedy regret in the many-armed Bayesian regime.

ASG uses exactly the same empirical information as the greedy baseline; it neither estimates uncertainty explicitly nor uses per-arm temperatures that depend on the pull counts $N_i(t)$. [Cesa-Bianchi et al. \(2017\)](#) prove that any monotone temperature schedule shared across arms can suffer linear regret in some small-arm instance, and propose per-arm scalings as a remedy: their Boltzmann–Gumbel Exploration uses $\sigma/\sqrt{N_i(t)}$, so less-pulled arms have noisier scores and are favored for exploration. We study the unadjusted version because it matches the structure used in group-based methods such as GRPO, where the policy does not adapt its per-completion sampling temperature to historical counts.

Concrete schedules. For the analysis, it is convenient to choose a logarithmically increasing schedule of the form

$$\eta_t = \frac{c_\eta}{\delta} \log(t \vee 2), \tag{5}$$

where $c_\eta > 1$ is a fixed constant and δ is a threshold parameter that will be tuned as a function of (T, m) in the regret bound. Intuitively, δ sets the scale of what counts as a meaningfully suboptimal arm in the proof. With this choice, $\exp(-\eta_t \delta) = (t \vee 2)^{-c_\eta}$, so the probability weight placed on arms that are worse by at least δ decays polynomially in time. In particular,

$$\sum_{t=1}^T \exp(-\eta_t \delta) = O(1),$$

Algorithm 2 Annealed Softmax Greedy (ASG)

Require: Arms $[m]$, horizon T , nonnegative nondecreasing inverse-temperature schedule $\{\eta_t\}_{t \geq 1}$ with $\eta_t \rightarrow \infty$

1: **Initialization:**

2: **for** $t = 1, \dots, m$ **do**

3: Pull arm $A_t = t$; observe reward $X_{t,1}$

4: **end for**

5: **Softmax phase:**

6: **for** $t = m + 1, \dots, T$ **do**

7: Compute empirical means $\hat{\mu}_{i,t-1}$ for all $i \in [m]$

8: Sample arm $A_t = i$ with probability

$$p_t(i) := \frac{\exp(\eta_t \hat{\mu}_{i,t-1})}{\sum_{j=1}^m \exp(\eta_t \hat{\mu}_{j,t-1})} \quad (6)$$

9: Observe reward $X_{A_t, N_{A_t}(t)}$

10: **end for**

which is exactly what we need to make the total softmax “leakage” summable. In the linear-tail case ($\beta = 1$), for example, the relevant choice is $\delta \asymp (\log T)/m$, leading to the $\tilde{O}(m + T/m)$ Bayes regret rate (Theorem 5.2).

5 Main Results

This section presents the main Bayes regret bound for ASG.

Before stating our main theorem, we recall the result for the greedy policy (Algorithm 1). In the Bernoulli many-armed setting with a linear upper tail, Bayati et al. (2020) show that empirical-mean greedy achieves near-optimal Bayes regret.

Proposition 5.1 (Greedy benchmark in the Bernoulli, $\beta = 1$ regime (Bayati et al., 2020, Theorem (Bernoulli))). *Assume Bernoulli rewards and a 1-regular prior Γ (Definition 3.1 with $\beta = 1$). Then Greedy (Algorithm 1) achieves*

$$\text{BR}_{T,m}(\text{Greedy}) \leq \tilde{O}\left(m + \frac{T}{m}\right),$$

and in particular for $m = \Theta(\sqrt{T})$ one has $\text{BR}_{T,m}(\text{Greedy}) = \tilde{O}(\sqrt{T})$.

5.1 ASG matches greedy up to a leakage term

We now turn to our main result for ASG. At a high level, the theorem shows that ASG behaves like greedy plus an additional *softmax leakage* penalty: the probability mass assigned away from the empirically best arms. The key technical point is that in the many-armed, thick-tail regime, this leakage remains controlled because when softmax samples an arm other than the empirically best, that arm tends to be another near-optimal one rather than a clearly inferior one. The extra randomization introduced by softmax therefore does not change the leading Bayes regret scaling.

The bound below decomposes regret into several interpretable terms. The terms m and $T\delta$ are the same coarse baseline contributions that already appear in greedy-style many-armed analyses.

The logarithmic term $m(1 + \log(1/\delta))$ reflects the cost of integrating over the upper tail. The summation term $\frac{1}{\delta} \sum_{t=1}^T \exp(-\eta_t \delta)$ is the softmax-specific leakage term, controlled by the cooling schedule. Finally, the exponentially small term $T \exp(-cm\delta)$ corresponds to the bad event that too few of the m arms are near-optimal. In this sense, ASG is essentially greedy up to a leakage penalty that becomes small under suitable annealing.

Theorem 5.2 (ASG Bayes regret, Bernoulli, linear tail). *Assume Bernoulli rewards, and assume the prior Γ is 1-regular (Definition 3.1 with $\beta = 1$) with regularity constants $(c_0, C_0, \varepsilon_0)$. Fix any $\delta \in (0, \delta_0]$ with $\delta_0 \leq \min\{1/8, \varepsilon_0/8\}$ (which also ensures $\delta < 1/6$, the condition required by Lemma A.2). Run ASG (Algorithm 2) with any nonnegative nondecreasing schedule $\{\eta_t\}_{t \geq 1}$.*

Then there exist constants $c, C > 0$ (depending only on c_0, C_0, ε_0 and on the Bernoulli crossing constant in Lemma A.2) such that

$$\text{BR}_{T,m}(\text{ASG}) \leq C \left[m + T\delta + m(1 + \log(1/\delta)) + \frac{1}{\delta} \sum_{t=1}^T \exp(-\eta_t \delta) + T \exp(-cm\delta) \right]. \quad (7)$$

In particular, choose

$$\delta = \min \left\{ \delta_0, A \frac{\log(T \vee 2)}{m} \right\}, \quad \eta_t = \frac{c_\eta}{\delta} \log(t \vee 2) \quad (c_\eta > 1), \quad (8)$$

with $A > 1/c$. Then

$$\text{BR}_{T,m}(\text{ASG}) = \tilde{O} \left(m + \frac{T}{m} \right),$$

and for $m = \Theta(\sqrt{T})$ one obtains $\text{BR}_{T,m}(\text{ASG}) = \tilde{O}(\sqrt{T})$.

The proof is given in Section A. The main takeaway is that ASG inherits the same leading-order regret scaling as greedy in the linear-tail many-armed regime. In this sense, the additional softness of the policy preserves the near-greedy rate; it only introduces a leakage term that can be made summable by a suitable logarithmic cooling schedule.

Remark 5.3 (General β (statement-level)). The linear-tail case $\beta = 1$ is the cleanest to state, but the same proof strategy extends to general $\beta > 0$, paralleling the corresponding generalization of the greedy benchmark in Bayati et al. (2020). Repeating the proof with Lemma A.3 replaced by the general- β analogue $p_\delta \geq c\delta^\beta$ (which follows from the same crossing argument together with the β -regularity bound $\Gamma([1-\delta, 1]) \geq c_0\delta^\beta$), and with the corresponding many-good-arms bound $\mathbb{P}(M(\delta) < r) \leq \exp(-cm\delta^\beta)$, yields a bound of the form

$$\text{BR}_{T,m}(\text{ASG}) \leq \tilde{O} \left(m + T\delta + m \cdot \mathfrak{M}_\beta(\delta) + \frac{1}{\delta^\beta} \sum_{t=1}^T e^{-\eta_t \delta} + T e^{-cm\delta^\beta} \right),$$

where $\mathfrak{M}_\beta(\delta)$ captures the tail-moment scaling from Lemma A.8. In particular,

$$\mathfrak{M}_1(\delta) = 1 + \log(1/\delta), \quad \mathfrak{M}_\beta(\delta) = O(1) \text{ for } \beta > 1, \quad \mathfrak{M}_\beta(\delta) = \Theta(\delta^{-(1-\beta)}) \text{ for } \beta \in (0, 1).$$

Optimizing this bound with $\delta \asymp (\log T/m)^{1/\beta}$ (paralleling the $\beta = 1$ specialization) yields rates of the form $\tilde{O}(m + T(\log T/m)^{1/\beta})$ for $\beta \geq 1$, matching the corresponding generalization of the greedy benchmark in Bayati et al. (2020). For $\beta < 1$, the same procedure yields a similar expression with an additional \mathfrak{M}_β contribution. The optimized rate is *not* $\tilde{O}(m + T/m)$ outside $\beta = 1$. The qualitative message is preserved: thicker upper tails make softmax leakage cheaper, and thinner tails give correspondingly slower rates.

6 Simulation Experiments

We complement the theoretical analysis with simulations on Bernoulli bandits. The section is also motivated by the bridge to GRPO: in practice, group-based policy optimization adds a KL penalty toward a reference base policy, and that base policy is rarely well-specified. We therefore design four experiments to (i) sanity-check the algorithm classes in the small-arm regime, (ii) verify the mechanism of Theorem 5.2 in the many-armed regime, (iii) test how a KL anchor toward an informative prior interacts with the regret behavior, and (iv) stress-test the behavior under prior misspecification.

Throughout, we compare four algorithm classes. Let $\hat{\mu}_{i,t}$ denote the empirical mean reward of arm i at time t (Section 3), and let S_i and F_i denote the running number of successes and failures observed from arm i , so $N_i = S_i + F_i$ is the per-arm pull count.

