



**HAL**  
open science

## Approximate information maximization for bandit games

Alex Barbier-Chebbah, Christian L. Vestergaard, Jean-Baptiste Masson,  
Etienne Boursier

### ► To cite this version:

Alex Barbier-Chebbah, Christian L. Vestergaard, Jean-Baptiste Masson, Etienne Boursier. Approximate information maximization for bandit games. 2024. hal-04246907v4

**HAL Id: hal-04246907**

**<https://hal.science/hal-04246907v4>**

Preprint submitted on 28 Nov 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

---

# Approximate information maximization for bandit games

---

**Alex Barbier–Chebbah**

Institut Pasteur, Université Paris Cité  
Epimethee INRIA Paris France  
CNRS UMR 3571, Paris France  
alex.barbier-chebbah@pasteur.fr

**Christian L. Vestergaard**

Institut Pasteur, Université Paris Cité  
Epimethee INRIA Paris France  
CNRS UMR 3571, Paris France

**Jean-Baptiste Masson**

Institut Pasteur, Université Paris Cité  
Epimethee INRIA Paris France  
CNRS UMR 3571, Paris France

**Etienne Boursier**

INRIA, Université Paris Saclay  
LMO, Orsay, France

## Abstract

Entropy maximization and free energy minimization are general physics principles for modeling dynamic systems. Notable examples include modeling decision-making within the brain using the free-energy principle, optimizing the accuracy-complexity trade-off when accessing hidden variables with the information bottleneck principle [Tishby et al., 2000], and navigation in random environments using information maximization [Vergassola et al., 2007]. Building on this principle, we propose a new class of bandit algorithms that maximize an approximation to the information of a key variable within the system. To this end, we develop an approximated, analytical physics-based representation of the entropy to forecast the information gain of each action and greedily choose the one with the largest information gain. This method yields strong performances in classical bandit settings. Motivated by its empirical success, we prove its asymptotic optimality for the multi-armed bandit problem with Gaussian rewards. Since it encompasses the system’s properties in a single, global functional, this approach can be efficiently adapted to more complex bandit settings. This calls for further investigation of information maximization approaches for bandit problems.

## 1 Introduction

Multi-armed bandit problems have attracted wide attention in the past decades. They embody the challenge of balancing exploration and exploitation and have been applied to various settings such as online recommendation [Bresler et al., 2014], medical trials [Thompson, 1933], dynamic pricing [Den Boer, 2015], and reinforcement learning-based decision making [Silver et al., 2016, Ryzhov et al., 2012]. Besides the classic stochastic version of the multi-armed bandit problem, many subsequent extensions have been developed, providing richer models for specific applications. These extensions include linear bandits [Li et al., 2010], many-armed bandits [Bayati et al., 2020], and pure exploration problems such as thresholding bandits [Locatelli et al., 2016] or top-K bandits [Kalyanakrishnan et al., 2012, Kaufmann et al., 2016].

In the classic setting, an agent chooses an arm at each time step and observes a stochastic reward. Since they only observe the payoff of the chosen arm, the agent should regularly explore suboptimal arms. This is often referred to as the exploration-exploitation trade-off. An agent can exploit its

current knowledge to optimize gains by drawing the current empirically best arm or exploring other arms to potentially increase future gains.

Optimal strategies are characterized asymptotically by the Lai and Robbins bound [Lai et al., 1985]. Among them, upper confidence bound [UCB, Auer, 2000, Garivier and Cappé, 2011] methods greedily pull the arm maximizing some tuned confidence index; Thompson sampling [Kaufmann et al., 2012a, Agrawal and Goyal, 2013] relies on sampling mean rewards from a posterior distribution and chooses the arm with the largest random sample; deterministic minimum empirical divergence [DMED, Honda and Takemura, 2010] builds on a balance between the maximum likelihood of an arm being the best and the posterior expectation of the regret.

Even if these approaches efficiently utilize current available information, they do not aim directly to acquire more information. We highlight, however, the information directed sampling approach (IDS) of Russo and Van Roy [2014], which relies on a measure of the information gain of the optimal actions. By leveraging an information measure that consistently captures the specific problem structure, IDS can address general classes of problems, particularly those with a complex information structure where classic bandit methods fall short. Surprisingly, IDS can even outperform UCB and Thompson sampling in classic bandit problems. However, like DMED, this method explicitly balances information gain with expected losses induced by exploration, and the efficiency of pure information-maximizing strategies thus remains to be proven.

Information-maximization approaches provide a decision-making strategy in which the agent tries to maximize information about one or more relevant stochastic variables. The information-maximizing principle has shown to be efficient in a broad range of domains [Helias and Dahmen, 2020, Parr et al., 2022, Hernández-Lobato et al., 2015, Vergassola et al., 2007] where decisions have to be taken in fluctuating or unknown environments. These domains include, e.g., robotics applications [Zhang et al., 2015], where the ability to share approximate information improves collective decisions, and the search for olfactory sources in turbulent flows [Masson, 2013, Reddy et al., 2022].

In the specific setting of bandit problems, information maximization has shown promising empirical results, and heuristic arguments support its asymptotic optimality [Reddy et al., 2016, Barbier-Chebbah et al., 2023]. As IDS, they leverage information structure to provide a versatile decision framework with the capability to address various bandits settings. However, the efficiency of such a “pure exploration” strategy in terms of regret minimization has yet to be proven, and it has been previously argued that it would result in a linear regret [Russo and Van Roy, 2014]. Moreover, current information-based algorithms often rely on complex numerical integration, leading to high computational costs, a significant challenge that information-based methods must overcome. In this context, we aim to leverage new strategies derived from information maximization principles focusing on global observables, *i.e.* variables depending on more than one bandit arm, that alleviate the computational burden of numerical evaluation of complex functionals and to rigorously prove their efficiency.

**Contributions.** Our main contribution is introducing a new class of asymptotically optimal algorithms that rely on approximations of a functional representing the current information of interest about the whole bandit system. This approach is based on the entropy of the posterior mean value of the best arm, for which we provide an approximate expression to enable robust, easily tunable, and extendable algorithms with a direct analytical formulation. We focus here on the multi-armed bandit problem with Gaussian rewards, for which we derive a simple approximate information maximization algorithm (AIM) and provide an upper bound on its pseudo-regret, ensuring that AIM is asymptotically optimal. The information from each arm is incorporated in a unique entropy functional, which shows promise for tackling more complex bandit settings such as linear bandit or many-armed bandits. Thus, our main motivation is to design an analytical functional-based algorithmic principle, which can potentially address problems with more correlated information structures in the future. Additionally, another strength of AIM lies in its short-time behavior, where it shows strong performances as we illustrate numerically for both Bernoulli and Gaussian rewards.

**Organization.** In Section 2, we briefly review the  $K$ -armed bandit setting. Section 3 presents the general principle of information maximization approaches, originally inspired by both the information bottleneck principle and navigation in turbulent plumes. Section 4 upper bounds the regret of AIM, showing it attains Lai and Robbin’s asymptotic bound. In Section 5, the performance of AIM is

numerically compared with known baselines on multiple examples. Finally, Section 6 discusses extensions of AIM to various bandit settings.

## 2 Setting

We consider the classic  $K$ -armed stochastic bandit game. In each round  $t$ , the agent selects an arm  $a_t \in [K] = \{1, \dots, K\}$  among a set of  $K$  choices solely based on the rewards of the previously pulled arms. The chosen arm  $k$  then returns a stochastic reward  $X_t(k)$ , drawn independently of the previous rounds, according to a distribution  $\nu_k$  of mean  $\mu_k$ . We denote by  $N_k(t)$  the number of times the arm  $k$  has been pulled. When clear from context, we omit the dependence on  $t$  for simplicity.

The goal of the agent is to maximize its cumulative reward, or equivalently, to minimize its pseudo-regret up to round  $T$ , defined as

$$R(T) = \mu^* T - \sum_{t=1}^T \mathbb{E}[\mu_{a_t}], \quad (1)$$

where  $\mu^* = \max_{i \in [K]} \mu_i$ . Hence, the agent will optimize its choice of  $a_t$  relying on the previous observations up to  $t$ . For a large family of reward distributions, the asymptotic pseudo-regret is lower-bounded for any uniformly good policy by

$$\liminf_{T \rightarrow \infty} \frac{R(T)}{\ln(T)} \geq \sum_{k, \mu_k < \mu^*} \frac{\mu^* - \mu_k}{D_{\text{KL}}(\nu_k \| \nu_{k^*})}, \quad (2)$$

where  $k^* \in \operatorname{argmax}_{i \in [K]} \mu_i$ , and  $D_{\text{KL}}(\nu_k \| \nu_{k^*})$  denotes the Kullback-Leibler divergence between the reward distributions of the arms  $k$  and  $k^*$  [Lai et al., 1985]. In the particular case of Gaussian rewards with equal variances, i.e.,  $\nu_i = \mathcal{N}(\mu_i, \sigma^2)$ , the Kullback-Leibler divergence is  $D_{\text{KL}}(\nu_k \| \nu_{k^*}) = (\mu^* - \mu_k)^2 / (2\sigma^2)$ .

## 3 Information maximization strategies

Here, we introduce entropy-based, information maximization strategies and their underlying physical principles. We then detail approximations leading to an analytical and simplified entropy functional, which is the basis of the AIM algorithm.

### 3.1 Algorithm design principle: physical intuition

We aim to design a functional encompassing the current available information of the full system. Inspired by the information maximization principle [Vergassola et al., 2007, Reddy et al., 2016] which has revealed effective in taxis strategies where the agent needs to find an emitting odour source [Martinez et al., 2014, Cardé, 2021, Murlis et al., 1992], we rely on an entropic functional for policy decision. More precisely, we choose  $S_{\max}$ , the entropy of the posterior distribution of the value of the maximal mean reward, denoted  $p_{\max}$ .

The algorithm relies on an arbitrary prior distribution on the arm mean rewards. With independent arm priors, the posterior distribution of the value of the maximal mean reward can be expressed as

$$p_{\max}(\theta) d\theta = d\mathbb{P} \left( \max_k \mu_k = \theta \mid \mathcal{F}_{t-1} \right) = \sum_{k=1}^K d\mathbb{P}(\mu_k = \theta \mid \mathcal{F}_{t-1}) \prod_{j \neq k} \mathbb{P}(\mu_j \leq \theta \mid \mathcal{F}_{t-1}), \quad (3)$$

where  $\mathcal{F}_{t-1} = \sigma(X_1(a_1), \dots, X_{t-1}(a_{t-1}))$  denotes the filtration associated to the observations up to time  $t-1$ . The associated entropy reads

$$S_{\max} = - \int_{\Theta} p_{\max}(\theta) \ln p_{\max}(\theta) d\theta, \quad (4)$$

where  $\Theta = [\mu_{\inf}, \mu_{\sup}]$  is the support of  $p_{\max}$  (which depends on the nature of the game and can be infinite). Note that, as exemplified by Equation (3),  $p_{\max}$  includes the arms' priors and directly

depends on the reward distributions.<sup>1</sup> The entropy,  $S_{\max}$ , is a measure of the information carried by all arms in a single functional, providing a global state description of the game.

Our policy aims to minimize the entropy of  $p_{\max}$ . For that, it greedily chooses the arm providing the largest expected decrease in entropy, conditioned on the current knowledge of the game,

$$\operatorname{argmin}_{k \in [K]} \mathbb{E} [S_{\max}(t) - S_{\max}(t-1) \mid \mathcal{F}_{t-1}, a_t = k]. \quad (5)$$

Similar to Thompson sampling, it relies on a Bayesian representation. Yet, it distinguishes itself by providing a deterministic decision procedure given past observations. We stress that  $S_{\max}$  quantifies the available information about the average reward of the best arm. This choice contrasts with using the entropy of the probability of the best arm, which is known to overexplore and is suboptimal for regret minimization [Reddy et al., 2016]. Because of this suboptimality, approaches based on the information on the best arm fix this concern by including the expected regret in the functional to favor exploitation [Russo and Van Roy, 2014]. Furthermore, we argue that by the definition of  $p_{\max}$ , the information carried by the arms’ posteriors is sufficiently mixed to ensure an optimal behaviour, as proved in Section 4. Since the policy aims to maximize the information about the best arm’s mean, it mainly pulls the current best arm to learn more about its value. On the contrary, policies aiming to identify the best arm pull worse empirical arms more often because they are only concerned about the arms’ order.

The information maximization policy based on Equation (5) has been empirically shown to be competitive with state-of-the-art algorithms [Reddy et al., 2016] and robust to variations of the prior [Reddy et al., 2016, Barbier-Chebbah et al., 2023] in classic bandit games. However, while Equation (5) can be numerically evaluated, it cannot be computed in closed form, preventing the gradient from being analytically tractable. This makes intricate to theoretically bound the regret even in the two-armed setting and it also prevents the policy’s extension to more complicated bandit settings. Additionally, it induces a high computational cost [a trait shared with IDS Russo and Van Roy, 2014], which becomes disadvantageous when considering a large number of arms and at large times (when  $p_{\max}$  is peaked), where one has to manage vanishing numerical precision, making the numerical integration even longer. Finally, the integral form of  $S_{\max}$  prevents fine-tuning, which could prove crucial for achieving or surpassing the empirical state-of-the-art performances.

A second simplified and analytical functional mirroring  $S_{\max}$  has to be derived to address these concerns. This analytical result strengthens the information maximization principle, both by providing novel algorithms that are analytical, tractable and computationally efficient while conserving the main advantages of the exact entropy [Reddy et al., 2016] and by making theoretical analysis tractable.

### 3.2 Main elements of the entropy analytical approximation

Here, we devise a set of approximations of  $p_{\max}$  and  $S_{\max}$  to get a tractable analytical algorithm. Given that the best empirical arm and the worse empirical arms have notably distinct contributions to  $p_{\max}$  (Figure 1(a)), we approximate  $p_{\max}$  while considering the current arms’ order. We sort them based on their current posterior means, labelling the highest one as  $M_t$  (with an empirical reward of  $\hat{\mu}_{M_t}$ ) and  $\mathcal{A}_t = [K] \setminus \{M_t\}$  the set of worse empirical arms. Of course,  $M_t$  might differ from the actual optimal arm  $k^*$  due to the randomness in the observed rewards. We focus on approximating Equations (3) and (4) when the best empirical arm has already been extensively drawn more often than the other arms.

The entropy is then decomposed into two tractable terms corresponding to distinct behaviors of  $p_{\max}(\theta)$  when  $\theta$  varies:

$$\tilde{S}_{\max} = \tilde{S}_{\text{body}} + \tilde{S}_{\text{tail}}, \quad (6)$$

The first term,  $\tilde{S}_{\text{body}}$ , approximates the contribution around the mode of  $p_{\max}$ , while the second term,  $\tilde{S}_{\text{tail}}$ , quantifies the information carried by the tail of  $p_{\max}$  (corresponding to high rewards, see Figure 1). Each of these terms then corresponds to a part of the entropy where the dominant term of Equation (3) is distinct (see Appendix A.I for details).

---

<sup>1</sup>In the remainder of the paper, we consider an improper uniform prior over  $\mathbb{R}$ , as often considered with Gaussian rewards.

More precisely, by denoting  $p_i(\theta)d\theta = d\mathbb{P}(\mu_i = \theta \mid \mathcal{F}_t)$  the mean posterior density of the associated arm  $i$ , the tail term is approximated as:

$$\tilde{S}_{\text{tail}} = - \sum_{m \in \mathcal{A}_t} \int_{\tilde{\mu}_{\text{eq},m}}^{\mu_{\text{sup}}} p_m(\theta) \ln p_m(\theta) d\theta. \quad (7)$$

where  $\tilde{\mu}_{\text{eq},m}$ , given in Appendix A.6, approximates  $\bar{\mu}_{\text{eq},m}$ , the value of  $\theta$  where the empirical best arm  $M_t$  and the selected worse arm  $m$  have the same probability of being the best arm (see red and orange curves in Figure 1(b)). Here,  $p_m(\theta)$  is the posterior density of the current worse arm evaluated at  $\theta$ . Roughly, because the better empirical arm has been predominantly drawn,  $p_{M_t}(\theta)$  decays faster than  $p_m(\theta)$ , resulting in a tail term (see Equation (7)) whose main contribution is the worse empirical arm. The approximate entropy of the body component is:

$$\tilde{S}_{\text{body}} = - \int_{\Theta} \left(1 - \sum_{m \in \mathcal{A}_t} [1 - C_m(\theta)]\right) p_{M_t}(\theta) \ln p_{M_t}(\theta) d\theta \quad (8)$$

where  $C_i(\theta) = \mathbb{P}(\theta > \mu_i \mid \mathcal{F}_t)$  is the cumulative posterior probability of the mean of the arm  $i$ . Equation (8) is the leading-order term of the mode of  $p_{\text{max}}$ , which is mainly contributed to by the best empirical arm.

This approximation of Equation (3) is good when the best empirical arm has been extensively drawn compared to the worse empirical ones, corresponding to the situation encountered asymptotically for uniformly good algorithms. Surprisingly, the approximation captured by Equation (6) is still accurate enough outside this asymptotic regime to provide a high-performance decision scheme.

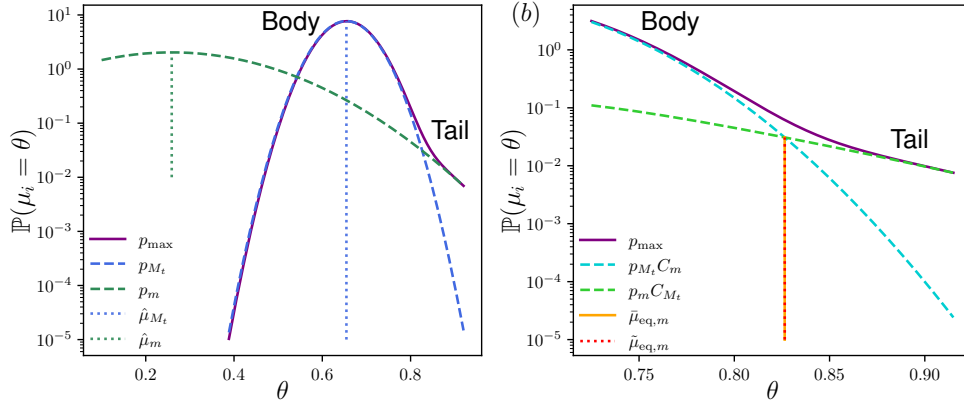


Figure 1: **(a)** Posterior distributions of a two-armed bandit with Gaussian rewards. The dotted lines represent the individual posterior distributions of each arm,  $p_{M_t}$  and  $p_m$ , while the continuous line represents the posterior of the maximum mean reward of all arms,  $p_{\text{max}}$  (Equation (3)). **(b)** Zoom of **(a)** around the point  $\bar{\mu}_{\text{eq},m}$  where both arms have the same posterior probability of being the best one.  $p_{M_t} C_m$  ( $p_m C_{M_t}$ ) is the probability that the maximal value is given by the better (worse) empirical arm, and  $\tilde{\mu}_{\text{eq},m}$  is the approximation to  $\bar{\mu}_{\text{eq},m}$  given in Appendix A.6.