- **Classic Thompson Sampling (TS).** Under Beta–Bernoulli conjugacy, given a prior $\text{Beta}(\alpha_0, \beta_0)$ on μ_i , the posterior after S_i successes and F_i failures is $\text{Beta}(\alpha_0 + S_i, \beta_0 + F_i)$. At each step, sample $\theta_i \sim \text{Beta}(\alpha_0 + S_i, \beta_0 + F_i)$ independently for each arm and pull $\arg \max_i \theta_i$.
- **Empirical-mean greedy (Algorithm 1).** Pull $\arg \max_i \hat{\mu}_{i,t}$.
- **Constant-temperature softmax (variant of Algorithm 2).** Sample arms with probability $\pi(i) \propto \exp(\eta \cdot \hat{\mu}_{i,t})$, where $\eta > 0$ is a fixed inverse temperature. Unlike ASG, which uses an increasing schedule η_t , we hold η constant across rounds and vary it across runs.
- **Softmax + KL penalty.** A KL-regularized variant that anchors to a reference policy $\pi_0(i) \propto \exp(\eta \cdot \hat{\mu}_{i,0})$ derived from the prior means, selecting arms via $\pi(i) \propto \pi_0(i) \exp(\eta \cdot \hat{\mu}_{i,t})$. This models the KL-regularized objective structure common in GRPO-style updates (Section 1).

All experiments use Beta–Bernoulli conjugacy and average results over independent trials (200 for Experiment 1, 100 for Experiments 2–4). Shaded bands around each curve in the figures show 95% confidence intervals computed as $\pm 1.96 \cdot \text{SE}$, where $\text{SE} = \text{std} / \sqrt{n_{\text{trials}}}$ is the standard error of the mean across trials.

6.1 Experiment 1: Small-arm sanity check and effect of inverse temperature

The first experiment is a sanity check in the small-arm regime: with few arms, all empirical-mean-based methods accumulate roughly linear regret while Classic TS achieves the lowest regret, consistent with the known near-optimality of TS for moderate, fixed m . The experiment also illustrates the role of the inverse-temperature parameter η in the constant-temperature softmax.

We fix $m = 10$ arms, horizon $T = 1,000$, and an uninformative prior $\text{Beta}(1, 1)$ (uniform on $[0, 1]$). We vary the inverse temperature $\eta \in \{1, 5, 10, 20, 50\}$ for both the constant-temperature Softmax and Softmax+KL algorithms.

Figure 1 shows cumulative regret over time. At low inverse temperature ($\eta = 1$), the Softmax policy explores nearly uniformly and accumulates regret linearly, behaving essentially as a random policy. As η increases, the policy concentrates mass on arms with higher empirical means, and regret curves flatten toward the Greedy baseline. By $\eta = 50$, the Softmax algorithm is nearly indistinguishable from empirical-mean greedy. Classic TS achieves the lowest regret in this small- m regime, consistent with its known near-optimality for fixed, moderate m .

Under the uninformative prior, the Softmax+KL algorithm reduces to the standard Softmax algorithm (since π_0 is uniform), confirming that the KL penalty is inert when the prior carries no

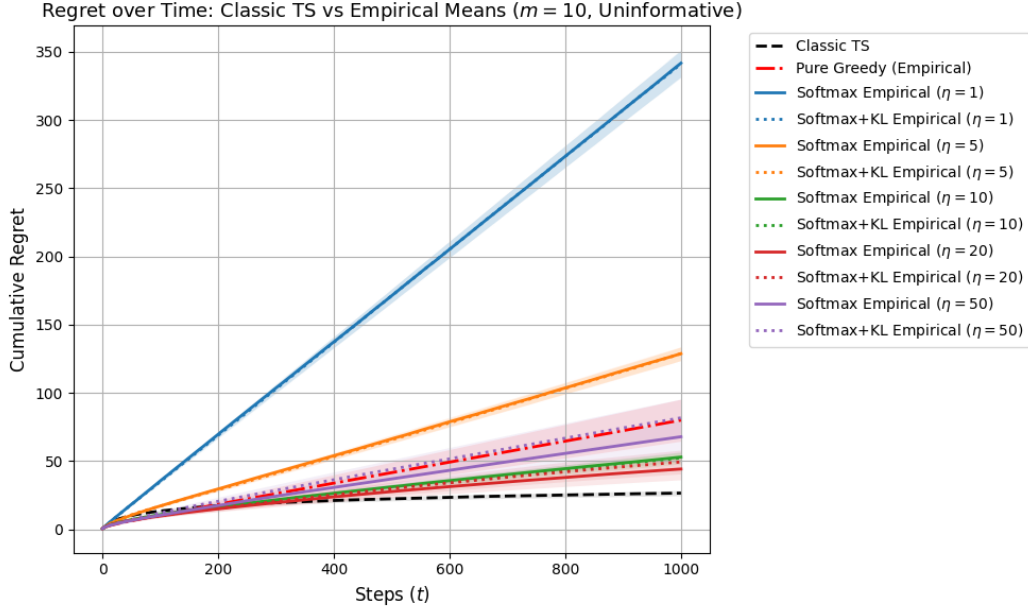


Figure 1: Cumulative regret over $T = 1,000$ steps with $m = 10$ arms and uninformative Beta(1, 1) prior. Higher inverse temperature η drives the Softmax policy toward greedy behavior. Under a uniform prior, Softmax+KL coincides with Softmax (overlapping curves).

useful information. This is visible in the overlapping solid and dotted curves of the same color in Figure 1.

Theorem 5.2 requires $\eta_t \rightarrow \infty$ (i.e. increasing inverse temperature over time) to make the leakage term $\sum_t \exp(-\eta_t \delta)$ summable. Experiment 1 illustrates the static analogue: higher constant η reduces leakage but with diminishing returns, and the transition from exploration-dominated to exploitation-dominated regret is smooth.

6.2 Experiment 2: Scaling with the number of arms (uninformative prior)

We vary $m \in \{10, 50, 100, 200, 500, 1,000\}$ with $T = 5,000$ and an uninformative Beta(1, 1) prior. Figure 2 reports the final cumulative regret R_T as a function of m .

The qualitative pattern is a performance inversion between Classic TS and the greedy-type algorithms as m grows. For small m , Classic TS dominates all other methods. However, as m increases past roughly $m \approx 200$, the Softmax algorithms with high inverse temperature ($\eta = 20, 50$) and empirical-mean greedy achieve lower regret than Classic TS, whose regret grows approximately linearly in m . Softmax ($\eta = 50$) and Greedy exhibit sublinear scaling, with regret plateauing or growing slowly.

This is the empirical counterpart of the many-armed regime analyzed in the paper. With $m = 1,000$ i.i.d. Uniform[0, 1] arms, the best arm has $\mu_* \approx 1 - 1/m$ with high probability (Lemma A.9), and many arms cluster near μ_* . The β -regularity condition (Definition 3.1) holds with $\beta = 1$ for the uniform distribution, placing the experiment in the regime of Theorem 5.2. Classic TS wastes pulls on clearly suboptimal arms because its posterior sampling randomization scales with the action space; greedy and high- η softmax exploit the abundance of near-optimal arms, the mechanism formalized in Lemmas A.4 and A.5.

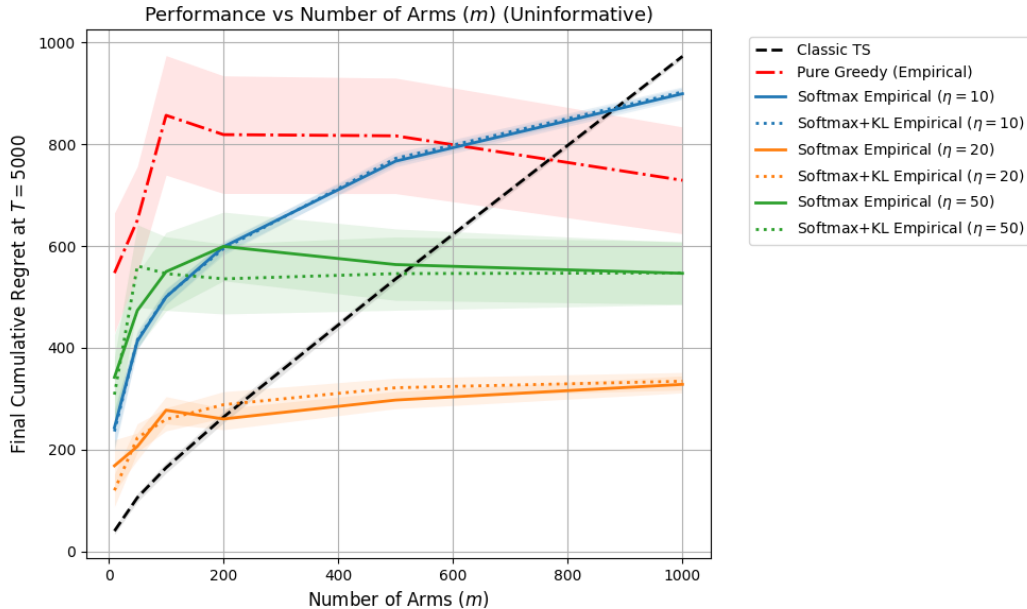


Figure 2: Final cumulative regret at $T = 5,000$ vs. number of arms m (uninformative prior). Greedy-type policies scale sublinearly, consistent with the $\tilde{O}(\sqrt{T})$ rate of Theorem 5.2 when $m = \Theta(\sqrt{T})$. Classic TS regret grows roughly linearly in m .

6.3 Experiment 3: Informative priors

We repeat the scaling experiment of Section 6.2 but initialize the prior using information correlated with the true arm probabilities:

$$\alpha_{0,i} = 1 + c\mu_i, \quad \beta_{0,i} = 1 + c(1 - \mu_i), \quad c = 2.$$

This yields informative Beta priors whose means $\alpha_{0,i}/(\alpha_{0,i} + \beta_{0,i})$ approximate the true μ_i .

Figure 3 shows that informative priors improve all exploitative algorithms. Pure Greedy achieves near-zero regret across all values of m , since the prior already identifies the best arm with high accuracy. The high- η Softmax algorithms ($\eta = 50$) and the Softmax+KL algorithms benefit substantially as well. The KL-regularized variant outperforms standard Softmax at every η , because the reference policy π_0 encodes accurate prior information that the KL penalty preserves. In effect, the KL anchor compounds the prior-mean signal with the empirical-mean signal.

This experiment illustrates the role of prior quality in the β -regular framework. When the prior accurately reflects the reward landscape, the “always-good” arms of Lemma A.4 are identified almost immediately. The KL-regularized variant, while not analyzed in our theoretical framework, suggests that anchoring to an informative prior can further suppress the leakage term, a direction we leave for future work.

6.4 Experiment 4: Misspecified priors

We stress-test the exploitative strategies by corrupting the informative prior with Gaussian noise ($\sigma = 0.1$) applied to the true probabilities before constructing the prior parameters. This misspecification means the prior’s “best arm” is likely not the true best arm.

Figure 4 shows that all exploitative algorithms suffer increased regret relative to the informative-prior setting, because they partially commit to the prior’s (now noisy) ranking. The KL anchor

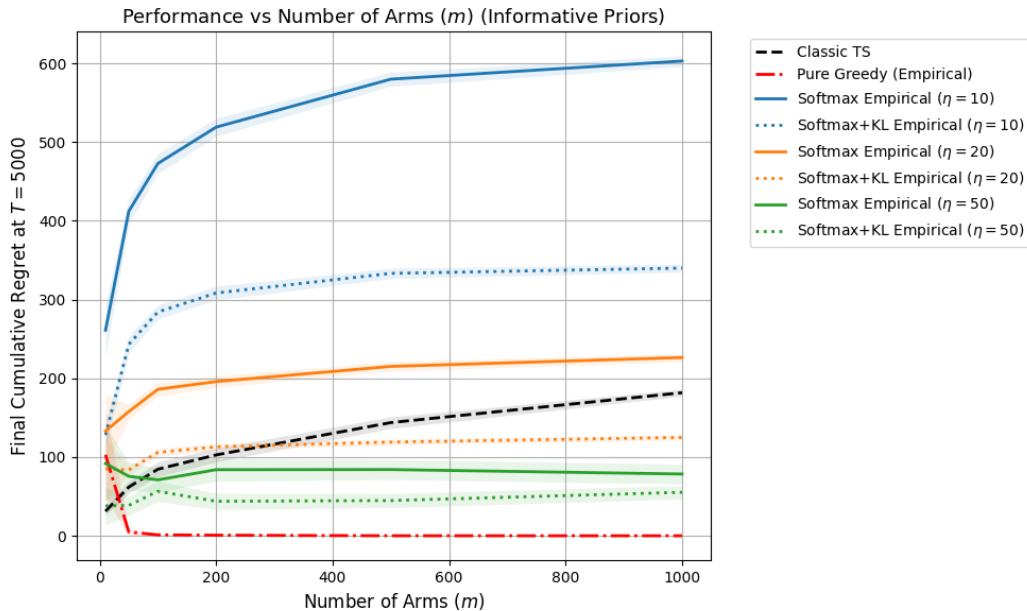


Figure 3: Final cumulative regret vs. m with informative priors ($c = 2$). Exploitative methods achieve lower regret; the KL-regularized Softmax benefits from the compounding of prior-mean and empirical-mean signals.

still helps at low η (Softmax+KL sits below plain Softmax for $\eta \in \{10, 20\}$) because anchoring to a partially-correct reference adds useful concentration when the unanchored softmax is too flat. At high η ($\eta = 50$), the empirical-mean signal already drives the softmax to concentrate, and the KL contribution shrinks toward zero. Classic TS, while still paying a large exploration cost for high m , recovers gracefully because its posterior sampling naturally adapts as empirical evidence accumulates.

This experiment probes the boundaries of our theoretical guarantees. Theorem 5.2 assumes a well-specified i.i.d. prior satisfying β -regularity, and the guarantee relies on many arms being genuinely near-optimal *a priori*. When the prior is misspecified, the empirical means are biased, and the always-good event $\mathcal{G}_i(\delta)$ (Equation 9) may fail to hold for the arms the policy concentrates on. The misspecified-prior experiment thus delineates the regime in which uncertainty-agnostic softmax policies are effective: they require that the prior (or the initial signal from the environment) provides a reasonable ordering of arm quality. This parallels the RLVR observation discussed in Section 1: softmax-style reweighting succeeds when the base policy already covers near-optimal completions, but cannot compensate for fundamental coverage failures.

6.5 Summary of empirical findings

Across experiments, three observations stand out. First, in the many-armed, uninformative-prior regime (Experiment 2), high- η softmax and empirical-mean greedy outperform Classic TS as m grows, consistent with the mechanism of Theorem 5.2: the abundance of near-optimal arms under a β -regular prior makes explicit exploration unnecessary. Second, the inverse-temperature parameter η smoothly interpolates between uniform exploration and greedy exploitation (Experiment 1), corresponding to the leakage–exploitation tradeoff captured by the $\sum_t \exp(-\eta_t \delta)$ term in (7). Third, the KL anchor amplifies the reference policy’s signal, with the marginal benefit

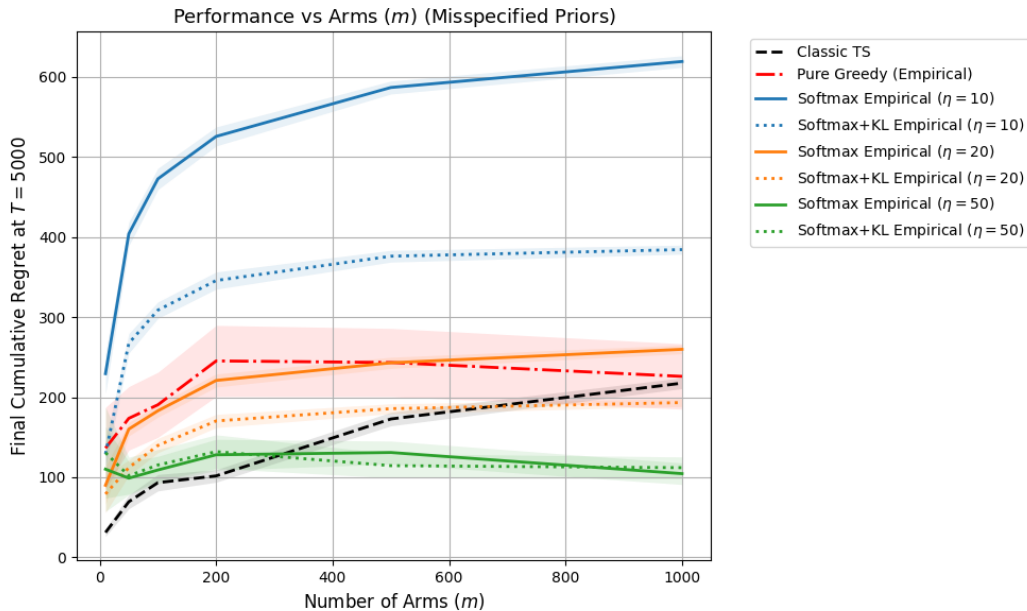


Figure 4: Final cumulative regret vs. m with misspecified priors ($\sigma = 0.1$ noise). Exploitative algorithms suffer from prior misspecification; the KL penalty amplifies this failure. Classic TS retains enough exploration to partially recover.

decreasing as η grows: it has no effect under a uniform reference (Experiment 1), helps at all η under an informative reference (Experiment 3), and still helps at low-to-moderate η under a mildly misspecified reference, with the benefit shrinking at high η where the empirical-mean signal alone drives the policy (Experiment 4).

7 Discussion

This paper studies whether annealed softmax over empirical mean rewards, despite ignoring epistemic uncertainty, can achieve near-optimal Bayes regret in the many-armed Bayesian Bernoulli bandit. Theorem 5.2 answers this positively in the 1-regular case of the upper-tail framework introduced by Bayati et al. (2020) for empirical-mean greedy: ASG attains the same near-optimal rate, up to logarithmic factors. The analysis extends to general $\beta > 0$ (Remark 5.3).