For Gaussian reward distributions, one can derive an analytical expression for  $\tilde{\mu}_{\text{eq}}$  (see Appendix A.2 for details), Equations (7) and (8) can be computed exactly (see Appendix A.4). However, at this stage, even if we have already obtained a closed-form expression for  $S_{\text{max}}$ , it remains too involved to directly compute its exact (discrete) gradient for our decision policy. To finally derive a simplified gradient, we retain only the asymptotic terms of Equations (7) and (8) and of the obtained gradient (see Appendix A.5 for derivation details). Finally, the expression of our approximate difference of gradients of the entropy, whether calculated along a given worse empirical arm  $k$  or along the best empirical arm, reads:

$$\begin{aligned} \Delta_{M_t,k} &= \frac{1}{2} \ln\left(\frac{N_{M_t}}{N_{M_t} + 1}\right) + \frac{1}{2N_{M_t}} \min\left(\frac{1}{2} \sum_m \text{erfc}(\delta \tilde{\mu}_{\text{eq},m}), 1 - \frac{1}{K}\right) \\ &+ Q(N_k^{-1}, \ln(N_{M_t}), \delta \tilde{\mu}_{\text{eq},k}) e^{-\delta \tilde{\mu}_{\text{eq},k}^2} + \sum_m P(N_m^{1/2}, N_{M_t}^{-1}, \ln(N_{M_t}), \delta \tilde{\mu}_{\text{eq},m}) e^{-\delta \tilde{\mu}_{\text{eq},m}^2} \end{aligned} \quad (9)$$

where  $Q$  and  $P$  are polynomials given in Appendix A.6 and  $\delta\tilde{\mu}_{\text{eq},i} = \frac{\sqrt{N_i}(\tilde{\mu}_{\text{eq},i} - \hat{\mu}_i)}{\sqrt{2\sigma^2}}$  are standardized variables with  $\tilde{\mu}_{\text{eq},i}$  given in Appendix A.6. In words,  $\Delta_{M_t,k}$  approximates the difference

$$\Delta_{M_t,k} \approx \mathbb{E}[S_{\max}(t+1) \mid \mathcal{F}_t, a_{t+1} = M_t] - \mathbb{E}[S_{\max}(t+1) \mid \mathcal{F}_t, a_{t+1} = k],$$

which is directly related to greedily maximizing the entropy decrease, described in Equation (5). The decision procedure can be summarized as follows: if  $\Delta_{M_t,k}$  is negative for all  $k \in \mathcal{A}_t$ , the better empirical arm is chosen as it reduces the most the expected value of the approximate entropy. Inversely, if at least one value  $\Delta_{M_t,k}$  is positive, the arm  $k$  maximizing  $\Delta_{M_t,k}$  is chosen.

In conclusion, we have derived an analytical expression for the information available about the maximum expected reward of all arms. We isolated an analytically tractable gradient acting as a decision procedure that eluded previous approximated information derivations [Barbier-Chebbah et al., 2023]. Our scheme leads to an efficient numerical implementation by eliminating numerical integrals, substantially improving the computational speed of information maximization, a crucial challenge for information methods, which is also stressed by Russo and Van Roy [2014] for the IDS algorithm. We now provide the full implementation of AIM and bound its regret in the next section.

### 3.3 Approximate information maximization algorithm

The pseudo-code for the AIM algorithm is presented in Alg. 1 below.

---

#### Algorithm 1: AIM Algorithm for $K$ Gaussian arms

---

```

Draw each arm once, observe reward  $X_t(t)$  and update statistics  $\hat{\mu}_t$ 
for  $t = K + 1$  to  $T$  do                                     // Arm selection
    if  $N_{M_t} \leq N_m$  then  $a_t \leftarrow M_t$ 
    else
         $M_t \leftarrow \operatorname{argmax}_{k \in [K]} \hat{\mu}_k$ 
        Evaluate  $m = \operatorname{argmax}_{k \in \mathcal{A}_t} \Delta_{M_t,k}$  following Equation (9)
        if  $\Delta_{M_t,m} \leq 0$  then  $a_t \leftarrow M_t$ 
        else  $a_t \leftarrow m$ 
        Pull  $a_t$  and observe  $X_t(a_t)$ 
     $\hat{\mu}_{a_t} \leftarrow \frac{\hat{\mu}_{a_t} N_{a_t} + X_t(a_t)}{N_{a_t} + 1}, N_{a_t} \leftarrow N_{a_t} + 1$            // Update statistics

```

---

The best empirical arm is drawn by default if there exists one empirical arm  $m$  that has been drawn more frequently  $N_{M_t} \leq N_m$ . In such a case, both entropy components in Equation (6) are mainly contributed to by  $M_t$ .

## 4 Regret bound

This section provides theoretical guarantees on the performance of AIM. More precisely, Theorem 1 below states that AIM is asymptotically optimal on the multi-armed bandits problem with Gaussian rewards.

**Theorem 1.** *For Gaussian reward distributions with variance  $\sigma^2$ , the regret of AIM satisfies for any mean vector  $\boldsymbol{\mu} \in \mathbb{R}^K$*

$$\limsup_{T \rightarrow \infty} \frac{R(T)}{\ln(T)} \leq \sum_{k, \mu_k < \mu^*} \frac{2\sigma^2}{\mu^* - \mu_k},$$

where  $\mu^* = \max_{k \in [K]} \mu_k$ .

With Gaussian rewards, the asymptotic regret of AIM thus exactly reaches the lower bound of Lai et al. [1985] given by Equation (2). A non-asymptotic version of Theorem 1 is given by Theorem 2 in Appendix B. We briefly sketch the proof idea below and refer to Appendix B for the complete proof.

**Sketch of the proof.** We assume for sake of clarity in this sketch that  $\mu_1 > \mu_k$  for any  $k \geq 2$ . The structure of the proof is similar to the one found in Kaufmann et al. [2012a]. In particular, the first

main step shows that the optimal arm is pulled at least  $\sqrt{t}$  times with high probability. This result holds because otherwise, the contribution of arm 1 to the tail of the distribution would dominate the contribution of other arms in the approximate information. In that case, pulling the first arm would naturally lead to a larger decrease in entropy, which ensures that the optimal arm is always pulled a significant amount of times.

Then, we only need to work in the asymptotic regime where arm 1 is pulled at least  $\sqrt{t}$  times and we aim at bounding the number of pulls on the arm  $k \geq 2$ . Additionally, we restrict ourselves to a large number (in  $\log(T)$ ) of pulls on arm  $k$  and automatically count the pulls before that point in the regret. As a consequence, we can show that with high probability:

$$\hat{\mu}_{M_t} \geq \mu^* - \sqrt{\frac{6\sigma^2 \ln t}{\sqrt{t}}} \quad \text{and} \quad \hat{\mu}_k \leq \mu_k + \varepsilon$$

for some arbitrary  $\varepsilon > 0$ . An important property of entropy is that it approximates the behaviour of the bound of Lai et al. [1985]. More precisely, in the asymptotic regime, the difference of the entropy increments behaves as

$$\Delta_{M_t, k} \approx -\frac{1}{2N_{M_t}} + Q(N_k)e^{-\frac{N_k(\mu_1 - \mu_k)^2}{2\sigma^2}} + \sum_{i \neq M_t} P(N_i)e^{-\frac{N_i(\mu_1 - \mu_i)^2}{2\sigma^2}}, \quad (10)$$

where  $Q_k$  and  $P_i$  are polynomials that also depend on extra variables (see Equation 9). Manipulating these polynomial terms altogether is intricate, but we can still show that if the arm  $k$  is pulled, this means the term  $e^{-\frac{N_k(\mu_1 - \mu_k)^2}{2\sigma^2}}$  somewhat dominates the other exponents in the sum of Equation (10). This then implies that  $N_k$  is of order at most  $\frac{2\sigma^2 \ln T}{(\mu^* - \mu_k)^2}$ , as arm  $k$  is only pulled if  $\Delta_{M_t, k} \geq 0$ .  $\square$

Our policy is deterministic at each time step while displaying a logarithmic regret, showing that intuitions from Russo and Van Roy [2014] of linear regrets for stationary (in the sense they only depend on the posterior distribution) deterministic algorithms was inexact. Moreover, our regret bound is frequentist, in opposition to the Bayesian regret bound obtained for IDS [Russo and Van Roy, 2014]. As a consequence, AIM does not need a well-specified prior: using a uniform prior, as done in our work, is a well-suited choice. Also, the required form of the entropy for the proof is general. The algorithm yields an optimal regret as long as we are guaranteed that the optimal arm is pulled a significant amount of times with high probability and that the asymptotic regime behaves as Equation (10). Hence, Theorem 1 will hold for a large family of entropy approximations [and likely for generalizations to free energies too, as in Masson, 2013] as long as the approximation is accurate enough to not yield trivial behaviors in the short time regime. Additionally, the approximate framework devised here allows fine-tuning the formulas to improve short-time performance all the while ensuring asymptotic optimality by keeping the correct asymptotic terms.

## 5 Experiments

This section investigates the empirical performance of AIM (Alg. 1) on numerical examples. All details of the numerical experiments are given in Appendix D.

We start by considering two arms with Gaussian rewards [Honda and Takemura, 2010] of unit variance and means  $\mu_k$  drawn uniformly from  $[0, 1]$ . Figure 2 compares the Bayesian regret (i.e., the regret averaged over all values of  $(\mu_1, \mu_2)$  in  $[0, 1] \times [0, 1]$ ) of Alg. 1 with the state-of-the-art algorithms UCB-tuned, Thompson sampling, Thompson sampling+, KLUCB++, and MED [Kaufmann et al., 2012b, Pilarski et al., 2021, Cappé et al., 2013, Jin et al., 2022, Honda and Takemura, 2011, Ménard and Garivier, 2017]. We refer to Appendix D.4 for an overview and detailed descriptions of these bandit algorithms. The Bayesian regret of AIM empirically scales as  $\log(T)$ . Its long-time performance matches Thompson sampling, as implied by Theorem 1, while relying on a (conditionally) deterministic decision process. Additionally, AIM outperforms Thompson sampling at both short and intermediate times (see Appendix D.5.3 for finer measurements). AIM particularly outperforms Thompson sampling when the arms are difficult to distinguish due to their mean rewards being close (see examples in Appendix D.5.1 with single instance regret experiments).

AIM yields strong performance in both two-armed Gaussian and 50-armed Gaussian rewards case, as predicted by our theoretical analysis. We now aim to extend our method to other bandit settings.



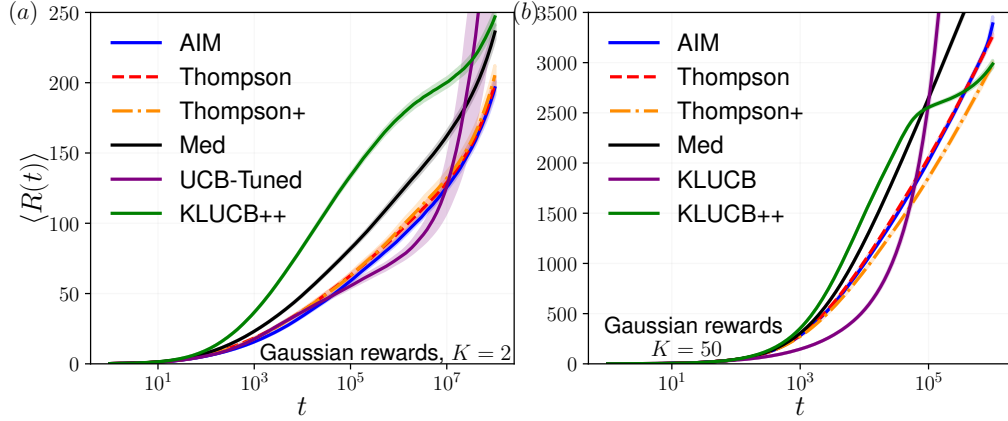


Figure 2: Evolution of the Bayesian regret for (a) 2-armed and (b) 50-armed bandit with Gaussian rewards under a uniform mean prior. Regret is averaged over 8000 for (a) and 2000 runs for (b) Confidence intervals show the standard deviation.

Figure 3 presents the performance of AIM when adapted to Bernoulli rewards [Pilarski et al., 2021] with arm means drawn uniformly in  $[0, 1]$ . This adaptation is described in detail in Section 6 below. The performance of AIM is comparable to Thompson sampling here. Additionally, AIM performs comparably to Thompson sampling for close mean rewards (see Appendix D.5.2). Additionally, for 50 arms with Bernoulli rewards AIM’s short-time efficiency is comparable to Thompson sampling, and it is significantly more efficient at intermediary times while showing the same logarithmic scaling at long times as Thompson sampling.

Hence, our algorithm shows strong empirical performances compared to state-of-the-art baselines for both Bernoulli and Gaussian rewards while providing outstanding effectiveness when facing multiple arms with Bernoulli rewards. Experiments suggest that AIM displays the same typical worst-case regret as Thompson sampling (which is minimax optimal up to  $\sqrt{\ln K}$  for sub-Gaussian rewards), but proving a theoretical bound remains challenging and left for future work. Of note, similar observations are drawn in Appendices D.5.1 and D.5.2 for non-Bayesian versions of the regret, with fixed bandit instances. These observations support the robustness of AIM and its potential for extensions to more complex bandit settings.

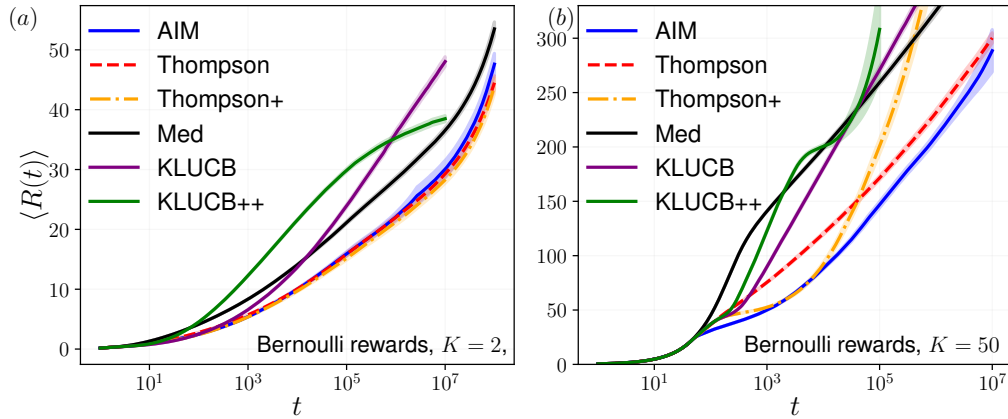


Figure 3: Evolution of the Bayesian regret for (a) 2-armed and (b) 50-armed bandit with Bernoulli rewards under a uniform mean prior. The regret is averaged over 16000 runs for (a) and 2000 runs for (b). Confidence intervals show the standard deviation.

## 6 Extensions

We apply our information maximisation approach to Bernoulli bandits both with two and with many arms, where it shows strong empirical performances (see Figure 3 above). This section describes the extensions of AIM to this case and discusses potential extensions to more general bandit settings.

**Exponential family bandits.** Since Equation (3) explicitly relies on the arms’ posterior distributions, information maximization methods can be directly extended to various reward distributions. In particular, when the reward distributions belong to the exponential family [see Korda et al., 2013, and Appendix C.1 for details on such distributions], an asymptotic and analytical expression of the entropy can be derived for the case of uniform priors (see Appendix C for more details), yielding

$$\tilde{S}_{\max} = \frac{1}{2} \ln(2\pi\tilde{\sigma}_i^2) \left[ 1 - \frac{e^{-N_m \text{KL}(\hat{\theta}_{m_t}, \tilde{\theta}_{\text{eq}})}}{N_m \partial_2 \text{KL}(\hat{\theta}_{m_t}, \tilde{\theta}_{\text{eq}}) \sqrt{2\pi\tilde{\sigma}_m^2}} \right] + \frac{\text{KL}(\hat{\theta}_{m_t}, \tilde{\theta}_{\text{eq}}) e^{-N_m \text{KL}(\hat{\theta}_{m_t}, \tilde{\theta}_{\text{eq}})}}{\partial_2 \text{KL}(\hat{\theta}_{m_t}, \tilde{\theta}_{\text{eq}}) \sqrt{2\pi\tilde{\sigma}_m^2}}. \quad (11)$$

Here  $\text{KL}(\hat{\theta}_i, \tilde{\theta}_{\text{eq}})$  is the Kullback-Leibler divergence between the reward distribution parameterized by  $\hat{\theta}_i$  and  $\tilde{\theta}_{\text{eq}}$  where  $\theta$  is the family parameter, and  $\partial_2 \text{KL}$  denotes its derivative w.r.t. the second variable. All the steps leading to Equations (7) and (8) in Section 3 are not specific to Gaussian rewards. The main difference lies in their asymptotic simplifications obtained afterwards with Laplace’s method. Our implementation of AIM to Bernoulli rewards (a specific case of the exponential family) with Equation (11) shows comparable performance to state-of-the-art algorithms (see Figure 2), supporting its adaptability to general settings. We believe that AIM should be optimal for all exponential family reward distributions and general prior distributions and that similar proof techniques can be used (see Appendix C for a detailed discussion). However, significant work still remains to ensure that the asymptotic regime, where all arms have been sufficiently drawn, is reached for any reward distribution and will be addressed in future work.

**Other bandit settings.** Here, we provide a quick overview of several other bandit settings for which approximate information maximization, adapted to the specific bandit problem, should provide efficient algorithms. First, we emphasize that AIM’s partitioning between body and tail components remains relevant even when dealing with heavy-tailed [Lee et al., 2023] or non-parametric reward distributions [Baudry et al., 2020]. It should thus be able to provide strong guarantees in these settings, similarly to Thompson sampling. Secondly, let us stress that information can also be quantified for unpulled arms, which may prove crucial when facing large numbers of arms. The agent could quantify the information of the “reservoir” of unpulled arms to anticipate the information gained from exploring these unpulled arms. Additionally, if the agent has access to the remaining time, it can not only evaluate the expected information gain when pulling an arm for a single round but also evaluate the information gain of multiple pulls of the same arm. We believe that such a consideration might be pivotal when facing many arms, since the limited amount of time does not allow to pull all the arms [Bayati et al., 2020] sufficiently. Thirdly, in linear bandits, where arms are correlated with each other [Li et al., 2010], AIM will be efficient because pulling a specific direction provides information on correlated directions, the shared information gain could be leveraged by information-based methods to yield strong performances. Finally, we could consider pure exploration problems [Bubeck et al., 2011, Locatelli et al., 2016, Kalyanakrishnan et al., 2012] where the agent’s goal is directly linked to an information gain, thus making the information maximization principle an inherent candidate when a suitable entropy is derived from the underlying bandit structure and problem objective.

A last advantage of AIM lies in its possible extension to multiple constraints that would be introduced using Lagrange multipliers (or borrowed from physics reasoning by defining free energy), further improving its adaptability to various settings and specific requirements.

## 7 Conclusion

This paper introduces a new algorithm class, Approximate Information Maximization (AIM), which leverages approximate information maximization of the whole bandit system to achieve optimal regret performances. This approach builds on the entropy of the posterior of the arms’ maximal mean, from which we extract a simplified and analytical functional at the core of the decision scheme. It enables easily tunable and tractable algorithms, which we prove to be optimal for multi-armed

Gaussian bandits. Numerical experiments for Bernoulli rewards with two or several arms emphasize the robustness and efficiency of AIM. An additional strength of AIM lies in its efficiency at short times and when the arms have close mean rewards where it outperforms existing state-of-the art. Further research should focus on adjusting the information maximization framework to more complex bandit settings, including many-armed bandits, linear bandits and thresholding bandits, where appropriately selected information measures can efficiently apprehend the games' structure and correlations.