The mechanism is the abundance of near-optimal arms under β -regularity: with high probability, many arms maintain empirical means close to the optimum throughout learning. When softmax samples an arm other than the empirically best, that arm tends to be another near-optimal one rather than a clearly inferior one, so the price of softmax randomization is paid on near-optimal arms and contributes little to regret. Our proof builds on the analysis strategy of Bayati et al. (2020) for greedy, adding control of a softmax leakage term that quantifies the cost of placing some probability away from the empirical argmax.

This perspective does not contradict the classical negative results for Boltzmann exploration. In stochastic bandits with a small number of arms, monotone temperature schedules can amplify early noise and persistently concentrate on suboptimal arms (Cesa-Bianchi et al., 2017). Our theorem identifies a different regime, the many-armed Bayesian one with sufficient upper-tail mass, in which the same softmax policy achieves near-greedy Bayes regret. The difference is the geometry of the action space induced by the prior, not the policy class.

Two limitations remain. First, the model is narrow: we study i.i.d. Bayesian Bernoulli bandits, with no context, no shared structure across actions, and no sequential state dynamics. Second, the guarantees are Bayesian rather than minimax, and rely on upper-tail regularity of the prior; if near-optimal arms are rare, the mechanism weakens.

The most direct extension is to structured or contextual action spaces, and to sequential RL settings where “many near-optimal actions” might be replaced by many near-optimal trajectories or completions. More broadly, an open question is to characterize the boundary of the phenomenon: precisely when does upper-tail abundance make uncertainty-agnostic reweighting effective, and when do the classical Boltzmann failure modes re-emerge?

References

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 2312–2320, 2011. URL <https://proceedings.neurips.cc/paper/2011/hash/e1d5be1c7f2f456670de3d53c7b54f4a-Abstract.html>.
- Shipra Agrawal and Navin Goyal. Analysis of Thompson sampling for the multi-armed bandit problem. In *Proceedings of the 25th Annual Conference on Learning Theory (COLT)*, volume 23 of *Proceedings of Machine Learning Research*, pages 39.1–39.26. PMLR, 2012. URL <https://proceedings.mlr.press/v23/agrawal12.html>.
- Jean-Yves Audibert, Rémi Munos, and Csaba Szepesvári. Tuning bandit algorithms in stochastic environments. In *International Conference on Algorithmic Learning Theory (ALT)*, volume 4754 of *Lecture Notes in Computer Science*, pages 150–165. Springer, 2007. doi: 10.1007/978-3-540-75225-7_15. URL https://doi.org/10.1007/978-3-540-75225-7_15.
- Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2-3):235–256, 2002. doi: 10.1023/A:1013689704352. URL <https://doi.org/10.1023/A:1013689704352>.
- Hamsa Bastani, Mohsen Bayati, and Khashayar Khosravi. Mostly exploration-free algorithms for contextual bandits. *Management Science*, 67(3):1329–1349, 2020. doi: 10.1287/mnsc.2020.3605. URL <https://doi.org/10.1287/mnsc.2020.3605>.
- Mohsen Bayati, Nima Hamidi, Ramesh Johari, and Khashayar Khosravi. Unreasonable effectiveness of greedy algorithms in multi-armed bandit with many arms. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/12d16adf4a9355513f9d574b76087a08-Abstract.html>.
- Donald A. Berry, Robert W. Chen, Alan Zame, David C. Heath, and Larry A. Shepp. Bandit problems with infinitely many arms. *The Annals of Statistics*, 25(5):2103–2116, 1997. doi: 10.1214/aos/1069362389. URL <https://doi.org/10.1214/aos/1069362389>.
- Thomas Bonald and Alexandre Proutière. Two-target algorithms for infinite-armed bandits with Bernoulli rewards. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 2184–2192, 2013. URL <https://proceedings.neurips.cc/paper/2013/hash/fc2c7c47b918d0c2d792a719dfb602ef-Abstract.html>.

- Sébastien Bubeck and Nicolo Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012. doi: 10.1561/2200000024. URL <https://doi.org/10.1561/2200000024>.
- Alexandra Carpentier and Michal Valko. Simple regret for infinitely many armed bandits. In *International Conference on Machine Learning (ICML)*, volume 37 of *Proceedings of Machine Learning Research*, pages 1133–1141. PMLR, 2015. URL <https://proceedings.mlr.press/v37/carpentier15.html>.
- Nicolò Cesa-Bianchi, Claudio Gentile, Gábor Lugosi, and Gergely Neu. Boltzmann exploration done right. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. URL <https://arxiv.org/abs/1705.10257>.
- Arghya Roy Chaudhuri and Shivaram Kalyanakrishnan. Quantile-regret minimisation in infinitely many-armed bandits. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 425–434, 2018. URL <http://auai.org/uai2018/proceedings/papers/169.pdf>.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code, 2021. URL <https://arxiv.org/abs/2107.03374>.
- Ganqu Cui, Yuchen Zhang, Jiacheng Chen, Lifan Yuan, Zhi Wang, Yuxin Zuo, Haozhan Li, Yuchen Fan, Huayu Chen, Weize Chen, Zhiyuan Liu, Hao Peng, Lei Bai, Wanli Ouyang, Yu Cheng, Bowen Zhou, and Ning Ding. The entropy mechanism of reinforcement learning for reasoning language models, 2025. URL <https://arxiv.org/abs/2505.22617>.
- DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in LLMs via reinforcement learning, 2025. URL <https://arxiv.org/abs/2501.12948>. Published in Nature 645, 633–638 (2025), doi:10.1038/s41586-025-09422-z.
- John C. Gittins. Bandit processes and dynamic allocation indices. *Journal of the Royal Statistical Society: Series B (Methodological)*, 41(2):148–164, 1979. doi: 10.1111/j.2517-6161.1979.tb01068.x. URL <https://doi.org/10.1111/j.2517-6161.1979.tb01068.x>.
- Botao Hao, Tor Lattimore, and Csaba Szepesvári. Adaptive exploration in linear contextual bandit. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 108 of *Proceedings of Machine Learning Research*, pages 3536–3545. PMLR, 2020. URL <https://proceedings.mlr.press/v108/hao20b.html>.
- Sampath Kannan, Jamie H. Morgenstern, Aaron Roth, Bo Waggoner, and Zhiwei Steven Wu. A smoothed analysis of the greedy algorithm for the linear contextual bandit problem. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 2227–2236, 2018. URL <https://proceedings.neurips.cc/paper/2018/hash/2cfd4560539f887a5e420412b370b361-Abstract.html>.
- Emilie Kaufmann. On Bayesian index policies for sequential resource allocation. *The Annals of Statistics*, 46(2):842–865, 2018. doi: 10.1214/17-AOS1569. URL <https://doi.org/10.1214/17-AOS1569>.
- Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1):4–22, 1985. doi: 10.1016/0196-8858(85)90002-8. URL [https://doi.org/10.1016/0196-8858\(85\)90002-8](https://doi.org/10.1016/0196-8858(85)90002-8).

- Tor Lattimore and Csaba Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020. doi: 10.1017/9781108571401. URL <https://doi.org/10.1017/9781108571401>.
- Mingjie Liu, Shizhe Diao, Ximing Lu, Jian Hu, Xin Dong, Yejin Choi, Jan Kautz, and Yi Dong. Prorl: Prolonged reinforcement learning expands reasoning boundaries in large language models, 2025a. URL <https://arxiv.org/abs/2505.24864>.
- Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding R1-zero-like training: A critical perspective, 2025b. URL <https://arxiv.org/abs/2503.20783>.
- Youssef Mroueh. Reinforcement learning with verifiable rewards: GRPO’s effective loss, dynamics, and success amplification, 2025. URL <https://arxiv.org/abs/2503.06639>.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. 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, 2022. URL <https://arxiv.org/abs/2203.02155>. Advances in Neural Information Processing Systems 35 (NeurIPS 2022).
- Ruotian Peng, Yi Ren, Zhouliang Yu, Weiyang Liu, and Yandong Wen. Simko: Simple pass@K policy optimization, 2025. URL <https://arxiv.org/abs/2510.14807>.
- Manish Raghavan, Aleksandrs Slivkins, Jennifer Wortman Vaughan, and Zhiwei Steven Wu. The externalities of exploration and how data diversity helps exploitation. In *Conference on Learning Theory (COLT)*, volume 75 of *Proceedings of Machine Learning Research*, pages 1724–1738. PMLR, 2018. URL <https://proceedings.mlr.press/v75/raghavan18a.html>.
- Daniel Russo and Benjamin Van Roy. Learning to optimize via information-directed sampling. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 1583–1591, 2014a. URL <https://papers.nips.cc/paper/5463-learning-to-optimize-via-information-directed-sampling>.
- Daniel Russo and Benjamin Van Roy. Learning to optimize via posterior sampling. *Mathematics of Operations Research*, 39(4):1221–1243, 2014b. doi: 10.1287/moor.2014.0650. URL <https://doi.org/10.1287/moor.2014.0650>.
- Daniel Russo and Benjamin Van Roy. An information-theoretic analysis of Thompson sampling. *Journal of Machine Learning Research*, 17(68):1–30, 2016. URL <https://jmlr.org/papers/v17/14-087.html>.
- Daniel Russo and Benjamin Van Roy. Satisficing in time-sensitive bandit learning, 2018. URL <https://arxiv.org/abs/1803.02855>.
- John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, and Pieter Abbeel. Trust region policy optimization. In *International Conference on Machine Learning (ICML)*, volume 37 of *Proceedings of Machine Learning Research*, pages 1889–1897. PMLR, 2015. URL <https://proceedings.mlr.press/v37/schulman15.html>.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>.

- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y.K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL <https://arxiv.org/abs/2402.03300>.
- Zafir Stojanovski, Oliver Stanley, Joe Sharratt, Richard Jones, Abdulhakeem Adefioye, Jean Kadour, and Andreas Köpf. Reasoning gym: Reasoning environments for reinforcement learning with verifiable rewards, 2025. URL <https://arxiv.org/abs/2505.24760>. NeurIPS 2025 Spotlight.
- William R. Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933. doi: 10.2307/2332286. URL <https://doi.org/10.2307/2332286>.
- Yizao Wang, Jean-Yves Audibert, and Rémi Munos. Algorithms for infinitely many-armed bandits. In *Advances in Neural Information Processing Systems 21 (NeurIPS)*, pages 1729–1736, 2008. URL <https://proceedings.neurips.cc/paper/2008/hash/49ae49a23f67c759bf4fc791ba842aa2-Abstract.html>.
- Xumeng Wen, Zihan Liu, Shun Zheng, Shengyu Ye, Zhirong Wu, Yang Wang, Zhijian Xu, Xiao Liang, Junjie Li, Ziming Miao, Jiang Bian, and Mao Yang. Reinforcement learning with verifiable rewards implicitly incentivizes correct reasoning in base LLMs, 2025. URL <https://arxiv.org/abs/2506.14245>.
- Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, et al. DAPO: An open-source LLM reinforcement learning system at scale, 2025. URL <https://arxiv.org/abs/2503.14476>.
- Yang Yu. Pass@k metric for RLVR: A diagnostic tool of exploration, but not an objective, 2025. URL <https://arxiv.org/abs/2511.16231>.
- Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Shiji Song, and Gao Huang. Does reinforcement learning really incentivize reasoning capacity in LLMs beyond the base model? 2025. URL <https://arxiv.org/abs/2504.13837>.
- Guanning Zeng, Zhaoyi Zhou, Daman Arora, and Andrea Zanette. Shrinking the variance: Shrinkage baselines for reinforcement learning with verifiable rewards, 2025. URL <https://arxiv.org/abs/2511.03710>.
- Kaiyan Zhang, Yuxin Zuo, Bingxiang He, Youbang Sun, Runze Liu, Che Jiang, Yuchen Fan, Kai Tian, Guoli Jia, Pengfei Li, Yu Fu, Xingtai Lv, Yuchen Zhang, Sihang Zeng, Shang Qu, Haozhan Li, Shijie Wang, Yuru Wang, Xinwei Long, Fangfu Liu, Xiang Xu, Jiase Ma, Xuekai Zhu, Ermo Hua, Yihao Liu, Zonglin Li, Huayu Chen, Xiaoye Qu, Yafu Li, Weize Chen, Zhenzhao Yuan, Junqi Gao, Dong Li, Zhiyuan Ma, Ganqu Cui, Zhiyuan Liu, Biqing Qi, Ning Ding, and Bowen Zhou. A survey of reinforcement learning for large reasoning models, 2025. URL <https://arxiv.org/abs/2509.08827>.
- Heyang Zhao, Chenlu Ye, Quanquan Gu, and Tong Zhang. Sharp analysis for KL-regularized contextual bandits and RLHF, 2025. URL <https://openreview.net/forum?id=TE63KPCXWt>. NeurIPS 2025 poster; arXiv:2411.04625.

A Proofs

The proof technique follows the many-armed analysis of [Bayati et al. \(2020\)](#), adapted to control the softmax leakage probability. Even when the score vector is already informative, Boltzmann exploration still assigns some probability mass to inferior arms ([Cesa-Bianchi et al., 2017](#)). Under a 1-regular prior and many available arms, we show that the presence of many persistently good arms suppresses this leakage enough to recover the many-armed regret rate.

Throughout, \mathcal{F}_t denotes the history σ -field generated by actions and observations up to time t .

A.1 A Bernoulli “never-crossing” tail event

Definition A.1 (Never-crossing probability $q_\theta(\mu)$). Fix $\theta \in (0, 1)$ and $\mu \in [0, 1]$. Let $\{X_s\}_{s \geq 1}$ be i.i.d. Bernoulli(μ) and let $\hat{\mu}(n) := \frac{1}{n} \sum_{s=1}^n X_s$. Define

$$q_\theta(\mu) := \mathbb{P}(\hat{\mu}(n) > \theta \text{ for all } n \geq 1),$$

with the convention $q_\theta(\mu) = 0$ for $\mu \leq \theta$.

Lemma A.2 (Bernoulli crossing bound ([Bayati et al., 2020](#), Lemma 4.2)). *Let $X_s \sim \text{Bernoulli}(\mu)$ i.i.d. and fix $\theta > 2/3$. If $\mu \geq (1 + \theta)/2$, then*

$$q_\theta(\mu) \geq C_{\text{Bern}} := \frac{e^{-1/2}}{3}.$$

A.2 Many always-good arms under a β -regular tail

We now formalize the event that the prior produces many “always-good” arms, which will be used to control softmax leakage.

Fix $\delta \in (0, 1/8]$ and define two thresholds

$$\theta := 1 - 2\delta, \quad \theta' := 1 - 3\delta.$$

For each arm $i \in [m]$, define the *always-good* event

$$\mathcal{G}_i(\delta) := \left\{ \mu_i \geq 1 - \delta \right\} \cap \left\{ \hat{\mu}_i(n) > \theta \quad \forall n \geq 1 \right\}, \quad (9)$$

where $\hat{\mu}_i(n)$ denotes the empirical mean of the *first n rewards* of arm i (a property of the arm’s reward sequence, independent of the policy).

Let

$$M(\delta) := \sum_{i=1}^m \mathbf{1}\{\mathcal{G}_i(\delta)\}$$

be the number of always-good arms.

Lemma A.3 (Expected mass of always-good arms). *Assume Γ is 1-regular (Definition 3.1 with $\beta = 1$). There exists $\delta_0 > 0$ and a constant $c > 0$ such that for all $\delta \in (0, \delta_0]$,*

$$p_\delta := \mathbb{P}(\mathcal{G}_i(\delta)) = \mathbb{E}_{\mu \sim \Gamma} [\mathbf{1}\{\mu \geq 1 - \delta\} q_{1-2\delta}(\mu)] \geq c\delta.$$

Proof. By Lemma A.2, for $\delta < 1/6$ and $\mu \geq 1 - \delta$, we have $q_{1-2\delta}(\mu) \geq C_{\text{Bern}}$ since $1 - \delta \geq (1 + (1 - 2\delta))/2$. Therefore

$$p_\delta \geq C_{\text{Bern}} \cdot \mathbb{P}(\mu \geq 1 - \delta).$$

By 1-regularity, $\mathbb{P}(\mu \geq 1 - \delta) \geq c_0 \delta$ for all sufficiently small δ , hence $p_\delta \geq C_{\text{Bern}} c_0 \delta$. \square

Lemma A.4 (Many always-good arms with high probability). *Under the assumptions of Lemma A.3, there exist constants $c, c' > 0$ such that for all sufficiently small δ ,*

$$\mathbb{P}\left(M(\delta) \geq r(\delta)\right) \geq 1 - \exp(-c' m \delta), \quad r(\delta) := \max\left\{1, \left\lfloor \frac{mp_\delta}{2} \right\rfloor\right\}.$$

Proof. Across arms, the pairs $(\mu_i, (X_{i,s})_{s \geq 1})$ are i.i.d., so the indicators $\mathbf{1}\{\mathcal{G}_i(\delta)\}$ are i.i.d. Bernoulli with mean p_δ . Thus $M(\delta) \sim \text{Binomial}(m, p_\delta)$.

If $mp_\delta \geq 2$, then $r(\delta) = \lfloor mp_\delta/2 \rfloor$ and a multiplicative Chernoff bound gives

$$\mathbb{P}(M(\delta) < r(\delta)) \leq \mathbb{P}(M(\delta) < mp_\delta/2) \leq \exp(-mp_\delta/8).$$

If $mp_\delta < 2$, then $r(\delta) = 1$ and

$$\mathbb{P}(M(\delta) < r(\delta)) = \mathbb{P}(M(\delta) = 0) = (1 - p_\delta)^m \leq \exp(-mp_\delta).$$

In either case,

$$\mathbb{P}(M(\delta) < r(\delta)) \leq \exp(-c m p_\delta).$$

Lemma A.3 gives $p_\delta \geq c'\delta$, which yields the stated bound after adjusting constants. \square

A.3 A key inequality: softmax leakage controlled by r good arms

We next bound how often a suboptimal arm can be sampled when many good arms maintain a score margin.