## References

- Shipra Agrawal and Navin Goyal. Thompson Sampling for Contextual Bandits with Linear Payoffs. In *Proceedings of the 30th International Conference on Machine Learning*, pages 127–135. PMLR, May 2013.
- P. Auer. Using upper confidence bounds for online learning. In *Proceedings 41st Annual Symposium on Foundations of Computer Science*, pages 270–279, Redondo Beach, CA, USA, 2000. IEEE Comput. Soc. ISBN 978-0-7695-0850-4. doi: 10.1109/SFCS.2000.892116.
- Alex Barbier-Chebbah, Christian L. Vestergaard, and Jean-Baptiste Masson. Approximate information for efficient exploration-exploitation strategies, July 2023.
- Dorian Baudry, Emilie Kaufmann, and Odalric-Ambrym Maillard. Sub-sampling for efficient non-parametric bandit exploration. In *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS’20*, pages 5468–5478, Red Hook, NY, USA, December 2020. Curran Associates Inc. ISBN 978-1-71382-954-6.
- Mohsen Bayati, Nima Hamidi, Ramesh Johari, and Khashayar Khosravi. Unreasonable effectiveness of greedy algorithms in multi-armed bandit with many arms. *Advances in Neural Information Processing Systems*, 33:1713–1723, 2020.
- Guy Bresler, George H Chen, and Devavrat Shah. A latent source model for online collaborative filtering. *Advances in neural information processing systems*, 27, 2014.
- Sébastien Bubeck, Rémi Munos, and Gilles Stoltz. Pure exploration in finitely-armed and continuous-armed bandits. *Theoretical Computer Science*, 412(19):1832–1852, April 2011. ISSN 0304-3975. doi: 10.1016/j.tcs.2010.12.059.
- Olivier Cappé, Aurélien Garivier, Odalric-Ambrym Maillard, Rémi Munos, and Gilles Stoltz. Kullback–Leibler upper confidence bounds for optimal sequential allocation. *The Annals of Statistics*, 41(3):1516–1541, June 2013. ISSN 0090-5364, 2168-8966. doi: 10.1214/13-AOS1119.
- Ring T. Cardé. Navigation Along Windborne Plumes of Pheromone and Resource-Linked Odors. *Annual Review of Entomology*, 66(1):317–336, 2021. doi: 10.1146/annurev-ento-011019-024932.
- Arnoud V. Den Boer. Dynamic pricing and learning: historical origins, current research, and new directions. *Surveys in operations research and management science*, 20(1):1–18, 2015.
- Aurélien Garivier. Informational confidence bounds for self-normalized averages and applications. In *2013 IEEE Information Theory Workshop (ITW)*, pages 1–5. IEEE, 2013.
- Aurélien Garivier and Olivier Cappé. The kl-ucb algorithm for bounded stochastic bandits and beyond. In *Proceedings of the 24th annual conference on learning theory*, pages 359–376. JMLR Workshop and Conference Proceedings, 2011.
- Moritz Helias and David Dahmen. *Statistical Field Theory for Neural Networks*, volume 970 of *Lecture Notes in Physics*. Springer International Publishing, Cham, 2020. ISBN 978-3-030-46443-1 978-3-030-46444-8. doi: 10.1007/978-3-030-46444-8.
- José Miguel Hernández-Lobato, Michael A. Gelbart, Matthew W. Hoffman, Ryan P. Adams, and Zoubin Ghahramani. Predictive entropy search for Bayesian optimization with unknown constraints. In *Proceedings of the 32nd International Conference on Machine Learning - Volume 37, ICML’15*, pages 1699–1707, Lille, France, July 2015. JMLR.org.
- Junya Honda and Akimichi Takemura. An Asymptotically Optimal Bandit Algorithm for Bounded Support Models. In Adam Tauman Kalai and Mehryar Mohri, editors, *COLT 2010 - The 23rd Conference on Learning Theory, Haifa, Israel, June 27-29, 2010*, pages 67–79. Omnipress, 2010.
- Junya Honda and Akimichi Takemura. An asymptotically optimal policy for finite support models in the multiarmed bandit problem. *Mach Learn*, 85(3):361–391, December 2011. ISSN 1573-0565. doi: 10.1007/s10994-011-5257-4.

- Tianyuan Jin, Pan Xu, Xiaokui Xiao, and Anima Anandkumar. Finite-Time Regret of Thompson Sampling Algorithms for Exponential Family Multi-Armed Bandits. *Advances in Neural Information Processing Systems*, 35:38475–38487, December 2022.
- Shivaram Kalyanakrishnan, Ambuj Tewari, Peter Auer, and Peter Stone. PAC Subset Selection in Stochastic Multi-armed Bandits. *Proceedings of the 29th International Conference on Machine Learning, ICML 2012*, 1, January 2012.
- Emilie Kaufmann, Nathaniel Korda, and Rémi Munos. Thompson sampling: An asymptotically optimal finite-time analysis. In *International conference on algorithmic learning theory*, pages 199–213. Springer, 2012a.
- Emilie Kaufmann, Nathaniel Korda, and Rémi Munos. Thompson Sampling: An Asymptotically Optimal Finite-Time Analysis. In Nader H. Bshouty, Gilles Stoltz, Nicolas Vayatis, and Thomas Zeugmann, editors, *Algorithmic Learning Theory*, Lecture Notes in Computer Science, pages 199–213, Berlin, Heidelberg, 2012b. Springer. ISBN 978-3-642-34106-9. doi: 10.1007/978-3-642-34106-9\_18.
- Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of best-arm identification in multi-armed bandit models. *J. Mach. Learn. Res.*, 17(1):1–42, January 2016. ISSN 1532-4435.
- Nathaniel Korda, Emilie Kaufmann, and Remi Munos. Thompson sampling for 1-dimensional exponential family bandits. In *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 1, NIPS’13*, pages 1448–1456, Red Hook, NY, USA, December 2013. Curran Associates Inc.
- Tze Leung Lai, Herbert Robbins, et al. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4–22, 1985.
- Jongyeong Lee, Junya Honda, Chao-Kai Chiang, and Masashi Sugiyama. Optimality of thompson sampling with noninformative priors for pareto bandits. *arXiv preprint arXiv:2302.01544*, 2023.
- Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pages 661–670, 2010.
- Andrea Locatelli, Maurilio Gutzeit, and Alexandra Carpentier. An optimal algorithm for the Thresholding Bandit Problem. In *Proceedings of The 33rd International Conference on Machine Learning*, pages 1690–1698. PMLR, June 2016.
- Dominique Martinez, Lotfi Arhidi, Elodie Demondion, Jean-Baptiste Masson, and Philippe Lucas. Using Insect Electroantennogram Sensors on Autonomous Robots for Olfactory Searches. *JoVE (Journal of Visualized Experiments)*, (90):e51704, August 2014. ISSN 1940-087X. doi: 10.3791/51704.
- Jean-Baptiste Masson. Olfactory searches with limited space perception. *Proceedings of the National Academy of Sciences*, 110(28):11261–11266, July 2013. doi: 10.1073/pnas.1221091110.
- Pierre Ménard and Aurélien Garivier. A minimax and asymptotically optimal algorithm for stochastic bandits. In *Proceedings of the 28th International Conference on Algorithmic Learning Theory*, pages 223–237. PMLR, October 2017.
- John Murlis, Joseph S. Elkinton, and Ring T. Cardé. Odor Plumes and How Insects Use Them. *Annual Review of Entomology*, 37(1):505–532, January 1992. doi: 10.1146/annurev.en.37.010192.002445.
- Edward W. Ng and Murray Geller. A table of integrals of the Error functions. *J. RES. NATL. BUR. STAN. SECT. B. MATH. SCI.*, 73B(1):1, January 1969. ISSN 0098-8979. doi: 10.6028/jres.073B.001.
- Thomas Parr, Giovanni Pezzulo, and Karl J. Friston. *Active Inference: The Free Energy Principle in Mind, Brain, and Behavior*. The MIT Press, March 2022. ISBN 978-0-262-36997-8. doi: 10.7551/mitpress/12441.001.0001.

- Sebastian Pilarski, Slawomir Pilarski, and Dániel Varró. Optimal Policy for Bernoulli Bandits: Computation and Algorithm Gauge. *IEEE Transactions on Artificial Intelligence*, 2(1):2–17, February 2021. ISSN 2691-4581. doi: 10.1109/TAI.2021.3074122.
- Gautam Reddy, Antonio Celani, and Massimo Vergassola. Infomax Strategies for an Optimal Balance Between Exploration and Exploitation. *Journal of Statistical Physics*, 163(6):1454–1476, April 2016. doi: 10.1007/s10955-016-1521-0.
- Gautam Reddy, Venkatesh N. Murthy, and Massimo Vergassola. Olfactory Sensing and Navigation in Turbulent Environments. *Annual Review of Condensed Matter Physics*, 13:191–213, March 2022. doi: 10.1146/annurev-conmatphys-031720-032754.
- Daniel Russo and Benjamin Van Roy. Learning to Optimize via Information-Directed Sampling. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.
- Ilya O. Ryzhov, Warren B. Powell, and Peter I. Frazier. The Knowledge Gradient Algorithm for a General Class of Online Learning Problems. *Operations Research*, 60(1):180–195, February 2012. ISSN 0030-364X. doi: 10.1287/opre.1110.0999.
- David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, January 2016. ISSN 1476-4687. doi: 10.1038/nature16961.
- 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.
- Naftali Tishby, Fernando C. Pereira, and William Bialek. The information bottleneck method, April 2000.
- Massimo Vergassola, Emmanuel Villermaux, and Boris I. Shraiman. ‘Infotaxis’ as a strategy for searching without gradients. *Nature*, 445(7126):406–409, January 2007. ISSN 1476-4687. doi: 10.1038/nature05464.
- Siqi Zhang, Dominique Martinez, and Jean-Baptiste Masson. Multi-Robot Searching with Sparse Binary Cues and Limited Space Perception. *Frontiers in Robotics and AI*, 2, 2015. ISSN 2296-9144.

# Appendix

## Table of Contents

---

<b>A Towards an analytical approximation of the entropy</b>	<b>14</b>
A.1 The partitioning approximation . . . . .	14
A.2 Asymptotics of the intersection point . . . . .	16
A.3 Closed-form expressions for the main mode’s contribution . . . . .	16
A.4 Closed form and asymptotic expression for the tail’s entropy . . . . .	18
A.5 Derivation of the increment for the closed-form expression of entropy . . . . .	18
A.6 Final expression for the increment comparison . . . . .	21
<b>B Proof of Theorem 1</b>	<b>21</b>
B.1 Auxiliary Lemmas . . . . .	23
<b>C Generalization of the information maximization approximation</b>	<b>27</b>
C.1 Asymptotic expression for exponential family rewards . . . . .	28
C.2 The partitioning approximation . . . . .	29
C.3 Asymptotic intersection point . . . . .	29
C.4 Generalization of the main mode’s contribution . . . . .	29
C.5 Generalized expression for the entropy tail . . . . .	30
C.6 Generalized form of the entropy approximation . . . . .	30
C.7 Derivation of the increment for the closed-form expression of entropy . . . . .	31
<b>D Numerical experiments</b>	<b>31</b>
D.1 Numerical settings . . . . .	31
D.2 AIM implementation details . . . . .	32
D.3 Information maximization approximation for Bernoulli rewards with more than two arms . . . . .	34
D.4 Overview of baseline bandit algorithms . . . . .	35
D.5 Additional experiments . . . . .	36

---

## A Towards an analytical approximation of the entropy

In this section, we recapitulate all the steps leading to the analytical expression constitutive of our AIM algorithm. We stress, that it involves exact derivations but also simplifications to considerably simplify the final form of AIM. Therefore, alternative approaches could lead to a slightly different version of AIM. However, our chosen method retains the essential features which emerges in the asymptotic regime while providing a simple version of AIM.

### A.1 The partitioning approximation

We start by commenting on the partition scheme and the approximations leading to the following body/tail expressions. We first recall the expression for  $p_{\max}(\theta)$ , with the arms ordered along  $M_t$ ,

$$p_{\max}(\theta) = \left( p_{M_t}(\theta) \prod_m^{A_t} C_m(\theta) + \sum_m^{A_t} C_{M_t}(\theta) p_m(\theta) \prod_{j \neq m}^{A_t} C_j(\theta) \right). \quad (12)$$

where we remind that  $C_m(\theta)$  is the cumulative posterior probability of the mean of arm  $m$

Because of its dependency along all arms, there is no unique dominant term in Equation (12), and distinct regimes emerge depending on  $\theta$  and the state of the game. Then, we assume to isolate distinct regimes contributing asymptotically to the entropy while significantly simplifying them. It will considerably simplify the derivation of an analytical expression for the body/tail components in the next section. The next paragraphs will then present heuristic arguments justifying our simplification scheme.

We start by rewriting the exact entropy expression isolating  $M_t$ :

$$\begin{aligned}
S_{\max} &= - \int_{\Theta} p_{M_t}(\theta) \prod_m^{\mathcal{A}_t} C_j(\theta) \ln \left( p_{M_t}(\theta) \prod_j^{\mathcal{A}_t} C_j(\theta) + \sum_m^{\mathcal{A}_t} p_m(\theta) C_{M_t}(\theta) \prod_{j \neq m}^{\mathcal{A}_t} C_j(\theta) \right) d\theta \\
&\quad - \sum_m^{\mathcal{A}_t} \int_{\Theta} p_i(\theta) C_{M_t}(\theta) \prod_{j \neq m}^{\mathcal{A}_t} C_j(\theta) \ln \left( p_{M_t}(\theta) \prod_j^{\mathcal{A}_t} C_j(\theta) + \sum_m^{\mathcal{A}_t} p_m(\theta) C_{M_t}(\theta) \prod_{j \neq m}^{\mathcal{A}_t} C_j(\theta) \right) d\theta.
\end{aligned} \tag{13}$$

Let us briefly comment on the different contributions to Equation (13). We aim to keep the leading orders of  $p_{\max}(\theta)$  when  $N_{M_t} \gg N_m \gg 1$  and  $\hat{\mu}_{M_t} > \hat{\mu}_m$  for all  $m$  in the set of current worse empirical arms  $\mathcal{A}_t$ . Here, the posterior distributions are assumed uni-modal. The first term is the leading order in the vicinity of the mode of  $\hat{\mu}_{M_t}$ . Also, since  $N_{M_t} > N_m$ ,  $p_{M_t}(\theta)$  is more concentrated than all  $p_m(\theta)$ , resulting in the dominance of the second term in the distribution's tail (*i.e.*, for high rewards).

We now decompose the entropy in the body/tail components defined in the main text, the first term of Equation (13) will form the body component, and the sum over all worse empirical arms will compose the tail. We now define  $\bar{\mu}_{\text{eq},m}$  the intersection point associated to the arm  $m$  verifying  $C_m(\bar{\mu}_{\text{eq},m})p_{M_t}(\bar{\mu}_{\text{eq},m}) = C_{M_t}(\bar{\mu}_{\text{eq},m})p_m(\bar{\mu}_{\text{eq},m})$ . Then, in the asymptotic regime,  $\bar{\mu}_{\text{eq},m}$  will verify  $p_m(\theta) \gg p_{M_t}(\theta)$  for  $\theta > \bar{\mu}_{\text{eq},m}$  and  $p_m(\theta) \ll p_{M_t}(\theta)$  for  $\theta < \bar{\mu}_{\text{eq},m}$ . Again, we will assume to neglect the transition regime where  $\bar{\mu}_{\text{eq},m} \sim \theta$  where both distributions are of the same order because it is narrow (in the asymptotic regime) and has very little influence on the total value of the entropy.

To get the body component, we consider the first term of Equation (13). We neglect all the inner terms inside the logarithm which is then dominated by  $M_t$ . Next, by noticing that  $C_i(\theta) \approx 1$  is in the vicinity of  $\hat{\mu}_{M_t}$ , we make a first-order expansion of the remaining product along all the worse empirical arms. Since the inner term of the body component is negligible for  $\theta > \min(\{\bar{\mu}_{\text{eq},m}, m \in \mathcal{A}_t\})$  (because of its dependency along  $p_{M_t}$ ), we ignore that our simplification is no more valid in this specific regime without loss of consistency. Taken together we obtain the body expression of the main text:

$$\tilde{S}_{\text{body}} = - \int_{\Theta} \left( 1 - \sum_m^{\mathcal{A}_t} [1 - C_m(\theta)] \right) p_{M_t}(\theta) \ln p_{M_t}(\theta) d\theta. \tag{14}$$

Then, we consider the additional terms (each denoted as  $m$ ) in Equation (13). First, each term of the sum is negligible to the first one for  $\theta < \bar{\mu}_{\text{eq},m}$ , we then only keep the upper part of the integral where  $\theta > \bar{\mu}_{\text{eq},m}$ . Because  $N_m \gg 1$  and  $\theta > \min(\{\bar{\mu}_{\text{eq},m}, i \neq M_t\}) > \hat{\mu}_{M_t} > \hat{\mu}_j$ , we approximate all the cumulative by one. Finally, to get a simplified expression for the increment, we assume to neglect all the posterior distributions except for  $p_i(\theta)$  inside the logarithm of the  $i$ -th term and approximates  $\bar{\mu}_{\text{eq},m}$  (see next section) which leads to the tail expression:

$$\tilde{S}_{\text{tail}} = - \sum_m^{\mathcal{A}_t} \int_{\bar{\mu}_{\text{eq},m}}^{\mu_{\text{sup}}} p_m(\theta) \ln p_m(\theta) d\theta. \tag{15}$$

note that some of these posterior distributions ( $j \neq m, M_t$ ) are not negligible compared to  $p_m(\theta)$  at a given  $\theta$ . However, this cross-information between current suboptimal arms is asymptotically negligible regarding the decision procedure (which largely resumes as balancing exploiting the best empirical solution compared to exploring worse empirical arms) while unnecessarily complicating the increment evaluation.

Finally, we obtain the full expression of the entropy approximation:

$$\tilde{S}_{\max} = - \int_{\Theta} \left( 1 - \sum_m^{\mathcal{A}_t} [1 - C_m(\theta)] \right) p_{M_t}(\theta) \ln p_{M_t}(\theta) d\theta - \sum_m^{\mathcal{A}_t} \int_{\bar{\mu}_{\text{eq},m}}^{\mu_{\text{sup}}} p_i(\theta) \ln p_i(\theta) d\theta. \tag{16}$$



## A.2 Asymptotics of the intersection point

In this section, we derive the asymptotic expression of the intersection point (defined above as  $\tilde{\mu}_{\text{eq},m}$ ) where the distributions  $C_m(\tilde{\mu}_{\text{eq},m})p_{M_t}(\tilde{\mu}_{\text{eq},m})$  and  $C_{M_t}(\tilde{\mu}_{\text{eq},m})p_m(\tilde{\mu}_{\text{eq},m})$  intersect (at their highest value if they intersect more than once). Here, we consider Gaussian rewards and the intersection between  $M_t$  and a given worse empirical arm denoted  $m$ . The exact equation verified by the intersection point  $\bar{\mu}_{\text{eq},m}$  is:

$$\frac{\sqrt{N_{M_t}} e^{-\frac{N_{M_t}(\bar{\mu}_{\text{eq},m} - \hat{\mu}_{M_t})^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}} \frac{1}{2} \left[ 1 + \operatorname{erf} \left( \frac{\sqrt{N_m}(\bar{\mu}_{\text{eq},m} - \hat{\mu}_m)}{\sqrt{2\sigma^2}} \right) \right] = \frac{\sqrt{N_m} e^{-\frac{N_m(\bar{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}} \frac{1}{2} \left[ 1 + \operatorname{erf} \left( \frac{\sqrt{N_{M_t}}(\bar{\mu}_{\text{eq},m} - \hat{\mu}_{M_t})}{\sqrt{2\sigma^2}} \right) \right]. \quad (17)$$

Taking the logarithm of Equation (17) and normalizing the last term leads to:

$$\frac{N_m(\bar{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{2\sigma^2} - \frac{N_{M_t}(\bar{\mu}_{\text{eq},m} - \hat{\mu}_{M_t})^2}{2\sigma^2} + \frac{1}{2} \ln \frac{N_{M_t}}{N_m} + \ln \left[ \frac{1 + \operatorname{erf} \left( \frac{\sqrt{N_m}(\bar{\mu}_{\text{eq},m} - \hat{\mu}_m)}{\sqrt{2\sigma^2}} \right)}{1 + \operatorname{erf} \left( \frac{\sqrt{N_{M_t}}(\bar{\mu}_{\text{eq},m} - \hat{\mu}_{M_t})}{\sqrt{2\sigma^2}} \right)} \right] = 0. \quad (18)$$

The distributions are uni-modal, and assuming that  $\hat{\mu}_{M_t} > \hat{\mu}_m$ ,  $N_{M_t} > N_m$  and recalling that  $\bar{\mu}_{\text{eq},m}$  is the highest intersection, we get that  $\bar{\mu}_{\text{eq},m} > \hat{\mu}_{M_t} > \hat{\mu}_m$ . Both error functions are then bounded in  $[0, 1]$ , making the last term bounded as well. We then approximate  $\bar{\mu}_{\text{eq},m}$  with  $\tilde{\mu}_{\text{eq},m}$  by neglecting the last term, which leads to the following solution:

$$\tilde{\mu}_{\text{eq},m} = \hat{\mu}_{M_t} + \frac{N_m(\hat{\mu}_{M_t} - \hat{\mu}_m)}{N_{M_t} - N_m} + \sqrt{\frac{N_{M_t}N_m}{(N_{M_t} - N_m)^2}(\hat{\mu}_{M_t} - \hat{\mu}_m)^2 + \frac{\sigma^2}{N_{M_t} - N_m} \ln \left( \frac{N_{M_t}}{N_m} \right)}. \quad (19)$$

Note that Equation (19) relies on both  $\hat{\mu}_{M_t} > \hat{\mu}_m$  and  $N_{M_t} > N_m$ . For  $N_{M_t} \leq N_m$ , even if the above  $\tilde{\mu}_{\text{eq},m}$  can be computed, it does not quantify the tail contribution. As a matter of fact, for  $N_{M_t} \leq N_m$ , the tail is always dominated by  $p_{M_t}$ , which means that it has already been included in the main mode  $\tilde{S}_{\text{body}}$ . Then, in this specific configuration, we take the contribution of the arm  $m$  to  $\tilde{S}_{\text{tail}}$  equals to 0 (in other words  $\tilde{\mu}_{\text{eq},m} = \mu_{\text{sup}}$ ).

## A.3 Closed-form expressions for the main mode's contribution

Here, we derive the  $\tilde{S}_{\text{body}}$  expression given in the main text for Gaussian rewards distribution. Inserting the Gaussian form of the posterior into Equation (8) gives:

$$\tilde{S}_{\text{body}} = - \int_{-\infty}^{+\infty} \frac{\sqrt{N_{M_t}} e^{-\frac{N_{M_t}(\theta - \hat{\mu}_{M_t})^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}} \left( -\frac{1}{2} \ln \left( \frac{2\pi\sigma^2}{N_{M_t}} \right) - \frac{N_{M_t}(\theta - \hat{\mu}_{M_t})^2}{2\sigma^2} \right) \times \left( 1 - \sum_m \frac{\mathcal{A}_t}{2} \left[ 1 - \operatorname{erf} \left( \frac{\sqrt{N_m}(\theta - \hat{\mu}_m)}{\sqrt{2\sigma^2}} \right) \right] \right) d\theta, \quad (20)$$

We integrate the constant part of the first term, denoted  $T_1$  by the use of the following identity [Ng and Geller, 1969]:

$$\int_{-\infty}^{\infty} \left[ 1 + \operatorname{erf} \left( \frac{\theta - \theta_1}{\sqrt{2V_1}} \right) \right] \frac{e^{-\frac{(\theta - \theta_2)^2}{2V_2}}}{\sqrt{2\pi V_2}} d\theta = \left[ 1 + \operatorname{erf} \left( \frac{\theta_2 - \theta_1}{\sqrt{2}\sqrt{V_2 + V_1}} \right) \right], \quad (21)$$

which leads to

$$\begin{aligned}
T_1 &= \frac{1}{2} \ln \left( \frac{2\pi\sigma^2}{N_{M_t}} \right) \int_{-\infty}^{+\infty} \frac{\sqrt{N_{M_t}} e^{-\frac{N_{M_t}(\theta - \hat{\mu}_{M_t})^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}} \left[ 1 - \frac{1}{2} \sum_m^{A_t} 1 - \operatorname{erf} \left( \frac{\sqrt{N_m}(\theta - \hat{\mu}_m)}{\sqrt{2\sigma^2}} \right) \right] d\theta \\
&= \frac{1}{2} \ln \left( \frac{2\pi\sigma^2}{N_{M_t}} \right) \left( 1 - \sum_m^{A_t} \frac{1}{2} \left[ 1 - \operatorname{erf} \left( \frac{\sqrt{N_m}(\hat{\mu}_{M_t} - \hat{\mu}_m)}{\sqrt{2\sigma^2(\frac{1}{N_{M_t}} + \frac{1}{N_m})}} \right) \right] \right) \\
&= \frac{1}{2} \ln \left( \frac{2\pi\sigma^2}{N_{M_t}} \right) \left( 1 - \sum_m^{A_t} \frac{1}{2} \operatorname{erfc} \left[ \frac{\sqrt{N_m}(\hat{\mu}_{M_t} - \hat{\mu}_m)}{\sqrt{2\sigma^2(\frac{1}{N_{M_t}} + \frac{1}{N_m})}} \right] \right)
\end{aligned} \tag{22}$$

Next, we separate the second term in two parts  $T_{2,1}$  and  $T_{2,2}$ , first :

$$T_{2,1} = \int_{-\infty}^{+\infty} \frac{\sqrt{N_{M_t}} e^{-\frac{N_{M_t}(\theta - \hat{\mu}_{M_t})^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}} \frac{N_{M_t}(\theta - \hat{\mu}_{M_t})^2}{2\sigma^2} d\theta = \frac{1}{2}. \tag{23}$$

Then, we integrate by parts the remaining term  $T_{2,2}$  to obtain:

$$\begin{aligned}
T_{2,2} &= - \sum_m^{A_t} \int_{-\infty}^{\infty} \frac{N_{M_t}^{3/2}(\theta - \hat{\mu}_{M_t})^2}{4\sigma^2} \frac{e^{-\frac{N_{M_t}(\theta - \hat{\mu}_{M_t})^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}} \left[ 1 - \operatorname{erf} \left( \frac{\sqrt{N_m}(\theta - \hat{\mu}_m)}{\sqrt{2\sigma^2}} \right) \right] d\theta \\
&= - \sum_m^{A_t} \int_{-\infty}^{\infty} \frac{1}{4} \frac{\sqrt{N_{M_t}} e^{-\frac{N_{M_t}(\theta - \hat{\mu}_{M_t})^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}} \left[ 1 - \operatorname{erf} \left( \frac{\sqrt{N_m}(\theta - \hat{\mu}_m)}{\sqrt{2\sigma^2}} \right) \right] \\
&\quad - \frac{(\theta - \hat{\mu}_{M_t}) \sqrt{N_{M_t} N_m} e^{-\frac{N_{M_t}(\theta - \hat{\mu}_{M_t})^2}{2\sigma^2} - \frac{N_m(\theta - \hat{\mu}_m)^2}{2\sigma^2}}}{2 \cdot 2\pi\sigma^2} d\theta \\
&= - \sum_{M_t}^{A_t} \frac{1}{4} \operatorname{erfc} \left( \frac{\hat{\mu}_{M_t} - \hat{\mu}_m}{\sqrt{2\sigma^2(\frac{1}{N_{M_t}} + \frac{1}{N_m})}} \right) - \frac{(\hat{\mu}_{M_t} - \hat{\mu}_m)\sigma^2}{2N_{M_t}\sqrt{2\pi}(\frac{\sigma^2}{N_{M_t}} + \frac{\sigma^2}{N_m})^{3/2}} e^{-\frac{(\hat{\mu}_{M_t} - \hat{\mu}_m)^2}{2\sigma^2(\frac{1}{N_{M_t}} + \frac{1}{N_m})}},
\end{aligned} \tag{24}$$

where we also rely on the identity of Equation (21).

Combining Equations (22) to (24) leads to the analytical expression of the body component.

$$\begin{aligned}
\tilde{S}_{\text{body}} &= \frac{1}{2} \ln \left( \frac{2\pi\sigma^2 e}{N_{M_t}} \right) \left[ 1 - \frac{1}{2} \sum_m^{A_t} \operatorname{erfc} \left( \frac{\sqrt{N_m N_{M_t}}(\hat{\mu}_{M_t} - \hat{\mu}_m)}{\sqrt{2\sigma^2(N_m + N_{M_t})}} \right) \right] \\
&\quad - \sum_m^{A_t} \frac{\sqrt{N_{M_t} N_m}^{3/2} (\hat{\mu}_{M_t} - \hat{\mu}_m)}{2\sigma\sqrt{2\pi}(N_{M_t} + N_m)^{3/2}} e^{-\frac{N_m N_{M_t} (\hat{\mu}_{M_t} - \hat{\mu}_m)^2}{2\sigma^2(N_m + N_{M_t})}}.
\end{aligned} \tag{25}$$

To finally get an asymptotic and simplified expression of the body component, we neglect the second term. Then, since  $\hat{\mu}_{\text{eq},m} \xrightarrow{N_{M_t} \rightarrow \infty} \hat{\mu}_{M_t}$  and  $N_m \ll N_{M_t}$  asymptotically, we approximate the first term as:

$$\tilde{S}_b = \frac{1}{2} \ln \left( \frac{2\pi\sigma^2 e}{N_{M_t}} \right) \left[ 1 - \frac{1}{2} \sum_m^{A_t} \operatorname{erfc} \left( \frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{\sqrt{2\sigma^2}} \right) \right]. \tag{26}$$

This last approximation will enable to provide an analytically tractable gradient without altering the asymptotic behavior expected at large times for the entropy measure.

#### A.4 Closed form and asymptotic expression for the tail's entropy

The contribution from the tail can be derived exactly and reads:

$$\begin{aligned}\tilde{S}_{\text{tail}} &= \sum_m^{A_t} \int_{\tilde{\mu}_{\text{eq},m}}^{\infty} \frac{\sqrt{N_m} e^{-\frac{N_m(\theta - \hat{\mu}_m)^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}} \left[ \frac{1}{2} \ln\left(\frac{2\pi\sigma^2}{N_m}\right) + \frac{N_m(\theta - \hat{\mu}_m)^2}{2\sigma^2} \right] d\theta \\ &= \sum_m^{A_t} \frac{1}{4} \ln\left(\frac{2\pi\sigma^2 e}{N_m}\right) \operatorname{erfc}\left(\frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{\sqrt{2\sigma^2}}\right) \\ &\quad + \frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{2\sqrt{2\pi\sigma^2}} e^{-\frac{N_m(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{2\sigma^2}}.\end{aligned}\quad (27)$$

To get a simplified analytical expression of the tail component, we only keep the second term since it dominates the others asymptotically,

$$\tilde{S}_t = \sum_m^{A_t} \frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{2\sqrt{2\pi\sigma^2}} e^{-\frac{N_m(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{2\sigma^2}}. \quad (28)$$

Taken altogether, Equations (26) and (28) lead to the desired simplified approximation of the entropy:

$$\tilde{S}_m = \tilde{S}_b + \tilde{S}_t. \quad (29)$$

#### A.5 Derivation of the increment for the closed-form expression of entropy

Since Equations (26) and (28) exhibit simple closed-form expressions, it becomes possible to derive an explicit expression of its expected increment. Here, we again consider continuous Gaussian reward distributions.

We start by deriving the increment along the better empirical arm,  $\Delta_{M_t} \tilde{S}_m$ . The posterior of the reward obtained at time  $t + 1$  is approximated as a Gaussian of variance  $\sigma^2$  and centred around  $\hat{\mu}_{M_t}$ , leading to:

$$\Delta_{M_t} \tilde{S}_m = \int_{-\infty}^{\infty} \frac{e^{-\frac{\mu^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}} \left[ \tilde{S}_m\left(\hat{\mu}_{M_t} + \frac{\mu}{N_{M_t+1}}, N_{M_t} + 1, \dots\right) - \tilde{S}_m\left(\hat{\mu}_{M_t}, N_{M_t}, \dots\right) \right] d\mu. \quad (30)$$

where the dots runs over all the worse empirical arms variables remaining constant when the best empirical arm is drawn at time  $t + 1$ .

For the sake of simplicity, we neglect the variations of all the subdominant terms inside all  $\tilde{\mu}_{\text{eq},m}$  meaning we approximate them as  $\tilde{\mu}_{\text{eq},m}(\hat{\mu}_{M_t} + \frac{\mu}{N_{M_t+1}}, N_{M_t} + 1, \hat{\mu}_m, N_m) \approx \tilde{\mu}_{\text{eq},m}(\hat{\mu}_{M_t}, N_{M_t}, \hat{\mu}_m, N_m) + \frac{\mu}{N_{M_t+1}}$ , after observing a reward  $\mu$  when pulling the arm  $M_t$  for the  $(N_{M_t} + 1)$ th time.

By use of the identity Equation (21), the gradient of the body component  $\Delta_{M_t} \tilde{S}_b$  can be rewritten as:

$$\begin{aligned}\Delta_{M_t} \tilde{S}_b &= \frac{1}{2} \ln\left(\frac{2\pi\sigma^2 e}{N_{M_t} + 1}\right) \left[ 1 - \frac{1}{2} \sum_m^{A_t} \operatorname{erfc}\left(\frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{\sqrt{2\sigma^2} \sqrt{1 + \frac{N_m}{(N_{M_t} + 1)^2}}}\right) \right] \\ &\quad - \frac{1}{2} \ln\left(\frac{2\pi\sigma^2 e}{N_{M_t}}\right) \left[ 1 - \frac{1}{2} \sum_m^{A_t} \operatorname{erfc}\left(\frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{\sqrt{2\sigma^2}}\right) \right].\end{aligned}\quad (31)$$

The increment of the tail component along the better empirical arm can be calculated as:

$$\begin{aligned}
\Delta_{M_t} \tilde{S}_t &= \sum_m^{A_t} \int_{-\infty}^{\infty} \frac{e^{-\frac{\mu^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}} \frac{\sqrt{N_m} \left( \frac{\mu}{N_{M_t}+1} + \tilde{\mu}_{\text{eq},m} - \hat{\mu}_m \right)}{2\sqrt{2\pi\sigma^2}} e^{-\frac{N_m(\tilde{\mu}_{\text{eq},m} + \frac{\mu}{N_{M_t}+1} - \hat{\mu}_m)^2}{2\sigma^2}} d\mu \\
&\quad - \sum_m^{A_t} \frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{2\sqrt{2\pi\sigma^2}} e^{-\frac{N_m(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{2\sigma^2}} \\
&= \sum_m^{A_t} e^{-N_m \frac{(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{2\sigma^2 \left(1 + \frac{N_m}{(1+N_{M_t})^2}\right)}} \sqrt{\frac{N_m}{8\pi\sigma^2}} \frac{(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{\left(1 + \frac{N_m}{(N_{M_t}+1)^2}\right)^{3/2}} \\
&\quad - \sum_m^{A_t} \frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{2\sqrt{2\pi\sigma^2}} e^{-\frac{N_m(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{2\sigma^2}}.
\end{aligned} \tag{32}$$

Next, we consider the increment evaluation along a given worse empirical arm denoted by  $k$ ,

$$\Delta_k \tilde{S}_m = \int_{-\infty}^{\infty} \frac{e^{-\frac{\mu^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}} \left[ \tilde{S}_m(\dots, \hat{\mu}_k + \frac{\mu}{N_k+1}, N_k+1, \dots) - \tilde{S}_m(\dots, \hat{\mu}_k, N_k, \dots) \right] d\mu. \tag{33}$$

We here also neglect the variations of the subdominant term inside  $\tilde{\mu}_{\text{eq},m}$ . We start by considering the increment of the body component:

$$\begin{aligned}
\Delta_k \tilde{S}_b &= \frac{1}{2} \ln \left( \frac{2\pi\sigma^2 e}{N_{M_t}} \right) \left[ 1 - \frac{1}{2} \operatorname{erfc} \left( \frac{(N_k+1)(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)}{\sqrt{2\sigma^2(N_k+2)}} \right) \right] \\
&\quad - \frac{1}{2} \ln \left( \frac{2\pi\sigma^2 e}{N_{M_t}} \right) \left[ 1 - \frac{1}{2} \operatorname{erfc} \left( \frac{\sqrt{N_k}(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)}{\sqrt{2\sigma^2}} \right) \right].
\end{aligned} \tag{34}$$

Of note, all other terms in the sum independent of index  $k$  remain constant, showing no increment. Finally, we consider the associated tail component of the increment along  $k$ :

$$\begin{aligned}
\Delta_k \tilde{S}_t &= \int_{-\infty}^{\infty} \frac{e^{-\frac{\mu^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}} \frac{\sqrt{N_k+1} \left( \frac{\mu}{N_k+1} + \tilde{\mu}_{\text{eq},k} - \hat{\mu}_k \right)}{2\sqrt{2\pi\sigma^2}} e^{-\frac{(N_k+1)(\tilde{\mu}_{\text{eq},k} + \frac{\mu}{N_k+1} - \hat{\mu}_k)^2}{2\sigma^2}} d\mu \\
&\quad - \frac{\sqrt{N_k}(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)}{2\sqrt{2\pi\sigma^2}} e^{-\frac{N_k(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)^2}{2\sigma^2}} \\
&= e^{-\frac{(N_k+1)^2}{(N_k+2)} \frac{(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)^2}{2\sigma^2}} \frac{(1+N_k)^2 (\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)}{\sqrt{8\pi\sigma^2} (2+N_k)^{3/2}} \\
&\quad - \frac{\sqrt{N_k}(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)}{2\sqrt{2\pi\sigma^2}} e^{-\frac{N_k(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)^2}{2\sigma^2}}
\end{aligned} \tag{35}$$

As for the body increment, all other terms in the sum independent of index  $k$  remain constant, showing no increment.

Taken altogether, Equations (31), (32), (34) and (35) lead to the final analytical expression of the increment:

$$\begin{aligned}
\Delta_{M_t,k} = & \frac{1}{2} \ln\left(\frac{N_{M_t}}{N_{M_t}+1}\right) - \frac{1}{4} \ln\left(\frac{2\pi\sigma^2 e}{N_{M_t}+1}\right) \sum_m^{A_t} \operatorname{erfc} \left[ \frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{\sqrt{2\sigma^2(1 + \frac{N_m}{(N_{M_t}+1)^2})}} \right] \\
& + \frac{1}{4} \ln\left(\frac{2\pi\sigma^2}{N_{M_t}}\right) \sum_m^{A_t} \operatorname{erfc} \left[ \sqrt{N_m} \frac{\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m}{\sqrt{2\sigma^2}} \right] \\
& + \sum_m^{A_t} e^{-N_m \frac{(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{2\sigma^2(1 + \frac{N_m}{(1+N_{M_t})^2})}} \sqrt{\frac{N_m}{8\pi\sigma^2}} \frac{(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{(1 + \frac{N_m}{(N_{M_t}+1)^2})^{3/2}} \\
& - \sum_m^{A_t} \frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{\sqrt{8\pi\sigma^2}} e^{-\frac{N_m(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{2\sigma^2}} \\
& + \frac{1}{4} \ln\left(\frac{2\pi\sigma^2 e}{N_{M_t}}\right) \left[ \operatorname{erfc} \left( \frac{(N_k+1)(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)}{\sqrt{2\sigma^2(2+N_k)}} \right) - \operatorname{erfc} \left( \sqrt{N_k} \frac{\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k}{\sqrt{2\sigma^2}} \right) \right] \\
& - e^{-\frac{(N_k+1)^2}{(N_k+2)} \frac{(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)^2}{2\sigma^2}} \frac{(1+N_k)^2 (\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)}{\sqrt{8\pi\sigma^2(2+N_k)^{3/2}}} + \frac{\sqrt{N_k}(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)}{2\sqrt{2\pi\sigma^2}} e^{-\frac{N_k(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)^2}{2\sigma^2}}
\end{aligned} \tag{36}$$

To obtain a simplified expression, we expand to the first order each component of the different components of Equation (36) denoted  $T_1, T_2, T_3, T_4$ . The former is given by:

$$\begin{aligned}
T_1 = & -\frac{1}{4} \ln\left(\frac{2\pi\sigma^2 e}{N_{M_t}+1}\right) \sum_m^{A_t} \operatorname{erfc} \left[ \frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{\sqrt{2\sigma^2(1 + \frac{N_m}{(N_{M_t}+1)^2})}} \right] \\
& + \frac{1}{4} \ln\left(\frac{2\pi\sigma^2 e}{N_{M_t}}\right) \sum_{m \neq M_t}^K \operatorname{erfc} \left[ \sqrt{N_m} \frac{\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m}{\sqrt{2\sigma^2}} \right] \\
\approx & \sum_m^{A_t} \frac{1}{4N_{M_t}} \operatorname{erfc} \left( \frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{\sqrt{2\sigma^2}} \right) \\
& - \frac{1}{4} \ln\left(\frac{2\pi\sigma^2 e}{N_{M_t}}\right) \frac{N_m}{N_{M_t}^2} \frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{\sqrt{2\pi\sigma^2}} e^{-\frac{N_m(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{2\sigma^2}}.
\end{aligned} \tag{37}$$