Lemma A.5 (Pointwise leakage bound for softmax). *Fix time t , scores $\{s_j\}_{j=1}^m \subset \mathbb{R}$, an inverse temperature $\eta_t \geq 0$, and an integer $r \in \{1, \dots, m\}$. Let $p_t(i) \propto \exp(\eta_t s_i)$. If there is a subset $\mathcal{J} \subset [m]$ with $|\mathcal{J}| = r$ such that $s_j \geq s^*$ for all $j \in \mathcal{J}$, then for any i ,*

$$p_t(i) \leq \frac{\exp(\eta_t s_i)}{r \exp(\eta_t s^*)} = \frac{1}{r} \exp(-\eta_t (s^* - s_i)).$$

In particular, if $s^ - s_i \geq \Delta > 0$, then $p_t(i) \leq \frac{1}{r} e^{-\eta_t \Delta}$.*

Proof. Immediate from $\sum_{j=1}^m \exp(\eta_t s_j) \geq \sum_{j \in \mathcal{J}} \exp(\eta_t s_j) \geq r \exp(\eta_t s^*)$. \square

A.4 Bounding pull counts of suboptimal arms: spikes + leakage

We now use a simple “spikes + leakage” decomposition for empirical means.

Fix δ and thresholds $\theta = 1 - 2\delta$, $\theta' = 1 - 3\delta$ as above. For a fixed arm i , let

$$\widehat{\mu}_i(n) := \frac{1}{n} \sum_{s=1}^n X_{i,s}$$

be the empirical mean of the first n rewards of arm i .

Lemma A.6 (Spikes of empirical means are exponentially rare). *For any arm mean $\mu < \theta'$ and any $n \geq 1$,*

$$\mathbb{P}(\widehat{\mu}_i(n) \geq \theta' \mid \mu_i = \mu) \leq \exp(-2n(\theta' - \mu)^2).$$

Consequently,

$$\sum_{n=1}^{\infty} \mathbb{P}(\widehat{\mu}_i(n) \geq \theta' \mid \mu_i = \mu) \leq \frac{1}{2(\theta' - \mu)^2}.$$

Proof. The first claim is Hoeffding's inequality for Bernoulli averages. For the second, write $a = \theta' - \mu > 0$ and sum the geometric bound:

$$\sum_{n=1}^{\infty} e^{-2na^2} = \frac{e^{-2a^2}}{1 - e^{-2a^2}} = \frac{1}{e^{2a^2} - 1} \leq \frac{1}{2a^2},$$

using $e^x - 1 \geq x$ for $x \geq 0$. □

Good arms provide a uniform score floor. On the event $\mathcal{G}_i(\delta)$ in (9), we have $\widehat{\mu}_i(n) > \theta$ for every $n \geq 1$. Since each arm is pulled once during initialization, it follows that for every such arm and every $t \geq m + 1$,

$$\widehat{\mu}_{i,t-1} = \widehat{\mu}_i(N_i(t-1)) > \theta.$$

Thus, if $M(\delta) \geq r$, then throughout the softmax phase there are always at least r arms whose current empirical means exceed θ . This is the key simplification in the empirical-mean proof: no maturation argument is needed.

A.5 Pull-count bound conditional on having r good arms

We now state the pull-count lemma.

Lemma A.7 (Pull-count bound under r reference arms). *Fix an integer $r \in \{1, \dots, m\}$, fix δ , and set thresholds $\theta = 1 - 2\delta$, $\theta' = 1 - 3\delta$. For each $t \in \{m + 1, \dots, T\}$, define the \mathcal{F}_{t-1} -measurable event*

$$\mathcal{E}_t := \left\{ \exists \mathcal{J} \subset [m] \text{ with } |\mathcal{J}| = r \text{ such that } \widehat{\mu}_{j,t-1} \geq \theta \ \forall j \in \mathcal{J} \right\}.$$

Then for any arm i and any horizon T ,

$$\mathbb{E} \left[N_i(T) \mathbf{1} \left\{ \bigcap_{t=m+1}^T \mathcal{E}_t \right\} \right] \leq 1 + \sum_{n=1}^{\infty} \mathbb{P}(\widehat{\mu}_i(n) \geq \theta') + \frac{1}{r} \sum_{t=m+1}^T \exp(-\eta_t(\theta - \theta')). \quad (10)$$

Proof. Let $\tau_{i,n}$ be the time of the n -th pull of arm i , with $\tau_{i,n} = \infty$ if the n -th pull never occurs. If $\tau_{i,n+1} \leq T$, then at time $\tau_{i,n+1} - 1$ the arm has been observed exactly n times, hence

$$\widehat{\mu}_{i,\tau_{i,n+1}-1} = \widehat{\mu}_i(n).$$

Therefore

$$N_i(T) \mathbf{1} \left\{ \bigcap_{t=m+1}^T \mathcal{E}_t \right\} \leq 1 + \sum_{n=1}^{\infty} \mathbf{1} \left\{ \tau_{i,n+1} \leq T, \widehat{\mu}_i(n) \geq \theta' \right\} + \sum_{t=m+1}^T \mathbf{1} \left\{ A_t = i, \widehat{\mu}_{i,t-1} < \theta', \mathcal{E}_t \right\}.$$

Taking expectations, the first sum is bounded by

$$\sum_{n=1}^{\infty} \mathbb{P}(\widehat{\mu}_i(n) \geq \theta').$$

For the second sum, note that $\mathcal{E}_t \in \mathcal{F}_{t-1}$. Hence

$$\mathbb{E} \left[\mathbf{1} \{ A_t = i, \widehat{\mu}_{i,t-1} < \theta', \mathcal{E}_t \} \right] = \mathbb{E} \left[\mathbf{1} \{ \widehat{\mu}_{i,t-1} < \theta', \mathcal{E}_t \} \mathbb{P}(A_t = i \mid \mathcal{F}_{t-1}) \right].$$

On $\mathcal{E}_t \cap \{ \widehat{\mu}_{i,t-1} < \theta' \}$, Lemma A.5 gives

$$\mathbb{P}(A_t = i \mid \mathcal{F}_{t-1}) \leq \frac{1}{r} \exp(-\eta_t(\theta - \theta')).$$

Summing over t yields (10). □

A.6 A tail-moment lemma for integrating the spike bound

To convert (10) into Bayes regret, we need to integrate over μ the spike penalty. This is exactly where β -regularity yields logarithmic (for $\beta = 1$) or polynomial moment control.

Lemma A.8 (Tail-moment bound, general β). *Let Γ be β -regular (Definition 3.1) and let $\mu \sim \Gamma$. Fix $\delta \in (0, \varepsilon_0/8]$ and set $U := 1 - \mu$. Then there is a constant $C < \infty$ (depending only on $\beta, c_0, C_0, \varepsilon_0$) such that*

$$\mathbb{E} \left[\mathbf{1}\{U \geq 4\delta\} \left(1 + \frac{1}{U - 3\delta} \right) \right] \leq C \cdot \mathfrak{M}_\beta(\delta),$$

where

$$\mathfrak{M}_\beta(\delta) := \begin{cases} 1 + \log(1/\delta), & \beta = 1, \\ 1, & \beta > 1, \\ \delta^{-(1-\beta)}, & \beta \in (0, 1). \end{cases}$$

Proof. The proof is a straightforward integration-by-parts argument utilizing the β -regular bound $F(u) := \mathbb{P}(U \leq u) \leq C_0 u^\beta$ for $u \in (0, \varepsilon_0]$. Constants in the bounds below are allowed to depend on $(\beta, c_0, C_0, \varepsilon_0)$.

Writing

$$\mathbb{E} \left[\frac{1}{U - 3\delta} \mathbf{1}\{U \geq 4\delta\} \right] = \int_{4\delta}^1 \frac{1}{u - 3\delta} dF(u),$$

integration by parts yields the boundary contribution

$$\left[\frac{F(u)}{u - 3\delta} \right]_{4\delta}^1 = \frac{F(1)}{1 - 3\delta} - \frac{F(4\delta)}{\delta}.$$

Using the left-limit convention at the lower endpoint, the lower-boundary term is $-F((4\delta)^-)/\delta$, which is non-positive since the CDF F is non-negative; dropping it preserves the upper bound (and automatically absorbs any atom Γ may place at $u = 4\delta$). The boundary contribution is then at most $F(1)/(1 - 3\delta) = O(1)$. For the remaining integral, since $\delta \leq \varepsilon_0/8$ implies $u - 3\delta \geq \varepsilon_0/2$ on $[\varepsilon_0, 1]$, the contribution from $[\varepsilon_0, 1]$ is $O(1)$ (with constants depending on ε_0 , using $F(u) \leq 1$). On $[4\delta, \varepsilon_0]$, the bound $u - 3\delta \geq u/4$ gives

$$\int_{4\delta}^{\varepsilon_0} \frac{F(u)}{(u - 3\delta)^2} du \lesssim \int_{4\delta}^{\varepsilon_0} \frac{u^\beta}{u^2} du = \int_{4\delta}^{\varepsilon_0} u^{\beta-2} du,$$

which scales as $O(\log(1/\delta))$ for $\beta = 1$, $O(1)$ for $\beta > 1$, and $\Theta(\delta^{-(1-\beta)})$ for $\beta \in (0, 1)$. Adding the constant 1 term from $\mathbb{E}[\mathbf{1}\{U \geq 4\delta\}] \leq 1$ completes the proof. \square

A.7 Order statistics: size of the subtraction term

The next lemma records the size of the subtraction term in (3). It is not needed for the upper bound in Theorem 5.2, but it quantifies how much the exact identity can sharpen constants.

Lemma A.9 (Expected gap to 1 of the best of m arms). *If Γ is β -regular (Definition 3.1) and $\mu_* = \max_{i \leq m} \mu_i$, then*

$$\mathbb{E}[1 - \mu_*] \leq \frac{C}{m^{1/\beta}}$$

for a constant C depending only on $(\beta, c_0, C_0, \varepsilon_0)$. In particular, for $\beta = 1$, $\mathbb{E}[1 - \mu_*] \leq C/m$.

Proof. Let $U_i := 1 - \mu_i$ and $U_* := \min_{i \leq m} U_i = 1 - \mu_*$. By β -regularity, for small u we have $\mathbb{P}(U_i \leq u) \geq c_0 u^\beta$, hence

$$\mathbb{P}(U_* > u) = \mathbb{P}(U_1 > u)^m \leq (1 - c_0 u^\beta)^m \leq \exp(-c_0 m u^\beta) \quad (u \leq \varepsilon_0).$$

Integrating,

$$\mathbb{E}[U_*] = \int_0^\infty \mathbb{P}(U_* > u) du \leq \int_0^{\varepsilon_0} e^{-c_0 m u^\beta} du + (1 - \varepsilon_0) \mathbb{P}(U_* > \varepsilon_0) \leq C m^{-1/\beta},$$

using the change of variables $v = c_0 m u^\beta$. □

A.8 Proof of Theorem 5.2

Proof. Let

$$\theta := 1 - 2\delta, \quad \theta' := 1 - 3\delta, \quad p_\delta := \mathbb{P}(\mathcal{G}_i(\delta)), \quad r := \max\left\{1, \left\lfloor \frac{mp_\delta}{2} \right\rfloor\right\}.$$

Define the event

$$\mathcal{E} := \{M(\delta) \geq r\}.$$

By Lemma A.4,

$$\mathbb{P}(\mathcal{E}^c) \leq \exp(-cm\delta).$$

For each $t \in \{m+1, \dots, T\}$ let \mathcal{E}_t be the event from Lemma A.7. If \mathcal{E} occurs, then there are at least r arms j such that $\mathcal{G}_j(\delta)$ holds. For each such arm and each $t \geq m+1$,

$$\hat{\mu}_{j,t-1} = \hat{\mu}_j(N_j(t-1)) > \theta,$$

so $\mathcal{E} \subseteq \bigcap_{t=m+1}^T \mathcal{E}_t$. Hence Lemma A.7 implies that for every arm i ,

$$\mathbb{E}[N_i(T) \mathbf{1}\{\mathcal{E}\}] \leq 1 + \sum_{n=1}^{\infty} \mathbb{P}(\hat{\mu}_i(n) \geq \theta') + \frac{1}{r} \sum_{t=m+1}^T e^{-\eta_t \delta}. \quad (11)$$

Write $U_i := 1 - \mu_i$. Decompose the surrogate regret on \mathcal{E} as

$$\tilde{R}_T \mathbf{1}\{\mathcal{E}\} = \sum_{i=1}^m U_i N_i(T) \mathbf{1}\{\mathcal{E}\} \mathbf{1}\{U_i < 4\delta\} + \sum_{i=1}^m U_i N_i(T) \mathbf{1}\{\mathcal{E}\} \mathbf{1}\{U_i \geq 4\delta\}.$$

For the near-optimal arms,

$$\sum_{i=1}^m U_i N_i(T) \mathbf{1}\{\mathcal{E}\} \mathbf{1}\{U_i < 4\delta\} \leq 4\delta \sum_{i=1}^m N_i(T) = 4\delta T,$$

hence

$$\mathbb{E}\left[\sum_{i=1}^m U_i N_i(T) \mathbf{1}\{\mathcal{E}\} \mathbf{1}\{U_i < 4\delta\}\right] \leq 4T\delta. \quad (12)$$

For the remaining arms, we redo the pathwise decomposition from the proof of Lemma A.7, this time keeping the factor $U_i \mathbf{1}\{U_i \geq 4\delta\}$. On $\mathcal{E} \subseteq \bigcap_{t=m+1}^T \mathcal{E}_t$, that proof gives the pathwise inequality

$$N_i(T) \mathbf{1}\{\mathcal{E}\} \leq 1 + \sum_{n=1}^{\infty} \mathbf{1}\{\hat{\mu}_i(n) \geq \theta'\} + \sum_{t=m+1}^T \mathbf{1}\{A_t = i, \hat{\mu}_{i,t-1} < \theta', \mathcal{E}_t\}.$$

Multiply both sides by the nonnegative quantity $U_i \mathbf{1}\{U_i \geq 4\delta\}$, sum over i , and take expectation. The first two pathwise terms give the spike part on the right-hand side of (13) below; these will be handled by conditioning on μ_i . For the leakage part (the last sum), use the pathwise bound $U_i \mathbf{1}\{U_i \geq 4\delta\} \leq 1$ before taking expectation, and then repeat the conditioning step from the proof of Lemma A.7: on $\mathcal{E}_t \cap \{\widehat{\mu}_{i,t-1} < \theta'\}$, Lemma A.5 bounds $\mathbb{P}(A_t = i \mid \mathcal{F}_{t-1}) \leq \frac{1}{r} e^{-\eta_t \delta}$. The result is

$$\mathbb{E} \left[\sum_{i=1}^m U_i N_i(T) \mathbf{1}\{\mathcal{E}\} \mathbf{1}\{U_i \geq 4\delta\} \right] \leq \sum_{i=1}^m \mathbb{E} \left[\mathbf{1}\{U_i \geq 4\delta\} U_i \left(1 + \sum_{n=1}^{\infty} \mathbf{1}\{\widehat{\mu}_i(n) \geq \theta'\} \right) \right] + \frac{m}{r} \sum_{t=m+1}^T e^{-\eta_t \delta}. \quad (13)$$

Now condition on μ_i . Since the summands in the spike series are nonnegative, Tonelli's theorem justifies interchanging the expectation and the infinite sum. If $U_i \geq 4\delta$, then $\mu_i \leq 1 - 4\delta < \theta'$, so Lemma A.6 gives

$$\sum_{n=1}^{\infty} \mathbb{P}(\widehat{\mu}_i(n) \geq \theta' \mid \mu_i) \leq \frac{1}{2(\theta' - \mu_i)^2} = \frac{1}{2(U_i - 3\delta)^2}.$$

Therefore

$$\begin{aligned} \mathbb{E} \left[\mathbf{1}\{U_i \geq 4\delta\} U_i \left(1 + \sum_{n=1}^{\infty} \mathbf{1}\{\widehat{\mu}_i(n) \geq \theta'\} \right) \right] &= \mathbb{E} \left[\mathbf{1}\{U_i \geq 4\delta\} U_i \left(1 + \sum_{n=1}^{\infty} \mathbb{P}(\widehat{\mu}_i(n) \geq \theta' \mid \mu_i) \right) \right] \\ &\leq \mathbb{E} \left[\mathbf{1}\{U_i \geq 4\delta\} U_i \left(1 + \frac{1}{2(U_i - 3\delta)^2} \right) \right]. \end{aligned}$$

On $\{U_i \geq 4\delta\}$ we have $U_i - 3\delta \geq \delta$ and $U_i \leq 4(U_i - 3\delta)$, hence

$$U_i \left(1 + \frac{1}{2(U_i - 3\delta)^2} \right) \leq 4(U_i - 3\delta) + \frac{2}{U_i - 3\delta} \leq 6 \left(1 + \frac{1}{U_i - 3\delta} \right).$$

Applying Lemma A.8 with $\beta = 1$ yields

$$\mathbb{E} \left[\mathbf{1}\{U_i \geq 4\delta\} U_i \left(1 + \sum_{n=1}^{\infty} \mathbf{1}\{\widehat{\mu}_i(n) \geq \theta'\} \right) \right] \leq C(1 + \log(1/\delta)). \quad (14)$$

Summing (14) over i and combining with (13) gives

$$\mathbb{E} \left[\sum_{i=1}^m U_i N_i(T) \mathbf{1}\{\mathcal{E}\} \mathbf{1}\{U_i \geq 4\delta\} \right] \leq C m (1 + \log(1/\delta)) + \frac{m}{r} \sum_{t=m+1}^T e^{-\eta_t \delta}. \quad (15)$$

By Lemma A.3, $p_\delta \geq c\delta$ for all sufficiently small δ . If $mp_\delta \geq 2$, then $r = \lfloor mp_\delta/2 \rfloor$ and therefore

$$\frac{m}{r} \leq \frac{4}{p_\delta} \leq \frac{C}{\delta}.$$

If $mp_\delta < 2$, then $r = 1$ and

$$m < \frac{2}{p_\delta} \leq \frac{C}{\delta},$$

so again $m/r \leq C/\delta$. Combining (12) and (15),

$$\mathbb{E}[\widetilde{R}_T \mathbf{1}\{\mathcal{E}\}] \leq C \left[T\delta + m(1 + \log(1/\delta)) + \frac{1}{\delta} \sum_{t=1}^T e^{-\eta_t \delta} \right],$$

where we enlarged the sum from $t = m + 1$ to $t = 1$.