Next we consider the second component, which reads:

$$\begin{aligned}
T_2 = & \sum_m^{A_t} e^{-N_m \frac{(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{2\sigma^2(1 + \frac{N_m}{(1+N_{M_t})^2})}} \sqrt{\frac{N_m}{8\pi\sigma^2}} \frac{(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{(1 + \frac{N_m}{(N_{M_t}+1)^2})^{3/2}} \\
& - \sum_m^{A_t} \frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{2\sqrt{2\pi\sigma^2}} e^{-\frac{N_m(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{2\sigma^2}} \\
\approx & \sum_m^{A_t} e^{-N_m \frac{(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{2\sigma^2}} \sqrt{\frac{N_m}{8\pi\sigma^2}} (\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m) \left[ -\frac{3}{2} \frac{N_m}{N_{M_t}^2} + \frac{(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{2\sigma^2} \frac{N_m^2}{N_{M_t}^2} \right].
\end{aligned} \tag{38}$$

Next we consider the third term, denoted  $T_3$ , which reads:

$$\begin{aligned}
T_3 = & \frac{1}{4} \ln\left(\frac{2\pi\sigma^2 e}{N_{M_t}}\right) \left[ \operatorname{erfc} \left( \frac{(N_k+1)(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)}{\sqrt{2\sigma^2(2+N_k)}} \right) - \operatorname{erfc} \left( \sqrt{N_k} \frac{\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k}{\sqrt{2\sigma^2}} \right) \right] \\
\approx & -\frac{1}{4} \ln\left(\frac{2\pi\sigma^2 e}{N_{M_t}}\right) \frac{1}{N_k^2} \left( \frac{\sqrt{N_k}(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)}{\sqrt{2\pi\sigma^2}} \right) e^{-\frac{N_k(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)^2}{2\sigma^2}}
\end{aligned} \tag{39}$$

Finally, the last term  $T_4$  reads

$$\begin{aligned}
T_4 &= -e^{-\frac{(N_k+1)^2}{(N_k+2)^2} \frac{(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)^2}{2\sigma^2}} \frac{(1+N_k)^2 (\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)}{\sqrt{8\pi\sigma^2}(2+N_k)^{3/2}} + \frac{\sqrt{N_k}(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)}{2\sqrt{2\pi\sigma^2}} e^{-\frac{N_k(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)^2}{2\sigma^2}} \\
&\approx e^{-N_k \frac{(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)^2}{2\sigma^2}} \sqrt{\frac{N_k}{8\pi\sigma^2}} (\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k) \left[ \frac{1}{N_k} + \frac{1}{N_k} \frac{(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)^2}{2\sigma^2} \right].
\end{aligned} \tag{40}$$

Taken altogether, we finally obtain the following simplified increment :

$$\begin{aligned}
\tilde{\Delta}_{M_t,k} &= \frac{1}{2} \ln\left(\frac{N_{M_t}}{N_{M_t}+1}\right) + \frac{1}{2N_{M_t}} \sum_m^{\mathcal{A}_t} \frac{1}{2} \operatorname{erfc}\left(\frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{\sqrt{2\sigma^2}}\right) \\
&+ \sum_m^{\mathcal{A}_t} \frac{N_m^{3/2} (\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{\sqrt{2\pi\sigma^2} N_{M_t}^2} e^{-N_m \frac{(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{2\sigma^2}} \left[ \frac{1}{4} \ln\left(\frac{N_{M_t}}{2\pi\sigma^2 e}\right) - \frac{3}{4} + \frac{N_m(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{4\sigma^2} \right] \\
&+ \frac{\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k}{\sqrt{2\pi\sigma^2} N_k} e^{-N_k \frac{(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)^2}{2\sigma^2}} \left[ \frac{1}{4N_k} \ln\left(\frac{N_{M_t}}{2\pi\sigma^2 e}\right) + \frac{1}{2} + \frac{(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)^2}{4\sigma^2} \right]
\end{aligned} \tag{41}$$

By noticing that the sum of second term should account for the tail contribution along the body increment, it shouldn't be allowed to be superior to one. Then we assume to bound it by taking the minimum compared to  $1 - 1/K$ .

#### A.6 Final expression for the increment comparison

Taken altogether, it leads to expression used for AIM for multiple gaussian arms:

$$\begin{aligned}
\Delta_{M_t,k} &= \frac{1}{2} \ln\left(\frac{N_{M_t}}{N_{M_t}+1}\right) + \frac{1}{2N_{M_t}} \min\left(\sum_m^{\mathcal{A}_t} \frac{1}{2} \operatorname{erfc}\left(\frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{\sqrt{2\sigma^2}}\right), 1 - \frac{1}{K}\right) \\
&+ \sum_m^{\mathcal{A}_t} \frac{N_m^{3/2} (\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)}{\sqrt{2\pi\sigma^2} N_{M_t}^2} e^{-N_m \frac{(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{2\sigma^2}} \left[ \frac{1}{4} \ln\left(\frac{N_{M_t}}{2\pi\sigma^2 e}\right) - \frac{3}{4} + \frac{N_m(\tilde{\mu}_{\text{eq},m} - \hat{\mu}_m)^2}{4\sigma^2} \right] \\
&+ \frac{\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k}{\sqrt{2\pi\sigma^2} N_k} e^{-N_k \frac{(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)^2}{2\sigma^2}} \left[ \frac{1}{4N_k} \ln\left(\frac{N_{M_t}}{2\pi\sigma^2 e}\right) + \frac{1}{2} + \frac{(\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k)^2}{4\sigma^2} \right]
\end{aligned} \tag{42}$$

with

$$\tilde{\mu}_{\text{eq},i} = \hat{\mu}_{M_t} + \frac{N_i(\hat{\mu}_{M_t} - \hat{\mu}_i)}{N_{M_t} - N_i} + \sqrt{\frac{N_{M_t} N_i (\hat{\mu}_{M_t} - \hat{\mu}_i)^2}{(N_{M_t} - N_i)^2} + \frac{\sigma^2 \ln\left(\frac{N_{M_t}}{N_i}\right)}{N_{M_t} - N_i}}. \tag{43}$$

## B Proof of Theorem 1

This section provides the complete proof of Theorem 1. More precisely, it proves the more refined Theorem 2 below.

**Theorem 2.** *For any multi-armed bandits with Gaussian rewards of variance  $\sigma^2$  and mean vector  $\boldsymbol{\mu} \in \mathbb{R}^K$ , for any  $\varepsilon \in (0, \frac{1}{2})$ , there exists a constant  $C(\boldsymbol{\mu}, \varepsilon) \in \mathbb{R}$  depending solely on  $\boldsymbol{\mu}$  and  $\varepsilon$  such that for any  $T \in \mathbb{N}$*

$$R(T) \leq \sum_{k, \mu_k < \mu^*} \left[ \frac{2\sigma^2 \ln T}{(1-\varepsilon)(\mu^* - \mu_k)} + \frac{2\sigma^2 \ln \ln T}{(1-\varepsilon)(\mu^* - \mu_k)} \right] + C(\boldsymbol{\mu}, \varepsilon).$$

*Proof.* We denote in the whole proof  $\mathcal{M}^* = \{k \in [K] \mid \mu_k = \mu^*\}$ . For  $\Delta_k = \mu^* - \mu_k$ , the regret can then be written as

$$R(T) = \sum_{k, \Delta_k > 0} \Delta_k \mathbb{E}[N_k(T)].$$

We decompose this expectation in 4 terms as follows

$$\begin{aligned} \mathbb{E}[N_k(T)] &\leq \sum_{t=1}^T \mathbb{P}(\forall i \in \mathcal{M}^*, N_i(t) \leq \sqrt{t}) + \sum_{t=1}^T \mathbb{P}\left(\hat{\mu}_k(t) \geq \mu^* - \sqrt{\frac{6\sigma^2 \ln t}{\sqrt{t}}}, a_t = k\right) \\ &\quad + \sum_{t=1}^T \mathbb{P}\left(\exists i \in \mathcal{M}^*, \hat{\mu}_i(t) \leq \mu_i - \sqrt{\frac{6\sigma^2 \ln t}{N_i(t)}}\right) + \sum_{t=1}^T \mathbb{P}(\mathcal{E}_k(t)), \end{aligned}$$

where

$$\mathcal{E}_k(t) := \left\{ \exists i \in \mathcal{M}^*, N_i(t) \geq \sqrt{t} \text{ and } \hat{\mu}_i(t) \geq \mu^* - \sqrt{\frac{6\sigma^2 \ln t}{N_i(t)}} \geq \hat{\mu}_k(t), a_t = k \right\}.$$

This inequality comes simply by noticing the event  $\{a_t = k\}$  is included in the union of the 4 other events. Lemmas 1, 3 and 4 allow to respectively bound the first, second and third sums by a constant  $C(\boldsymbol{\mu})$  depending solely on  $\boldsymbol{\mu}$ , so that

$$\mathbb{E}[N_k(T)] \leq \sum_{t=1}^T \mathbb{P}(\mathcal{E}_k(t)) + C(\boldsymbol{\mu}).$$

Thanks to Lemma 5, there exist constants  $t(\boldsymbol{\mu}), n(\boldsymbol{\mu})$  depending solely on  $K$  and  $\Delta_k$  such that

$$\sum_{t=1}^T \mathbb{P}(\mathcal{E}_k(t)) \leq t(\boldsymbol{\mu}) + \sum_{t=1}^T \mathbb{P}(\mathcal{G}_1(t)) + \mathbb{P}(\mathcal{G}_2(t)),$$

where

$$\begin{aligned} \mathcal{G}_1(t) &= \{\mu_k - \hat{\mu}_k(t) \leq -\varepsilon \Delta_k, a_t = k\}, \\ \mathcal{G}_2(t) &= \{N_k(t) \leq \frac{2\sigma^2}{(1-2\varepsilon)^2 \Delta_k^2} (\ln t + \ln \ln t) + n(\boldsymbol{\mu}), a_t = k\}. \end{aligned}$$

Now, we bound individually the sum corresponding to each of these 2 events. The first one can be bounded using Hoeffding's inequality. Indeed, for independent random variables  $Z_k(n) \sim \mathcal{N}(\mu_k, \sigma^2)$ , it reads as:

$$\begin{aligned} \sum_{t=1}^T \mathbb{P}(\mu_k - \hat{\mu}_k(t) \leq -\varepsilon \Delta_k, a_t = k) &\leq \sum_{n=1}^T \mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n Z_k(i) - \mu_k \geq \varepsilon \Delta_k\right) \\ &\leq \sum_{n=1}^T e^{-\frac{n\varepsilon^2 \Delta_k^2}{2\sigma^2}} \\ &\leq \frac{1}{e^{\frac{\varepsilon^2 \Delta_k^2}{2\sigma^2}} - 1}. \end{aligned}$$

The bound of the second term is bounded as

$$\begin{aligned} \sum_{t=1}^T \mathbb{P}(\mathcal{E}_k(t)) &\leq \mathbb{E}\left[\sum_{t=1}^T \mathbb{1}\{N_k(t) \leq \frac{2\sigma^2}{(1-2\varepsilon)^2 \Delta_k^2} (\ln t + \ln \ln t) + n(\boldsymbol{\mu}), a_t = k\}\right] \\ &\leq \frac{2\sigma^2}{(1-2\varepsilon)^2 \Delta_k^2} (\ln T + \ln \ln T) + n(\boldsymbol{\mu}) + 1. \end{aligned}$$

Wrapping up everything finally yields that for some constant  $C(\boldsymbol{\mu}, \varepsilon)$  depending solely on  $\boldsymbol{\mu}, \varepsilon$ ,

$$R(T) \leq \frac{2\sigma^2}{(1-2\varepsilon)^2 \Delta_k} (\ln T + \ln \ln T) + C(\boldsymbol{\mu}, \varepsilon).$$

This concludes the proof of Theorem 2 with the reparameterization  $\varepsilon \leftarrow 1 - (1 - 2\varepsilon)^2$ .  $\square$

## B.1 Auxiliary Lemmas

Similarly to the proof of Thompson sampling, the first part of the proof shows that the optimal arm is at least pulled a polynomial number of times with high probability. We recall that we denote in this whole section  $\mathcal{M}^* = \arg \max_k \mu_k$ .

**Lemma 1.** *There exists a constant  $C_0(\boldsymbol{\mu})$  depending solely on the mean vector  $\boldsymbol{\mu}$  such that*

$$\sum_{t=1}^{\infty} \mathbb{P}(\forall i \in \mathcal{M}^*, N_i(t) \leq \sqrt{t}) \leq C_0(\boldsymbol{\mu}).$$

*Proof.* Let  $t_0(\boldsymbol{\mu})$  be a large constant that depends solely on  $\boldsymbol{\mu}$ . In the remaining of the proof, we assume at some points that  $t_0(\boldsymbol{\mu})$  is chosen large enough (but only larger than a threshold depending on  $\boldsymbol{\mu}$ ) such that some inequalities hold. We also assume in the following, without loss of generality, that  $\mu_1 = \mu^*$ , i.e.,  $1 \in \mathcal{M}^*$ .

Assume that for  $t \geq t_0(\boldsymbol{\mu})$ ,  $N_i(t) \leq \sqrt{t}$  for all  $i \in \mathcal{M}^*$ . Let then  $k$  be the most pulled arm at time  $t$ , i.e.,  $k \in \arg \max_j N_j(t)$  (if multiple arms maximise the number of pulls, we select the one such that its last pull happened the earliest). Necessarily  $N_k(t) \geq \frac{t}{K}$ . We can choose  $t_0(\boldsymbol{\mu})$  large enough so that  $\frac{t}{K} > \sqrt{t}$  and thus  $\Delta_k > 0$ . Let  $t' \leq t$  be the last time  $k$  was pulled. By design,  $k$  also maximised the number of pulls then, so that  $k \in \arg \max_j \hat{\mu}_j(t')$ . Moreover,  $N_i(t') \leq \sqrt{t}$  for all  $i \in \mathcal{M}^*$  and  $N_k(t') \geq \frac{t}{K} - 1$ . For  $t_0(\boldsymbol{\mu})$  large enough, this yields  $N_k(t') \geq N_i(t')$  for all  $i \in \mathcal{M}^*$  and  $a_{t'} = k$ . The arm  $k$  is thus pulled at time  $t'$ , in particular because  $S_k \geq S_1$  (i.e.,  $\Delta_{k,1} S \leq 0$ ), where

$$\begin{aligned} S_k &= \frac{1}{2} \ln \left( 1 + \frac{1}{N_k(t')} \right) - \frac{1}{2N_k(t')} \min \left( \frac{1}{2} \sum_{i \neq k} \operatorname{erfc} \left( \frac{\sqrt{N_i(t')} (\tilde{\mu}_{\text{eq},i} - \hat{\mu}_i)}{\sqrt{2\sigma^2}} \right), 1 - \frac{1}{K} \right), \\ S_1 &= \sum_{i \neq k} \sqrt{\frac{N_i(t')}{2\pi\sigma^2}} (\tilde{\mu}_{\text{eq},i} - \hat{\mu}_i) e^{-N_i(t') \frac{(\tilde{\mu}_{\text{eq},i} - \hat{\mu}_i)^2}{2\sigma^2}} \left[ \frac{1}{4} \ln \left( \frac{N_k(t')}{2\pi\sigma^2 e} \frac{N_i(t')}{N_k^2(t')} - \frac{3}{4} \frac{N_i(t')}{N_k^2(t')} + \frac{(\tilde{\mu}_{\text{eq},i} - \hat{\mu}_i)^2}{4\sigma^2} \frac{N_i^2(t')}{N_k^2(t')} \right) \right] \\ &\quad + e^{-N_1(t') \frac{(\tilde{\mu}_{\text{eq},1} - \hat{\mu}_1)^2}{2\sigma^2}} \sqrt{\frac{N_1(t')}{2\pi\sigma^2}} (\tilde{\mu}_{\text{eq},1} - \hat{\mu}_1) \left[ \frac{1}{4} \ln \left( \frac{N_k(t')}{2\pi\sigma^2 e} \frac{1}{N_1^2(t')} + \frac{1}{2N_1(t')} + \frac{1}{N_1(t')} \frac{(\tilde{\mu}_{\text{eq},1} - \hat{\mu}_1)^2}{4\sigma^2} \right) \right]. \end{aligned}$$

To simplify, note that  $S_k \leq \frac{1}{2N_k(t')}$ . Moreover since  $N_k(t') \geq \frac{t}{K} - 1 \geq 2\pi e^4 \sigma^2$  for a large enough choice of  $t_0(\boldsymbol{\mu})$ ,  $S_1$  can be easily lower bounded as

$$S_1 \geq \frac{1}{2} \frac{(\tilde{\mu}_{\text{eq},1} - \hat{\mu}_1)}{\sqrt{2\pi\sigma^2 N_1(t')}} e^{-\frac{N_1(t')(\tilde{\mu}_{\text{eq},1} - \hat{\mu}_1)^2}{2\sigma^2}}.$$

So we finally have the following inequality at time  $t'$ :

$$\frac{1}{N_2(t')} \geq \frac{(\tilde{\mu}_{\text{eq}} - \hat{\mu}_1)}{\sqrt{2\pi\sigma^2 N_1(t')}} e^{-\frac{N_1(t')(\tilde{\mu}_{\text{eq}} - \hat{\mu}_1)^2}{2\sigma^2}}. \quad (44)$$

Recall that  $N_2(t') \geq \frac{t}{K} - 1$ , so that Equation (44) can be rewritten as

$$N_1(t') \geq \left( \frac{t}{K} - 1 \right) \frac{\tilde{x}}{\sqrt{\pi}} e^{-\tilde{x}^2}, \quad (45)$$

where  $\tilde{x} = \frac{\sqrt{N_1(t')(\tilde{\mu}_{\text{eq}} - \hat{\mu}_1)}}{\sqrt{2\sigma^2}}$ . In the following, we will show that  $\tilde{x} \in [\tilde{x}_{\min}, \tilde{x}_{\max}] \subset \mathbb{R}_+$ . By analysing the variations of  $x \mapsto x e^{-x^2}$ , this will imply that

$$N_1(t') \geq \frac{\frac{t}{K} - 1}{\sqrt{\pi}} \min \{ \tilde{x}_{\min} e^{-\tilde{x}_{\min}^2}, \tilde{x}_{\max} e^{-\tilde{x}_{\max}^2} \}. \quad (46)$$

For the lower bound, the definition of  $\tilde{\mu}_{\text{eq},1}$  and the fact that  $N_1(t') \geq 1$  directly implies that

$$\tilde{x} \geq \sqrt{\frac{\ln \left( \frac{N_k(t')}{N_1(t')} \right)}{2 \left( \frac{N_k(t')}{N_1(t')} - 1 \right)}} = \Omega \left( \sqrt{\frac{\ln(t)}{t}} \right).$$



Moreover, by subadditivity of the square root:

$$\tilde{x} \leq \sqrt{\frac{N_1(t')}{2\sigma^2}}(\hat{\mu}_2 - \hat{\mu}_1) \left( 1 + \frac{N_1(t') + \sqrt{N_1(t')N_k(t')}}{N_k(t') - N_1(t')} \right) + \sqrt{\frac{N_1(t') \ln(\frac{N_k(t')}{N_1(t')})}{2(N_k(t') - N_1(t'))}} \quad (47)$$

$$\leq \sqrt{\frac{N_1(t')}{2\sigma^2}}(\hat{\mu}_k - \hat{\mu}_1) \left( 1 + Kt^{-\frac{1}{4}} \right) + \mathcal{O} \left( \frac{\sqrt{K \ln(t)}}{t^{\frac{1}{4}}} \right). \quad (48)$$

Let us now consider the events, for  $\Delta_{\min} = \min_{j, \Delta_j > 0} \Delta_j$ ,

$$\mathcal{H}_*(t) := \left\{ \exists i \in \mathcal{M}^*, \exists s \leq t, \hat{\mu}_i(s) - \mu_i \leq -\sqrt{\frac{2\sigma^2(\ln(t) - \ln \ln(t))}{N_i(s)}} - \frac{\Delta_{\min}}{3} \right\}, \quad (49)$$

$$\mathcal{H}_k(t) := \left\{ \exists s \leq t, \frac{t}{K} - 1 \leq N_k(s) \leq t \text{ and } \hat{\mu}_k(s) - \mu_k \geq \frac{\Delta_k}{3} \right\}. \quad (50)$$