Finally,

$$\mathbb{E}[\tilde{R}_T] = \mathbb{E}[\tilde{R}_T \mathbf{1}\{\mathcal{E}\}] + \mathbb{E}[\tilde{R}_T \mathbf{1}\{\mathcal{E}^c\}] \leq \mathbb{E}[\tilde{R}_T \mathbf{1}\{\mathcal{E}\}] + T \mathbb{P}(\mathcal{E}^c),$$

so

$$\mathbb{E}[\tilde{R}_T] \leq C \left[m + T\delta + m(1 + \log(1/\delta)) + \frac{1}{\delta} \sum_{t=1}^T e^{-\eta_t \delta} + T e^{-cm\delta} \right].$$

Using (3), we have

$$\text{BR}_{T,m}(\text{ASG}) \leq \mathbb{E}[\tilde{R}_T],$$

which proves (7).

For the specialization (8), set $\delta = \min\{\delta_0, A \log(T \vee 2)/m\}$ with $A > 1/c$, and $\eta_t = (c_\eta/\delta) \log(t \vee 2)$ with $c_\eta > 1$. Then $\sum_{t=1}^T e^{-\eta_t \delta} = O(1)$.

If $A \log(T \vee 2)/m \leq \delta_0$, then $\delta = A \log(T \vee 2)/m$ and $T e^{-cm\delta} = T^{1-cA} \leq 1$ since $cA > 1$. The other terms in (7) are m , $T\delta = A \log(T \vee 2) \cdot T/m$, $m(1 + \log(1/\delta)) = \tilde{O}(m)$, and $(1/\delta) \sum_t e^{-\eta_t \delta} = \tilde{O}(m)$, summing to $\tilde{O}(m + T/m)$.

If $A \log(T \vee 2)/m > \delta_0$ (i.e., $m < A \log(T \vee 2)/\delta_0$), the cap binds and $T/m > \delta_0 T / (A \log(T \vee 2)) = \tilde{\Omega}(T)$, so $\tilde{O}(m + T/m)$ already absorbs the trivial bound $\text{BR}_{T,m}(\text{ASG}) \leq T$.

In both cases, $\text{BR}_{T,m}(\text{ASG}) = \tilde{O}(m + T/m)$. \square

B Extended Related Work

This appendix expands on the condensed survey in Section 2, providing additional context and citations.

Stochastic bandits: optimism and posterior sampling. The modern theory of stochastic multi-armed bandits characterizes the exploration–exploitation trade-off through regret guarantees and instance-dependent lower bounds [Lai and Robbins \(1985\)](#); [Auer et al. \(2002\)](#); [Audibert et al. \(2007\)](#); [Bubeck and Cesa-Bianchi \(2012\)](#); [Lattimore and Szepesvári \(2020\)](#). A canonical Bayesian approach is posterior sampling (Thompson sampling), originally proposed in [Thompson \(1933\)](#) and analyzed in modern finite-time settings by, e.g., [Agrawal and Goyal \(2012\)](#); [Russo and Van Roy \(2016\)](#). Related Bayesian decision rules include Bayesian index policies [Kaufmann \(2018\)](#) and the classical Gittins index [Gittins \(1979\)](#). Our setting adopts the standard Beta–Bernoulli conjugate model, which enables a clean comparison between probability matching (posterior sampling) and softmax/Boltzmann action selection.

Boltzmann (softmax) exploration and its limitations. Boltzmann exploration (a.k.a. softmax or Gibbs action selection) is widely used in reinforcement learning and bandits as a simple randomized alternative to greedy choice. Despite its popularity, [Cesa-Bianchi et al. \(2017\)](#) show that, for stochastic K -armed bandits, *any monotone* temperature (learning-rate) schedule can be forced into suboptimal behavior: either it explores too long or it commits too early. They propose remedies including (i) tuned non-monotone schedules that require knowledge of the horizon and gaps and (ii) per-arm learning rates that explicitly track estimation uncertainty [Cesa-Bianchi et al. \(2017\)](#). This negative result motivates our focus on regimes where *uncertainty-aware* schedules may be unnecessary due to structural properties of the arm distribution.

Many-armed and infinite-armed bandits. A long line of work studies bandits with infinitely many arms and tail-based performance criteria, emphasizing how the distribution of arm qualities shapes achievable regret. Classical and modern examples include infinitely-many-armed formulations [Berry et al. \(1997\)](#); [Wang et al. \(2008\)](#); [Bonald and Proutière \(2013\)](#); [Carpentier and Valko \(2015\)](#); [Chaudhuri and Kalyanakrishnan \(2018\)](#). In the Bayesian many-armed regime, [Bayati et al. \(2020\)](#) formalize the idea that when the prior places substantial mass on near-optimal arms (via upper-tail regularity conditions), greedy-style policies can enjoy *free exploration*: discarding a poorly performing arm is likely to leave other near-optimal arms available. They show that subsampled greedy can achieve Bayesian regret scaling of order $\tilde{O}(\max\{m, T/m\})$ (up to prior-dependent exponents and logarithms), implying near-optimal performance with $m \asymp \sqrt{T}$ arms [Bayati et al. \(2020\)](#). Our analysis builds on this viewpoint, but focuses on softmax/Boltzmann policies that randomize *without* explicit epistemic-uncertainty bonuses.

Free exploration beyond non-contextual bandits. Complementary “free exploration” phenomena have been identified in contextual bandits, where exploration can be induced by natural diversity in contexts and data [Bastani et al. \(2020\)](#); [Kannan et al. \(2018\)](#); [Raghavan et al. \(2018\)](#); [Hao et al. \(2020\)](#); [Abbasi-Yadkori et al. \(2011\)](#). These works highlight that the need for explicit exploration is sensitive to structural assumptions (e.g., covariate diversity, smoothed analysis, or rich action sets). Our work studies an orthogonal mechanism: even *without* context diversity, the presence of many near-optimal arms (captured by prior tail regularity) can make epistemic-uncertainty-agnostic softmax policies achieve near-greedy Bayes regret.

Satisficing and multiple near-optimal actions. When near-optimal actions are plentiful, identifying the unique optimal action may be information-inefficient. This perspective is formalized in satisficing and information-theoretic approaches to bandits [Russo and Van Roy \(2014a,b, 2018\)](#). Conceptually, our “many near-optimal arms” regime is aligned with satisficing: regret can remain small even if the learner settles on an ε -optimal arm, provided such arms are sufficiently common under the prior.

RLVR and group-based policy optimization for LLM reasoning. The empirical motivation for this paper comes from the rapid adoption of reinforcement learning with verifiable (often binary) rewards for post-training large language models (LLMs) on reasoning-centric tasks. GRPO was introduced in DeepSeekMath as a memory-efficient alternative to PPO-style actor-critic training for verifiable rewards [Shao et al. \(2024\)](#); [Schulman et al. \(2017\)](#); it has since become a widely used baseline for RLVR training at scale [DeepSeek-AI \(2025\)](#); [Yu et al. \(2025\)](#). Several recent works analyze GRPO/RLVR training dynamics and relate them to KL-regularized reweighting. For example, [Mroueh \(2025\)](#) derive explicit optimal-policy forms for variants of GRPO under binary rewards and show how success probability can be amplified through iterative updates. At the same time, there is ongoing debate about whether RLVR expands a model’s reasoning *support* (high- k coverage) or mainly reweights probability mass within the base model’s existing support. [Yue et al. \(2025\)](#) report that RLVR often improves small- k performance while failing to improve—and sometimes degrading—large- k pass@ k , suggesting a “bounded-by-base” effect; subsequent work revisits this question and argues that RLVR can extend reasoning boundaries under specific protocols and metrics [Wen et al. \(2025\)](#); [Liu et al. \(2025a\)](#). Additional analyses emphasize training instabilities, entropy collapse, and variance reduction techniques in RLVR pipelines [Cui et al. \(2025\)](#); [Zeng et al. \(2025\)](#); [Liu et al. \(2025b\)](#); [Yu et al. \(2025\)](#).

Pass@ k as a tail metric and as an objective. Pass@ k was popularized as a functional-correctness evaluation metric for code generation and is an order-statistic probe of a model’s upper tail over solutions [Chen et al. \(2021\)](#). Recent RLVR-specific work argues that pass@ k is best interpreted as a diagnostic of exploration/coverage rather than a direct optimization target [Yu \(2025\)](#), and proposes alternative objectives to mitigate probability concentration and improve pass@ k at larger k [Peng et al. \(2025\)](#). Our bandit model adopts an explicit tail-regularity assumption on arm quality; this provides a stylized bridge to pass@ k phenomena by translating “good pass@ k ” into “many near-optimal arms with non-negligible mass.”

Connections to KL-regularized policy learning. A recurring theme in RLHF/RLVR is that KL regularization (to a reference or previous policy) acts as a trust-region or mirror-descent constraint, inducing stochastic policies that resemble Gibbs distributions [Schulman et al. \(2015, 2017\)](#); [Ouyang et al. \(2022\)](#). Very recent theory work in contextual bandits and RLHF argues that KL regularization alone can induce sufficient exploration under suitable coverage assumptions [Zhao et al. \(2025\)](#). Our contributions are complementary: we identify a distinct mechanism—prior tail mass / abundance of near-optimal arms—under which uncertainty-agnostic annealed softmax can achieve strong Bayesian regret, offering an alternative explanation for why soft reweighting can be effective even without explicit epistemic exploration.

Surveys and resources. For broader perspective on RL for large reasoning models, including RLVR training recipes, datasets, and open problems, see the recent survey [Zhang et al. \(2025\)](#). Reasoning Gym provides a large suite of procedurally generated, verifiable environments intended for RLVR research [Stojanovski et al. \(2025\)](#).