Assume in the following that  $\neg \mathcal{H}_*(t) \cap \neg \mathcal{H}_k(t)$ . This implies that

$$\hat{\mu}_k - \hat{\mu}_1 \leq -\frac{\Delta_k}{3} + \sqrt{\frac{2\sigma^2(\ln(t) - \ln \ln(t))}{N_1(s)}}. \quad (51)$$

In particular,

$$\sqrt{\frac{N_1(t')}{2\sigma^2}}(\hat{\mu}_2 - \hat{\mu}_1) \leq \sqrt{\ln(t) - \ln \ln(t)},$$

which implies that  $\tilde{x} \leq \sqrt{\ln(t) - \ln \ln(t)} + \mathcal{O} \left( K \sqrt{\ln(t)} t^{-\frac{1}{4}} \right)$ . Using the lower and upper bounds on  $\tilde{x}$ , we have thanks to Equation (46) that under  $\neg \mathcal{H}_*(t) \cap \neg \mathcal{H}_k(t)$ ,

$$N_1(t') = \Omega \left( \frac{\ln^{\frac{3}{2}}(t)}{K} \right).$$

For a large enough choice of  $t_0(\boldsymbol{\mu})$ , this last equality along with Equation (51) actually yield  $\hat{\mu}_k - \hat{\mu}_1 < 0$ , which contradicts the beginning of the proof ( $k$  being best empirical arm at time  $t'$ ). By contradiction, we thus showed the following event inclusion for  $t \geq t_0(\boldsymbol{\mu})$ :

$$\left\{ \forall i \in \mathcal{M}^*, N_i(t) \leq \sqrt{t} \right\} \subset \mathcal{H}_*(t) \cup \mathcal{H}_k(t). \quad (52)$$

Lemma 1 then follows, thanks to Lemma 2 below,

$$\sum_{t=1}^{\infty} \mathbb{P}(\forall i \in \mathcal{M}^*, N_i(t) \leq \sqrt{t}) \leq t_0(\boldsymbol{\mu}) + \sum_{t=t_0(\boldsymbol{\mu})+1}^{\infty} \mathbb{P}(\mathcal{H}_*(t)) + \mathbb{P}(\mathcal{H}_k(t)).$$

□

**Lemma 2.** For any  $b \in (0, 1)$  and the events  $\mathcal{H}_*(t), \mathcal{H}_k(t)$  defined in Equations (49) and (50), there exist constants  $c_1$  and  $c_2$  depending solely on  $\boldsymbol{\mu}$  such that

$$\sum_{t=1}^{\infty} \mathbb{P}(\mathcal{H}_*(t)) \leq c_1 \quad \text{and} \quad \sum_{t=1}^{\infty} \mathbb{P}(\mathcal{H}_k(t)) \leq c_2 \quad \text{for any } k \notin \mathcal{M}^*.$$

*Proof.* The two bounds directly result from Hoeffding's inequality. Consider independent random variables  $(Z_j(n))_{n \in \mathbb{N}, j \in [K]}$  where  $Z_j(n) \sim \mathcal{N}(\mu_j, \sigma^2)$ . Let us first bound the probability of  $\mathcal{H}_k(t)$ , which is simpler.

$$\begin{aligned} \mathbb{P}(\mathcal{H}_k(t)) &\leq \sum_{n=\lceil t-t^b-1 \rceil}^t \mathbb{P} \left( \sum_{i=1}^n (Z_k(i) - \mu_k) \geq \frac{n\Delta_k}{3} \right) \\ &\leq \sum_{n=\lceil \frac{t}{K}-1 \rceil}^t e^{-\frac{n\Delta_k^2}{18\sigma^2}} \\ &\leq \frac{e^{-\lceil \frac{t}{K}-1 \rceil \frac{\Delta_k^2}{18\sigma^2}}}{1 - e^{-\frac{\Delta_k^2}{18\sigma^2}}}. \end{aligned}$$

The second inequality of Lemma 2 then follows by noting that the last term is summable over  $t$ . For the second bound, we also have by Hoeffding's inequality

$$\begin{aligned} \mathbb{P}(\mathcal{H}_*(t)) &\leq \sum_{j \in \mathcal{M}^*} \sum_{n=1}^{\infty} \mathbb{P} \left( \sum_{i=1}^n (Z_j(i) - \mu_j) \leq -\sqrt{2n\sigma^2(\ln(t) - \ln \ln(t))} - \frac{n\Delta_{\min}}{3} \right) \\ &\leq \sum_{j \in \mathcal{M}^*} \sum_{n=1}^{\infty} \exp \left( -\ln(t) + \ln \ln(t) - \sqrt{2n\sigma^2(\ln(t) - \ln \ln(t))} \frac{\Delta_{\min}}{3\sigma^2} - \frac{n\Delta_{\min}^2}{18\sigma^2} \right) \\ &\leq |\mathcal{M}^*| \frac{\ln(t)}{t} \exp \left( -\sqrt{2(\ln(t) - \ln \ln(t))} \frac{\Delta_{\min}}{3\sqrt{\sigma^2}} \right) \sum_{n=1}^{\infty} e^{-\frac{n\Delta_{\min}^2}{18\sigma^2}}. \end{aligned}$$

The last sum is obviously finite. Moreover,  $\sqrt{2(\ln(t) - \ln \ln(t))} = \omega(\ln \ln(t))$ , so that  $\exp \left( -\sqrt{2(\ln(t) - \ln \ln(t))} \frac{\Delta_{\min}}{3\sqrt{\sigma^2}} \right) = \mathcal{O} \left( \frac{1}{\ln^\alpha(t)} \right)$  for any  $\alpha > 0$ . By comparison with series of the form  $\frac{1}{n \ln^\alpha(n)}$ , the term  $\frac{\ln(t)}{t} \exp \left( -\sqrt{2(\ln(t) - \ln \ln(t))} \frac{\Delta_{\min}}{3\sqrt{\sigma^2}} \right)$  is summable over  $t$ , which leads to the first bound of Lemma 2.  $\square$

**Lemma 3.** For any  $k \notin \mathcal{M}^*$ , there exists a constant  $C_1(\boldsymbol{\mu})$  depending solely on  $\boldsymbol{\mu}$  such that

$$\sum_{t=1}^{\infty} \mathbb{P} \left( \hat{\mu}_k(t) \geq \mu^* - \sqrt{\frac{6\sigma^2 \ln t}{t}}, a_t = k \right) \leq C_1(\boldsymbol{\mu}).$$

*Proof.* A union bound on the sum yields for any  $T \in \mathbb{N}$

$$\begin{aligned} \sum_{t=1}^T \mathbb{P} \left( \hat{\mu}_2(t) \geq \mu^* - \sqrt{\frac{6\sigma^2 \ln t}{t}}, a_t = k \right) \\ \leq \sum_{t=1}^T \sum_{n=0}^t \mathbb{P} \left( \hat{\mu}_k(t) \geq \mu^* - \sqrt{\frac{6\sigma^2 \ln t}{t}}, N_k(t) = n, N_k(t+1) = n+1 \right) \\ \leq \sum_{n=0}^T \sum_{t=n}^T \mathbb{P} \left( \underbrace{\hat{\mu}_k(t) \geq \mu^* - \sqrt{\frac{6\sigma^2 \min_{s \geq n} \ln s}{s}}}_{:= \mathcal{G}_1(t, n)}, N_k(t) = n, N_k(t+1) = n+1 \right). \end{aligned}$$

Now note that the  $\mathcal{G}_1(t, n)$  are disjoint for different  $t$ . In particular,

$$\sum_{t=n}^T \mathbb{P}(\mathcal{G}_1(t, n)) = \mathbb{P} \left( \exists t \in [n, T], \hat{\mu}_k(t) \geq \mu^* - \sqrt{\frac{6\sigma^2 \min_{s \geq n} \ln s}{s}}, N_2(t) = n \right).$$

For independent random variables  $Z_k(n) \sim \mathcal{N}(\mu_k, \sigma^2)$ , we have by independence of the  $X_t$  and  $a_t$ , and then by Hoeffding inequality:

$$\begin{aligned} \sum_{t=1}^T \mathbb{P} \left( \hat{\mu}_k(t) \geq \mu^* - \sqrt{\frac{6\sigma^2 \ln t}{t}}, a_t = k \right) &\leq 1 + \sum_{n=1}^T \mathbb{P} \left( \frac{1}{n} \sum_{i=1}^n (Z_k(i) - \mu_k) \geq \Delta_k - \sqrt{\frac{6\sigma^2 \min_{s \geq n} \ln s}{s}} \right), \\ &\leq 1 + \sum_{n=1}^T \exp \left( -\frac{n \left( \Delta_k - \sqrt{\frac{6\sigma^2 \min_{s \geq n} \ln s}{s}} \right)^2}{2\sigma^2} \right). \end{aligned}$$

Obviously, this sum can be bounded for any  $T \in \mathbb{N}$  by a constant solely depending on  $\Delta_k$ .  $\square$

**Lemma 4.** For any  $i \in [K]$ , there exists a universal constant  $C_2$  such that

$$\sum_{t=0}^{\infty} \mathbb{P} \left( \hat{\mu}_i(t) \leq \mu_i - \sqrt{\frac{6\sigma^2 \ln t}{N_i(t)}} \right) \leq C_2.$$

*Proof.* This is a direct consequence of Garivier [2013], which states that for Gaussian rewards with variance  $\sigma^2$ :

$$\mathbb{P}(N_i(t) \frac{(\hat{\mu}_i(t) - \mu_i)^2}{2\sigma^2} \geq (1+\alpha) \ln t) \leq 2 \left[ \frac{\ln t}{\ln(1+\eta)} \right] t^{-(1-\frac{\eta^2}{16})(1+\alpha)} \quad \text{for any } t \in \mathbb{N}^* \text{ and } \alpha, \eta > 0.$$

In particular with  $\alpha = 2 = \eta$ , this implies

$$\mathbb{P} \left( \hat{\mu}_i(t) \leq \mu_i - \sqrt{\frac{6\sigma^2 \ln t}{N_i(t)}} \right) \leq 2 \frac{\ln(t) + 1}{\ln(3)} t^{-\frac{9}{4}}.$$

This term is obviously summable so that there exists a constant  $C_2$  such that

$$\sum_{t=1}^{\infty} \mathbb{P} \left( \hat{\mu}_1(t) \leq \mu_1 - \sqrt{\frac{6\sigma^2 \ln t}{N_1(t)}} \right) \leq C_2.$$

Lemma 4 directly follows by the inclusion of the considered events.  $\square$

For any  $k \notin \mathcal{M}^*$ , Lemma 5 below gives an event inclusion for the event  $\mathcal{E}_k(t)$  that we recall here,

$$\mathcal{E}_k(t) := \left\{ \exists i \in \mathcal{M}^*, N_i(t) \geq \sqrt{t} \text{ and } \hat{\mu}_i(t) \geq \mu^* - \sqrt{\frac{6\sigma^2 \ln t}{N_i(t)}} \geq \hat{\mu}_k(t), a_t = k, \right\}.$$

**Lemma 5.** *There exist constants  $t(\boldsymbol{\mu})$  and  $n(\boldsymbol{\mu})$  depending solely on  $\boldsymbol{\mu}$  such that for any  $k \notin \mathcal{M}^*$ ,  $t \geq t(\boldsymbol{\mu})$  and  $\varepsilon \in (0, \frac{1}{3})$ ,*

$$\mathcal{E}_k(t) \subset \{\mu_k - \hat{\mu}_k(t) \leq -\varepsilon \Delta_k, a_t = k\} \cup \{N_k(t) \leq \frac{2\sigma^2}{(1-2\varepsilon)^2 \Delta_k^2} (\ln t + \ln \ln t) + n(\boldsymbol{\mu}), a_t = k\}.$$

*Proof.* Assume in the following that  $\mathcal{E}_k(t)$  holds for some  $t \geq t(\boldsymbol{\mu})$ . Let  $i \in [K]$  be an arm maximising the empirical mean at time  $t$ . Necessarily  $\hat{\mu}_i(t) \geq \mu^* - \sqrt{\frac{6\sigma^2 \ln t}{N_i(t)}}$ . Moreover,  $a_t = k$  so that  $i$  also maximises the number of pulls, in particular  $N_i(t) \geq \frac{t}{K}$ . Moreover, as we pull the arm  $k$ ,  $S_k \geq S_i$  where

$$S_i = \frac{1}{2} \ln \left( 1 + \frac{1}{N_i(t)} \right) - \frac{1}{2N_i(t)} \min \left( \frac{1}{2} \sum_{j \neq i} \operatorname{erfc} \left( \frac{\sqrt{N_j(t)} (\tilde{\mu}_{\text{eq},j} - \hat{\mu}_j)}{\sqrt{2\sigma^2}} \right), 1 - \frac{1}{K} \right), \quad (53)$$

$$S_k = g_k(t) Q_k(t) + \sum_{j \neq i} g_j(t) P_j(t),$$

where for all  $j \neq i$

$$\begin{aligned} g_j(t) &= \sqrt{\frac{N_j(t)}{2\pi\sigma^2}} (\tilde{\mu}_{\text{eq},j} - \hat{\mu}_j) e^{-N_j(t) \frac{(\tilde{\mu}_{\text{eq},j} - \hat{\mu}_j)^2}{2\sigma^2}}, \\ P_j(t) &= \left[ \frac{1}{4} \ln \left( \frac{N_i(t)}{2\pi\sigma^2 e} \frac{N_j(t)}{N_i^2(t)} - \frac{3}{4} \frac{N_j(t)}{N_i^2(t)} + \frac{(\tilde{\mu}_{\text{eq},j} - \hat{\mu}_j)^2 N_j(t)}{4\sigma^2 N_i^2(t)} \right) \right] \\ \text{and } Q_j(t) &= \left[ \frac{1}{4} \ln \left( \frac{N_i(t)}{2\pi\sigma^2 e} \frac{1}{N_j^2(t)} + \frac{1}{2N_j(t)} + \frac{1}{N_j(t)} \frac{(\tilde{\mu}_{\text{eq},j} - \hat{\mu}_j)^2}{4\sigma^2} \right) \right]. \end{aligned}$$

Also note that as we pull the arm  $k$ , we have for any  $j \leq i$ ,

$$g_j(t) Q_j(t) \leq g_k(t) Q_k(t). \quad (54)$$

As a consequence, we can write for any  $\delta > 0$  and  $\tilde{x}_j = \sqrt{\frac{N_j(t)}{2\sigma^2}} (\tilde{\mu}_{\text{eq},j} - \hat{\mu}_j)$ :

$$\begin{aligned} g_j(t) P_j(t) &\leq g_j(t) \left( \ln(t) \frac{N_j(t)}{4N_i^2(t)} \right) + \tilde{x}_j^3 e^{-\tilde{x}_j^2} \frac{N_j(t)}{2N_i^2(t)} \\ &\leq g_j(t) \left( 2 \ln(t) + \frac{1}{\delta} \right) \frac{N_j(t)}{2N_i^2(t)} + \frac{N_j(t)}{2N_i^2(t)} \delta, \end{aligned}$$

where we used the fact that  $\tilde{x}_j^3 e^{-\tilde{x}_j^2} \leq \frac{\tilde{x}_j e^{-\tilde{x}_j^2}}{\delta} + \delta$ . Moreover, note that for  $t(\boldsymbol{\mu})$  large enough,  $Q_j(t) \geq \frac{1}{2N_j(t)}$ . As a consequence,

$$\begin{aligned} g_j(t)P_j(t) &= g_j(t) \left( 2 \ln(t) + \frac{1}{\delta} \right) \frac{N_j(t)}{2N_i^2(t)Q_j(t)} Q_j(t) + \frac{N_j(t)}{2N_i^2(t)} \delta \\ &\leq Q_j(t)g_j(t) \left( 2 \ln(t) + \frac{1}{\delta} \right) + \frac{K\delta}{2t} \\ &\leq Q_k(t)g_k(t) \left( 2 \ln(t) + \frac{1}{\delta} \right) + \frac{K\delta}{2t}, \end{aligned}$$

where the last inequality comes from Equation (54). In particular,

$$S_k \leq K \left( 2 \ln(t) + \frac{1}{\delta} \right) Q_k(t)g_k(t) + \frac{K^2\delta}{2t}. \quad (55)$$

Also,  $S_i \geq \frac{1}{2t} - \frac{K^2}{4t^2}$  since  $N_i(t) \geq \frac{t}{K}$  and  $\ln(1+x) \geq x - \frac{x^2}{2}$  for  $x \in [0, 1]$ .

Now assume that  $\mu_k - \hat{\mu}_k(t) \geq -\varepsilon\Delta_k$ . It then holds

$$\begin{aligned} \tilde{\mu}_{\text{eq},k} - \hat{\mu}_k(t) &\geq \hat{\mu}_i(t) - \hat{\mu}_k(t) \\ &\geq \Delta_k + \mu_k - \hat{\mu}_k(t) - \sqrt{\frac{6\sigma^2 \ln t}{\sqrt{t}}} \\ &\geq (1 - \varepsilon)\Delta_k - \sqrt{\frac{6\sigma^2 \ln t}{\sqrt{t}}}. \end{aligned}$$

Again, we can choose  $t(\boldsymbol{\mu})$  large enough so that  $\tilde{\mu}_{\text{eq},k} - \hat{\mu}_k(t) \geq (1 - 2\varepsilon)\Delta_k$ . Moreover, note that the functions  $x \mapsto \frac{e^{-x^2}}{x}$ ,  $x \mapsto xe^{-x^2}$ ,  $x \mapsto x^3e^{-x^2}$  are all decreasing on an interval of the form  $[M, +\infty]$ . As a consequence, we can choose  $n(\boldsymbol{\mu})$  large enough so that  $\frac{\sqrt{n(\boldsymbol{\mu})((1-2\varepsilon)\Delta_k)^2}}{\sqrt{2\sigma^2}} \geq M$ . If  $N_k(t) \geq n(\boldsymbol{\mu})$ , we then have from Equation (55), for a constant  $c(K, \Delta_k)$  solely depending on  $K$  and  $\Delta_k$ :

$$S_k \leq c(K, \Delta_k) e^{-\frac{N_k(t)(1-2\varepsilon)^2\Delta_k^2}{2\sigma^2}} \left[ \ln t + \frac{1}{\delta} \right] + \frac{K^2\delta}{2t}.$$

The inequality  $S_k \geq S_i$  then implies, thanks to the above bounds:

$$c(K, \Delta_k) e^{-\frac{N_k(t)(1-2\varepsilon)^2\Delta_k^2}{2\sigma^2}} \left[ \ln t + \frac{1}{\delta} \right] \geq \frac{1 - K^2\delta}{2t} - \frac{K^2}{4t^2}.$$

In particular, for  $\delta = \frac{1}{2K^2}$ ,

$$N_k(t) \leq \frac{2\sigma^2}{(1-2\varepsilon)^2\Delta_k^2} (\ln t + \ln \ln t + \mathcal{O}(1)),$$

where the  $\mathcal{O}$  hides constants depending in  $K$  and  $\Delta_k$ . This concludes the proof of Lemma 5 as we just shown that if  $\mathcal{E}_k(t)$  holds, at least one of the two following events holds when  $N_k(t) \geq n(\boldsymbol{\mu})$ :

- $\mu_k - \hat{\mu}_k(t) \leq -\varepsilon\Delta_k$
- $N_k(t) \leq \frac{2\sigma^2}{(1-2\varepsilon)^2\Delta_k^2} (\ln t + \ln \ln t + \mathcal{O}(1))$ .

□

## C Generalization of the information maximization approximation

In this section, we will generalize the approach derived in Appendix A to bandit settings with a reward distribution belonging to the exponential family. We will retrace all the previous steps made in Appendix A, insisting on the differences with the Gaussian reward case. We will also discuss bandit settings with non-uniform priors and with more than two arms.

### C.1 Asymptotic expression for exponential family rewards

We derive an asymptotic expression for the one-dimensional canonical exponential family from which we will derive an analytical approximation of the entropy. We thus focus on a reward distribution density  $f$  with respect to some reference measure  $\nu$  belonging to some one-dimensional canonical exponential family, i.e.,

$$f(x|\theta) = A(x) \exp(T(x)\theta - F(\theta)), \quad (56)$$

where  $F$  is twice differentiable and strictly convex. Additionally, let us recall that the Kullback-Leibler divergence verifies: [Korda et al., 2013]:

$$\text{KL}(\theta, \theta') = F(\theta') - F(\theta) - F'(\theta)(\theta' - \theta), \quad (57)$$

where  $\text{KL}(\theta, \theta')$  is the Kullback-Leibler divergence between the reward distribution parameterized by  $\theta$  and the one parameterized by  $\theta'$ .

Given a prior  $\pi(\theta)$  and the reward realizations  $(x_1, \dots, x_n)$ , the associated posterior distribution on  $\theta$ , denoted  $p$ , reads:

$$p(\theta|x_1, \dots, x_n) = \frac{1}{C} \pi(\theta) \exp\left(\theta \sum_{k=1}^n T(x_k) - nF(\theta)\right), \quad (58)$$

where  $C = \int \pi(\theta) \exp(\theta \sum T(x_k) - nF(\theta)) d\theta$  is a normalization constant. Next, we derive the maximum a posteriori for the parameter  $\theta$ , denoted  $\hat{\theta}_l$ , which verifies:

$$\sum_{i=1}^n T(x_i) = nF'(\hat{\theta}_l) - \frac{\pi'(\hat{\theta}_l)}{\pi(\hat{\theta}_l)}. \quad (59)$$

At this stage, we assume that there exists such  $\hat{\theta}_l$  verifying Equation (59). In practice, for a reward distribution that does not meet this criteria, one can replace  $\hat{\theta}_l$  by a series  $\hat{\theta}_{n,l}$  which, for sufficiently large values of  $n$ , asymptotically conforms to the aforementioned definition. For example, a Bernoulli arm that consistently fails under a uniform prior, will result in an undefined  $\hat{\theta}_l$ . To address this, one may redefine  $\hat{\theta}_l$  such that  $(1 + \sum i = 1^n T(x_i)) = (n + 2)F'(\hat{\theta}_l)$ , effectively replacing the empirical mean in Equation (59) with the posterior mean.

Replacing the sum in Equation (58) leads to

$$\begin{aligned} p(\theta|x_1, \dots, x_n) &= \frac{1}{C} \pi(\theta) \exp\left(\theta nF'(\hat{\theta}_l) - \theta \frac{\pi'(\hat{\theta}_l)}{\pi(\hat{\theta}_l)} - nF(\theta)\right) \\ &= \frac{e^{n\hat{\theta}_l F'(\hat{\theta}_l) - nF(\hat{\theta}_l)}}{C} \pi(\theta) e^{-\theta \frac{\pi'(\hat{\theta}_l)}{\pi(\hat{\theta}_l)}} e^{-n\text{KL}(\hat{\theta}_l, \theta)} \\ &= \frac{1}{C_2} \pi(\theta) e^{-\theta \frac{\pi'(\hat{\theta}_l)}{\pi(\hat{\theta}_l)}} e^{-n\text{KL}(\hat{\theta}_l, \theta)}, \end{aligned} \quad (60)$$

where  $C_2$  also acts as a normalization constant of Equation (60). For  $n \gg 1$ , the distribution concentrates in the vicinity of  $\hat{\theta}_l$  from which we will derive the asymptotic scaling of  $C_2$ . We then integrate Equation (60) after a change of variable  $\theta(u) = \hat{\theta}_l + \frac{u}{\sqrt{n}}$ ,

$$\begin{aligned} 1 &= \int_{\Theta} p(\theta|x_1, \dots, x_n) d\theta = \int_{-(\theta_b - \hat{\theta}_l)\sqrt{n}}^{(\theta_b - \hat{\theta}_l)\sqrt{n}} \frac{1}{C_2 \sqrt{n}} \pi(\hat{\theta}_l + \frac{u}{\sqrt{n}}) e^{-(\hat{\theta}_l + \frac{u}{\sqrt{n}}) \frac{\pi'(\hat{\theta}_l)}{\pi(\hat{\theta}_l)}} e^{-n\text{KL}(\hat{\theta}_l, \hat{\theta}_l + \frac{u}{\sqrt{n}})} du \\ &\quad + \int_{\mu_{\text{inf}}}^{-\theta_b} p(\theta|x_1, \dots, x_n) d\theta + \int_{\theta_b}^{\mu_{\text{sup}}} p(\theta|x_1, \dots, x_n) d\theta \end{aligned} \quad (61)$$

Taking  $(\theta_b - \hat{\theta}_l) \sim n^{-b}$  with  $b < 1/2$ , we get rid of the tail components in the asymptotic limit. Secondly, by noticing that  $\frac{F''(\hat{\theta}_l)}{2} = \lim_{\theta \rightarrow \hat{\theta}_l} K(\hat{\theta}_l, \theta) / |\theta - \hat{\theta}_l|^2$  from Equation (57), we make an expansion to the lowest order of the Kullback-Leibler divergence, which gives:

$$1 = \lim_{\theta \rightarrow \hat{\theta}_l} \int_{-(\theta_b - \hat{\theta}_l)\sqrt{n}}^{(\theta_b - \hat{\theta}_l)\sqrt{n}} \frac{1}{C_2 \sqrt{n}} \pi(\hat{\theta}_l) e^{-\hat{\theta}_l \frac{\pi'(\hat{\theta}_l)}{\pi(\hat{\theta}_l)}} e^{-\frac{F''(\hat{\theta}_l)u^2}{2} \hat{\theta}_l + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right)} du. \quad (62)$$

Thus, we obtain:

$$C_2 \sim \frac{\sqrt{2\pi}}{\sqrt{nF''(\hat{\theta}_l)}} \pi(\hat{\theta}_l) e^{-\hat{\theta}_l \frac{\pi'(\hat{\theta}_l)}{\pi(\hat{\theta}_l)}}. \quad (63)$$

Of note, the gaussian limit also gives that  $\bar{\sigma}_i^2 \sim F''(\hat{\theta}_l)^{-1} N_i^{-1}$ .

Thus, we assume to develop an approximation scheme for a posterior distribution  $p_i$  asymptotically verifying:

$$p_i(\theta) \underset{N_i \rightarrow \infty}{\sim} \sqrt{\frac{1}{2\pi\bar{\sigma}_i^2}} H(\theta, \hat{\theta}_l) e^{-N_i \text{KL}(\hat{\theta}_l, \theta)}, \quad (64)$$

where  $H$  is a function accounting for the prior distribution. For the following, we take a uniform prior on  $\Theta$ , which leads to  $H(\theta, \hat{\theta}_l) = 1$ .

In the following, we will denote  $\hat{\mu}_{M_t}$  and  $\hat{\mu}_m$  as the maximum a posteriori estimates associated to their respective arms (instead of the empirical means).

## C.2 The partitioning approximation

Since all the steps leading to the partitioning approximation are independent of the type of reward distribution,  $\tilde{S}_{\text{tail}}$  and  $\tilde{S}_{\text{body}}$  have the same general form as given in Appendix A.1. Here, we consider all the distributions under  $\theta$  parameter for which we replace  $\hat{\mu}_{M_t}$ ,  $\hat{\mu}_m$  and  $\tilde{\mu}_{\text{eq}}$  by their equivalents  $\hat{\theta}_{m_t}$ ,  $\hat{\theta}_{M_t}$  and  $\tilde{\theta}_{\text{eq}}$ .

## C.3 Asymptotic intersection point

By use of Equation (64), the equation verified by the intersection point  $\bar{\theta}_{\text{eq}}$  asymptotically reads:

$$\frac{e^{-N_{M_t} \text{KL}(\hat{\theta}_{M_t}, \bar{\theta}_{\text{eq}})}}{\sqrt{2\pi\bar{\sigma}_{M_t}^2}} \int_{\mu_{\text{inf}}}^{\bar{\theta}_{\text{eq}}} \frac{e^{-N_m \text{KL}(\hat{\theta}_{m_t}, \theta')}}{\sqrt{2\pi\bar{\sigma}_m^2}} d\theta' = \frac{e^{-N_m \text{KL}(\hat{\theta}_{m_t}, \bar{\theta}_{\text{eq}})}}{\sqrt{2\pi\bar{\sigma}_m^2}} \int_{\mu_{\text{inf}}}^{\bar{\theta}_{\text{eq}}} \frac{e^{-N_{M_t} \text{KL}(\hat{\theta}_{M_t}, \theta')}}{\sqrt{2\pi\bar{\sigma}_{M_t}^2}} d\theta'. \quad (65)$$

Taking the logarithm of Equation (65) leads to

$$N_m \text{KL}(\hat{\theta}_{m_t}, \bar{\theta}_{\text{eq}}) - N_{M_t} \text{KL}(\hat{\theta}_{M_t}, \bar{\theta}_{\text{eq}}) + \frac{1}{2} \ln \frac{\bar{\sigma}_m^2}{\bar{\sigma}_{M_t}^2} + \ln \frac{\int_{\mu_{\text{inf}}}^{\bar{\theta}_{\text{eq}}} \sqrt{\bar{\sigma}_{M_t}^2} e^{-N_m \text{KL}(\hat{\theta}_{m_t}, \theta')} d\theta'}{\int_{\mu_{\text{inf}}}^{\bar{\theta}_{\text{eq}}} \sqrt{\bar{\sigma}_m^2} e^{-N_{M_t} \text{KL}(\hat{\theta}_{M_t}, \theta')} d\theta'} = 0. \quad (66)$$

Employing the same arguments as the ones exposed in Appendix A.2, we approximate  $\bar{\theta}_{\text{eq}}$  by neglecting the last term. Furthermore, in the considered asymptotic scaling regime ( $N_{M_t} \gg N_m$ ),  $\bar{\theta}_{\text{eq}}$  will be in the vicinity of  $\hat{\theta}_{M_t}$  where a Gaussian expansion of the Kullback-Leibler divergence is relevant (see Equation (61)). Thus, we approximate  $\text{KL}(\hat{\theta}_{m_t}, \bar{\theta}_{\text{eq}})$  by  $\text{KL}(\hat{\theta}_{m_t}, \hat{\theta}_{M_t})$  and expand  $\text{KL}(\hat{\theta}_{M_t}, \bar{\theta}_{\text{eq}})$  to lowest order in  $\tilde{\theta}_{\text{eq}}$  (with  $\bar{\sigma}_i^2 \sim F''(\hat{\theta}_l)^{-1} N_i^{-1}$ ), leading to:

$$\tilde{\theta}_{\text{eq}} = \hat{\theta}_{M_t} + \sqrt{2\bar{\sigma}_{M_t}^2 \left[ N_m \text{KL}(\hat{\theta}_{m_t}, \hat{\theta}_{M_t}) + \frac{1}{2} \ln \frac{\bar{\sigma}_m^2}{\bar{\sigma}_{M_t}^2} \right]}. \quad (67)$$

## C.4 Generalization of the main mode's contribution

We start by recalling the expression for the body component of the entropy:

$$\tilde{S}_{\text{body}} = - \int_{\Theta} p_{M_t}(\theta) C_m(\theta) \ln p_{M_t}(\theta) d\theta. \quad (68)$$

Without any additional information on the expression for  $\text{KL}$ , Equation (68) cannot be computed in a closed form. Thus, we will rely on the asymptotic scaling  $N_{M_t} \gg N_m \gg 1$  to provide a tractable

expression. First, we neglect variations of  $C_m(\theta)$  in Equation (68) integral by evaluating it at  $\tilde{\mu}_{\text{eq}}$ . Then, by noticing that the resulting integral is the entropy of the better empirical arm's mean, we approximate it by its leading order, proportional to  $\ln(2\pi\bar{\sigma}_{M_t}^2)/2$ :

$$\tilde{S}_{\text{body}} \approx \frac{1}{2} \ln(2\pi\bar{\sigma}_{M_t}^2) \left[ 1 - \int_{\tilde{\theta}_{\text{eq}}}^{\mu_{\text{sup}}} \frac{e^{-N_m \text{KL}(\hat{\theta}_{m_t}, \theta')}}{\sqrt{2\pi\bar{\sigma}_m^2}} d\theta' \right]. \quad (69)$$

We finally consider the last integral in Equation (69). By noticing that it is concentrated in the vicinity of  $\tilde{\mu}_{\text{eq}}$  for  $N_m \gg 1$ , we Taylor expand  $\text{KL}(\hat{\theta}_{m_t}, \theta')$  at  $\tilde{\mu}_{\text{eq}}$  to obtain:

$$\begin{aligned} \int_{\tilde{\theta}_{\text{eq}}}^{\mu_{\text{sup}}} \frac{e^{-N_m \text{KL}(\hat{\theta}_{m_t}, \theta')}}{\sqrt{2\pi\bar{\sigma}_m^2}} d\theta &\approx \frac{e^{-N_m \text{KL}(\hat{\theta}_{m_t}, \tilde{\mu}_{\text{eq}})}}{\sqrt{2\pi\bar{\sigma}_m^2}} \int_{\tilde{\theta}_{\text{eq}}}^{\mu_{\text{sup}}} e^{-N_m (\theta' - \tilde{\mu}_{\text{eq}}) \partial_2 \text{KL}(\hat{\theta}_{m_t}, \tilde{\mu}_{\text{eq}})} d\theta' \\ &\approx \frac{e^{-N_m \text{KL}(\hat{\theta}_{m_t}, \tilde{\mu}_{\text{eq}})}}{\sqrt{2\pi\bar{\sigma}_m^2} N_m \partial_2 \text{KL}(\hat{\theta}_{m_t}, \tilde{\mu}_{\text{eq}})}. \end{aligned} \quad (70)$$

Inserting Equation (70) into Equation (69) leads to the body component:

$$\tilde{S}_{\text{b}} = \frac{1}{2} \ln(2\pi\bar{\sigma}_{M_t}^2) \left[ 1 - \frac{e^{-N_m \text{KL}(\hat{\theta}_{m_t}, \tilde{\theta}_{\text{eq}})}}{\sqrt{2\pi\bar{\sigma}_m^2} N_m \partial_2 \text{KL}(\hat{\theta}_{m_t}, \tilde{\theta}_{\text{eq}})} \right]. \quad (71)$$

### C.5 Generalized expression for the entropy tail

We start by recalling the expression for the tail component:

$$\tilde{S}_{\text{tail}} = - \int_{\tilde{\theta}_{\text{eq}}}^{\mu_{\text{sup}}} p_m(\theta) \ln p_m(\theta) d\theta. \quad (72)$$

As for Equation (70), we Taylor expand  $\text{KL}(\hat{\theta}_{m_t}, \theta')$  at  $\tilde{\theta}_{\text{eq}}$  in the exponential term to obtain:

$$\tilde{S}_{\text{t}} = - \ln p_m(\tilde{\theta}_{\text{eq}}) \frac{e^{-N_m \text{KL}(\hat{\theta}_{m_t}, \tilde{\theta}_{\text{eq}})}}{\sqrt{2\pi\bar{\sigma}_m^2} N_m \partial_2 \text{KL}(\hat{\theta}_{m_t}, \tilde{\theta}_{\text{eq}})}. \quad (73)$$

Keeping the leading order of  $-\ln p_m(\tilde{\theta}_{\text{eq}}) \sim N_m \text{KL}(\hat{\theta}_{m_t}, \tilde{\theta}_{\text{eq}})$  leads to the expected tail expression used in the main text.

### C.6 Generalized form of the entropy approximation

To summarize, by combining Equations (71) and (73) we obtain an asymptotic expression for exponential family bandits with a uniform prior:

$$\tilde{S}_{\text{max}} = \frac{1}{2} \ln(2\pi\bar{\sigma}_{M_t}^2) \left[ 1 - \frac{e^{-N_m \text{KL}(\hat{\theta}_{m_t}, \tilde{\theta}_{\text{eq}})}}{N_m \partial_2 \text{KL}(\hat{\theta}_{m_t}, \tilde{\theta}_{\text{eq}}) \sqrt{2\pi\bar{\sigma}_m^2}} \right] + \frac{\text{KL}(\hat{\theta}_{m_t}, \tilde{\theta}_{\text{eq}}) e^{-N_m \text{KL}(\hat{\theta}_{m_t}, \tilde{\theta}_{\text{eq}})}}{\partial_2 \text{KL}(\hat{\theta}_{m_t}, \tilde{\theta}_{\text{eq}}) \sqrt{2\pi\bar{\sigma}_m^2}}. \quad (74)$$

Finally, depending implementation, we propose for convenience to replace in  $\tilde{S}_{\text{max}}$  the maximum a posteriori estimates of each arm by either their empirical mean, their mean posterior values or the maximum of the log-likelihood. This does not alter the algorithm's efficiency in practice, while it may simplify the implementation procedure for specific reward distributions.

Note that all these steps can be adapted to non-uniform priors (in particular by multiplying the tail by the prior effects evaluated in  $\tilde{\theta}_{\text{eq}}$ ). Finally, let us underline that our approximation scheme holds for any posterior distributions verifying Equation (64), a property we believe to be shared for more general reward distributions.

## C.7 Derivation of the increment for the closed-form expression of entropy

First, we stress there is no unique guideline to compute the expected increment of Equation (74), and multiple solutions emerge depending on the type of the reward distribution. In particular, if the reward distribution is continuous, one could integrate the increment as it has been done for Gaussian rewards above. But, if the integration cannot be solved analytically, one could approximate the increments by taking discrete reward values of the order of  $\pm\sigma$ . Similarly, if the reward takes discrete values, the increments are already discrete, but asymptotic simplifications or taking the continuous limit can also be considered.

Finally, if the increment evaluation is discrete or approximated, one could encounter rare cases where the algorithm gets trapped. It could occur when the algorithm observes a worse suboptimal arm close to the best empirical arm when it has already extensively been drawn. Because the entropy could increase drastically if an arm inversion occurs, the gradient signs may occasionally switch, leading to the failure of the minimization procedure. To prevent such cases, we change the decision procedure by maximizing the entropy variation rather than its direct minimization. An example is given for the implementation of Bernoulli rewards in the next section.

Lastly, depending on the reward distribution it may be straightforward to express the increments along the usual empirical or posterior mean as opposed to the family parameter  $\hat{\theta}$ . Often, this can be achieved through a basic variable transformation. The Bernoulli distribution example provided below serves as an illustration of this approach.

## D Numerical experiments

Here, we provide all the information regarding numerical experiments. This includes details on the numerical settings, implementation details for AIM in the investigated settings, an overview of investigated classical bandit algorithms, and additional experiments focusing on close-arm means.

### D.1 Numerical settings

In Figure 1, the posterior distributions are drawn with  $\hat{\mu}_{M_t} \approx 0.65$ ,  $N_1(t) = 374$ ,  $\hat{\mu}_m \approx 0.29$ ,  $N_m = 26$ , where  $\mu_i$  and  $N_i(t)$  are, respectively, the empirical mean and number of draws of arm  $i$  and have been obtained with the AIM algorithm.

For the Gaussian two-armed cases in Figure 2 the arm means are chosen from a uniform grid in  $(0, 1) \times (0, 1)$  using a Sobol sequence (we have avoided the values 0 and 1 but it has no impact on the obtained results). The regret is averaged over 8192 games and observed during  $10^8$  steps to attest the logarithmic scaling. For Bernoulli rewards with two-armed Figure 3, the regret is averaged over 16384 games and observed during  $10^8$  steps. It is worth noting that for Gaussian rewards, the prior information of arm means being only between 0 and 1 is not given to AIM nor to the Thompson sampling algorithm to allow a direct comparison.

For the fifty-armed case in Figures 2 and 3, the arm means are drawn from a uniform prior, and the regret is averaged over 2000 games and observed during  $10^6$  steps in Figure 2 and  $10^7$  in Figure 3.

For close arm means in Figures 4 and 5, the mean values are fixed with  $\mu_1 = 0.79$  and  $\mu_2 = 0.8$ , but this prior information is not given to the investigated algorithms. The regret is averaged over  $10^5$  games and observed during  $10^6$  steps.

For the two-armed cases in Figures 6 and 7, the arm means are chosen from a uniform grid in  $(0, 1) \times (0, 1)$  using a Sobol sequence. The regret is averaged over more than  $10^5$  games and observed during  $10^6$  steps to enhance measurement accuracy.

Finally, in the fifty-armed case in Figure 8, the arm means are drawn from a uniform prior, and the regret is averaged over  $4.10^4$  games and observed during  $5.10^4$  steps to enhance measurement accuracy.

Of note, for all the experiments, seed values are not shared throughout the algorithms. To obtain a sufficient number of runs, the code was parallelized on a cluster (asynchronously), with each run operating independently while ensuring that seed values are not common between runs. Because it



relies on an analytical expression, AIM shows an execution time of the same order of Thompson sampling (measured three times slower for two-armed Bernoulli rewards).

For completeness, an implementation of AIM for both Bernoulli and Gaussian rewards and more than two arms are given in the supplementary material (AIM Bernoulli bandits and AIM Gaussian folders).

## D.2 AIM implementation details

Here, we recap below AIM setups for the different settings evoked in the main text.

### D.2.1 Approximate information maximization for the two arm Gaussian rewards

Specifically for the two-armed case, one can simplify the expressions given in the main text.  $\tilde{\mu}_{\text{eq}}$  defined the value of  $\theta$  where both arms have the same probability of being the maximal one and reads

$$\tilde{\mu}_{\text{eq}} = \hat{\mu}_{M_t} + \frac{N_m(\hat{\mu}_{M_t} - \hat{\mu}_m)}{N_{M_t} - N_m} + \sqrt{\frac{N_{M_t}N_m(\hat{\mu}_{M_t} - \hat{\mu}_m)^2}{(N_{M_t} - N_m)^2} + \frac{\sigma^2}{N_{M_t} - N_m} \ln\left(\frac{N_{M_t}}{N_m}\right)}. \quad (75)$$

Hence, following identical approximations of the ones derived in Appendix A.1 for Gaussian reward distributions, the tail component given by Equation (7) simplifies into:

$$\tilde{S}_{\text{tail}} = \frac{1}{4} \ln\left(\frac{2\pi\sigma^2 e}{N_m}\right) \text{erfc}\left(\frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq}} - \hat{\mu}_m)}{\sqrt{2\sigma^2}}\right) + \frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq}} - \hat{\mu}_m)}{2\sqrt{2\pi\sigma^2}} e^{-\frac{N_m(\tilde{\mu}_{\text{eq}} - \hat{\mu}_m)^2}{2\sigma^2}}. \quad (76)$$

Similarly for  $\tilde{S}_{\text{body}}$ , we obtain:

$$\begin{aligned} \tilde{S}_{\text{body}} = \frac{1}{2} \ln\left(\frac{2\pi\sigma^2 e}{N_{M_t}}\right) \times & \left[1 - \frac{1}{2} \text{erfc}\left(\frac{\sqrt{N_m N_{M_t}}(\hat{\mu}_{M_t} - \hat{\mu}_m)}{\sqrt{2\sigma^2(N_m + N_{M_t})}}\right)\right] \\ & - \frac{\sqrt{N_{M_t} N_m^{3/2}}(\hat{\mu}_{M_t} - \hat{\mu}_m)}{2\sigma\sqrt{2\pi}(N_{M_t} + N_m)^{3/2}} e^{-\frac{N_m N_{M_t}(\hat{\mu}_{M_t} - \hat{\mu}_m)^2}{2\sigma^2(N_m + N_{M_t})}}. \end{aligned} \quad (77)$$

Finally, it allows us to derive the approximation of the gradient difference of the entropy for the two-armed case:

$$\begin{aligned} \Delta = \frac{1}{2} \ln\left(\frac{N_m}{N_m + 1}\right) + \frac{1}{4N_{M_t}} \text{erfc}\left(\frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq}} - \hat{\mu}_m)}{\sqrt{2\sigma^2}}\right) + \frac{\sqrt{N_m}(\tilde{\mu}_{\text{eq}} - \hat{\mu}_m)}{\sqrt{2\pi\sigma^2}} \times \\ e^{-\frac{N_m(\tilde{\mu}_{\text{eq}} - \hat{\mu}_m)^2}{2\sigma^2}} \left[ \frac{2N_{M_t}^2 - 3N_m^2}{4N_m N_{M_t}^2} + \frac{1}{4} \ln\left(\frac{N_{M_t}}{2\pi\sigma^2 e}\right) \frac{N_m^3 + N_{M_t}^2}{N_m^2 N_{M_t}^2} + \frac{(N_m^3 + N_{M_t}^2)(\tilde{\mu}_{\text{eq}} - \hat{\mu}_m)^2}{4\sigma^2 N_m N_{M_t}^2} \right], \end{aligned} \quad (78)$$

and the associated pseudo-code used for Figures 2, 4 and 6 is presented in Alg. 2 below.

### D.2.2 Approximate information maximization for Bernoulli rewards

We denote by  $\bar{\mu}_i$  the posterior mean, given by:

$$\mathbb{E}[X_{\mathcal{B}(r_i+1, N_i-r_i+1)}] = \frac{r_i + 1}{N_i + 2} = \bar{\mu}_i, \quad (79)$$

where  $r_i$  is the cumulative reward at time  $t$ ,  $N_i$  the number of draws and  $X_{\mathcal{B}(a,b)}$  follows a Beta distribution with parameters  $(a, b)$ . The variance verifies:

$$\begin{aligned} \text{Var}[X_{\mathcal{B}(r_i+1, N_i-r_i+1)}] &= \frac{r_i + 1}{N_i + 2} \frac{N_i - r_i + 1}{N_i + 2} \frac{1}{N_i + 3} \\ &= \frac{\bar{\mu}_i(1 - \bar{\mu}_i)}{N_i}, \end{aligned} \quad (80)$$

---

**Algorithm 2:** AIM Algorithm for 2 Gaussian arm

---

Draw each arm once; observe reward  $X_t(t)$  and update statistics

```
 $\hat{\mu}_t \leftarrow X_t(t), N_t \leftarrow 1 \forall t \in \{1, 2\}$   
for  $t = 3$  to  $T$  do  
  /* Arm selection */  
   $M_t \leftarrow \operatorname{argmax}_{k=1,2} \hat{\mu}_k, m \leftarrow \operatorname{argmin}_{k=1,2} \hat{\mu}_k;$   
  if  $N_{M_t} \leq N_m$  then  
     $a_t \leftarrow M_t$   
  else  
    Evaluate  $\Delta$  following Equation (9);  
    if  $\Delta < 0$  then  
       $a_t \leftarrow M_t$   
    else  
       $a_t \leftarrow m$   
  Pull  $a_t$  and observe  $X_t(a_t)$   
  /* Update statistics */  
   $\hat{\mu}_{a_t} \leftarrow \frac{\hat{\mu}_{a_t} N_{a_t} + X_t(a_t)}{N_{a_t} + 1}, N_{a_t} \leftarrow N_{a_t} + 1$ 
```

---

where  $\bar{N}_i = N_i + 3$ .

For Bernoulli rewards, we approximate the gradient as follows:

$$\begin{aligned} |\Delta_i \tilde{S}_{\max}| = & \left| \frac{\bar{\mu}_i(\bar{N}_i - 1) - 1}{\bar{N}_i - 3} \tilde{S}_{\max}\left(\frac{\bar{\mu}_i \bar{N}_i + 1 - \bar{\mu}_i}{\bar{N}_i}, \bar{N}_i + 1, \bar{\mu}_j, \bar{N}_j\right) \right. \\ & \left. + \frac{\bar{N}_i - 2 - \bar{\mu}_i(\bar{N}_i - 1)}{\bar{N}_i - 3} \tilde{S}_{\max}\left(\frac{\bar{\mu}_i(\bar{N}_i - 1)}{\bar{N}_i}, \bar{N}_i + 1, \bar{\mu}_j, \bar{N}_j\right) - \tilde{S}_{\max}(\bar{\mu}_i, \bar{N}_i, \bar{\mu}_j, \bar{N}_j) \right|, \end{aligned} \quad (81)$$

with  $\tilde{S}_{\max}$  given by Equation (74) expressed along  $\mu$  with  $\mu = \frac{e^\theta}{1+e^\theta}$ . For Bernoulli rewards the equation reads:

$$\begin{aligned} \tilde{S}_{\max}(\bar{\mu}_{M_t}, \bar{N}_{M_t}, \bar{\mu}_m, \bar{N}_m) = & \left( 1 - \frac{e^{-\bar{N}_m \operatorname{KL}(\bar{\mu}_m, \tilde{\mu}_{\text{eq}})}}{\sqrt{\bar{N}_m} \partial_2 \operatorname{KL}(\bar{\mu}_m, \tilde{\mu}_{\text{eq}}) \sqrt{2\pi \bar{\mu}_m (1 - \bar{\mu}_m)}} \right) \frac{1}{2} \ln \left( \frac{2\pi \bar{\mu}_{M_t} (1 - \bar{\mu}_{M_t})}{\bar{N}_{M_t}} \right) \\ & + \frac{\sqrt{\bar{N}_m} \operatorname{KL}(\bar{\mu}_m, \tilde{\mu}_{\text{eq}}) e^{-\bar{N}_m \operatorname{KL}(\bar{\mu}_m, \tilde{\mu}_{\text{eq}})}}{\partial_2 \operatorname{KL}(\bar{\mu}_m, \tilde{\mu}_{\text{eq}}) \sqrt{2\pi \bar{\mu}_m (1 - \bar{\mu}_m)}}, \end{aligned} \quad (82)$$

with  $\operatorname{KL}(\theta, \theta') = \theta \ln(\theta/\theta') + (1 - \theta) \ln([1 - \theta]/[1 - \theta'])$  and  $\sigma_i^2 = \frac{\bar{\mu}_i(1 - \bar{\mu}_i)}{\bar{N}_i}$ .

Briefly, the expected gradient is evaluated along arm  $i$  with a returned reward equal to 1 with probability  $\frac{\bar{\mu}_i(\bar{N}_i - 1) - 1}{\bar{N}_i - 3}$  (which is the empirical mean) or equal to 0 with probability  $1 - \frac{\bar{\mu}_i(\bar{N}_i - 1) - 1}{\bar{N}_i - 3}$ .

Of note, by adding absolute values, we seek to maximize the entropy variation rather than its direct minimization to avoid falling into an entrapment scenario (see Appendix C.7 for further discussion).

We draw some additional observations on the practical implementation of the code. First, in the gradient evaluation of  $\Delta_i$  following Equation (81) we may find a  $\tilde{\mu}_{\text{eq}}$  value to be undefined (because  $N_m > N_{M_t}$  or  $\tilde{\mu}_{\text{eq},i} > 1$ ), which is unusable for Bernoulli rewards. In this case,  $\tilde{\mu}_{\text{eq},i}$  is taken to be equal to 1, resulting in  $S_{\max} = \frac{1}{2} \ln \left( \frac{2\pi \bar{\mu}_{M_t} (1 - \bar{\mu}_{M_t})}{\bar{N}_{M_t}} \right)$ .

Second, at large times, noticing that the better empirical arm is drawn extensively, one can increment the algorithm by multiple steps at a time to speed up AIM. Indeed, let us assume that the better empirical arm is drawn  $T$  times successively while always returning a null reward, which is the worst scenario for the returned reward of the better empirical arm. Then, if the increment evaluation at  $t + T$  of Alg. 3 still returns  $M_t$ , then it ensures that all increment evaluations between  $[t, t + T]$  of

---

**Algorithm 3:** AIM Algorithm for 2 Bernoulli arm

---

Draw each arm once; observe reward  $X_t(t)$  and update statistics

```
 $\bar{\mu}_t \leftarrow \frac{X_t(t)+1}{3}, \bar{N}_t \leftarrow 4 \forall t \in \{1, 2\}$   
for  $t = 3$  to  $T$  do  
  /* Arm selection */  
   $M_t \leftarrow \operatorname{argmax}_{k=1,2} \bar{\mu}_k, m \leftarrow \operatorname{argmin}_{k=1,2} \bar{\mu}_k;$   
  if  $N_{M_t} \leq N_m$  then  
     $a_t \leftarrow M_t$   
  else  
    Evaluate  $\Delta = |\Delta_{M_t} \tilde{S}_{\max}| - |\Delta_m \tilde{S}_{\max}|$  following Equation (81);  
    if  $\Delta > 0$  then  
       $a_t \leftarrow M_t$   
    else  
       $a_t \leftarrow m$   
  Pull  $a_t$  and observe  $X_t(a_t)$   
  /* Update statistics */  
   $\bar{\mu}_{a_t} \leftarrow \frac{\bar{\mu}_{a_t}(\bar{N}_{a_t}-1)+X_t}{\bar{N}_{a_t}}, \bar{N}_{a_t} \leftarrow \bar{N}_{a_t} + 1$ 
```

---

Alg. 3 will always return  $M_t$  independently of its returned rewards. Then, using a dichotomy search on the variable  $T$ , we can diminish the number of increment evaluations of AIM at large times, thus improving AIM's performance.

### D.3 Information maximization approximation for Bernoulli rewards with more than two arms

We start by reminding the obtained entropy approximation for more than two arms:

$$\tilde{S}_{\max} = - \int_{\Theta} \left(1 - \sum_{i \neq M_t}^K [1 - C_i(\theta)]\right) p_{M_t}(\theta) \ln p_{M_t}(\theta) d\theta - \sum_{i \neq M_t}^K \int_{\tilde{\mu}_{\text{eq},i}}^{\mu_{\text{sup}}} p_i(\theta) \ln p_i(\theta) d\theta. \quad (83)$$

We first consider the increment along a worse empirical arm, which simplifies :

$$|\Delta_i \tilde{S}_{\max}| = \Delta_i \left[ - \int_{\Theta} C_i(\theta) p_{M_t}(\theta) \ln p_{M_t}(\theta) d\theta - \int_{\tilde{\mu}_{\text{eq},i}}^{\mu_{\text{sup}}} p_i(\theta) \ln p_i(\theta) d\theta \right], \quad (84)$$

which is exactly the increment evaluated in the two-armed case given in Equation (82).

Finally, we consider the increment along the better empirical arm. For simplicity, we neglect  $\tilde{\mu}_{\text{eq},i}$  variations for the increments evaluation. By use of Equation (82) we obtain

$$|\Delta_{M_t} S_{\max}| = \left| 1 - \sum_{i \neq M_t}^K \frac{e^{-\bar{N}_m \text{KL}(\bar{\mu}_i, \tilde{\mu}_{\text{eq}})}}{\sqrt{\bar{N}_m} \partial_2 \text{KL}(\bar{\mu}_i, \tilde{\mu}_{\text{eq}}) \sqrt{2\pi \bar{\mu}_i (1 - \bar{\mu}_i)}} \right| \left| \Delta_{M_t} H(\bar{\mu}_{M_t}, \bar{N}_{M_t}) \right|, \quad (85)$$

where

$$\left| \Delta_{M_t} H(\bar{\mu}_i, \bar{N}_i) \right| = \left| \frac{\bar{\mu}_i(\bar{N}_i - 1) - 1}{\bar{N}_i - 3} H\left(\frac{\bar{\mu}_i \bar{N}_i + 1 - \bar{\mu}_i}{\bar{N}_i}, \bar{N}_i + 1, \bar{\mu}_j, \bar{N}_j\right) + \frac{\bar{N}_i - 2 - \bar{\mu}_i(\bar{N}_i - 1)}{\bar{N}_i - 3} H\left(\frac{\bar{\mu}_i(\bar{N}_i - 1)}{\bar{N}_i}, \bar{N}_i + 1, \bar{\mu}_j, \bar{N}_j\right) - H(\bar{\mu}_i, \bar{N}_i, \bar{\mu}_j, \bar{N}_j) \right|, \quad (86)$$

with  $H(\bar{\mu}_{M_t}, \bar{N}_{M_t}) = \frac{1}{2} \ln \left( \frac{2\pi\bar{\mu}_{M_t}(1-\bar{\mu}_{M_t})}{\bar{N}_{M_t}} \right)$ .

---

**Algorithm 4:** AIM Algorithm for  $K > 2$  Bernoulli arm

---

Draw each arm once; observe reward  $X_t(t)$  and update statistics

```

 $\bar{\mu}_t \leftarrow \frac{X_t(t)+1}{3}, \bar{N}_t \leftarrow 4 \forall t \in \{1, \dots, K\}$ 
for  $t = K + 1$  to  $T$  do
  /* Arm selection */
   $M_t \leftarrow \operatorname{argmax}_{k=\{1, \dots, K\}} \bar{\mu}_k$ ; Evaluate  $\Delta_{M_t} \tilde{S}_{\max}$  following Equation (85);
  Evaluate  $m = \operatorname{argmax}_i (\Delta_i |\tilde{S}_{\max}|, i \neq M_t)$  with  $\Delta_i |\tilde{S}_{\max}|$  following Equation (81);
  if  $\Delta_{M_t} |\tilde{S}_{\max}| > \Delta_m |\tilde{S}_{\max}|$  then
     $a_t \leftarrow M_t$ 
  else
     $a_t \leftarrow m$ 
  Pull  $a_t$  and observe  $X_t(a_t)$ 
  /* Update statistics */
   $\bar{\mu}_{a_t} \leftarrow \frac{\bar{\mu}_{a_t}(\bar{N}_{a_t}-1)+X_t}{\bar{N}_{a_t}}, \bar{N}_{a_t} \leftarrow \bar{N}_{a_t} + 1$ 

```

---

Of note in the gradient evaluation of  $\Delta_i$  following Equation (81), if one finds a  $\tilde{\mu}_{\text{eq},i}$  value undefined (because  $N_m > N_{M_t}$  or  $\tilde{\mu}_{\text{eq},i} > 1$  which is unusable for Bernoulli reward), then,  $\mu_{\text{eq},i}$  is taken to be equal to 1 resulting in  $S_{\max} = \frac{1}{2} \ln \left( \frac{2\pi\bar{\mu}_{M_t}(1-\bar{\mu}_{M_t})}{\bar{N}_{M_t}} \right)$ . Finally, if  $M_t \leftarrow \operatorname{argmax}_{k=\{1, \dots, K\}} \bar{\mu}_k$  has multiple solutions, we suggest choosing the one displaying the lowest number of draws.

#### D.4 Overview of baseline bandit algorithms

Here, we briefly review several baseline algorithms and their chosen parameters to provide a benchmark of our information maximization method.

##### D.4.1 UCB-Tuned

This algorithm falls under the category of upper confidence bound (UCB) algorithms, which select the arm maximizing a proxy function typically defined as  $F_i = \hat{\mu}_i + B_i$ . For UCB-tuned,  $B_i$  is given by:

$$R_i = c(\mu_1, \mu_2) \sqrt{\frac{\ln(t)}{N_i(t)} \min \left( \frac{1}{4}, s_i(t) \right)}, \quad s_i(t) = \hat{\sigma}_i^2 + \sqrt{\frac{2 \ln(t)}{N_i(t)}}, \quad (87)$$

where  $\hat{\sigma}_i^2$  is the reward variance and  $c$  a hyperparameter. For Gaussian rewards, by testing various  $c$  values for uniform priors in Equation (87), we end up with  $c = 2.1$  and  $\hat{\sigma}_i^2 = \frac{\sigma^2}{N_i(t)}$ .

##### D.4.2 KL-UCB

This algorithm is another variant of the upper confidence bound (UCB) class specifically designed for bounded rewards. In particular, it is known to be optimal for Bernoulli distributed rewards [Garivier and Cappé, 2011, Cappé et al., 2013]. For KL-UCB,  $F_i$  is expressed as follows:

$$F_i = \max \left\{ \theta \in \Theta : N_i(t) \text{KL} \left( \frac{r_i(t)}{N_i(t)}, \theta \right) \leq \ln(t) + c(\mu_1, \mu_2) \ln(\ln(t)) \right\}, \quad (88)$$

where  $\Theta$  denotes the definition interval of the posterior distribution. By testing various  $c$  values for uniform priors, we end up with  $c(\mu_1, \mu_2) = 0.00001$  ( $c = 0$  for the 50-armed Gaussian setting and  $c = 10^{-6}$  for the 2-armed Bernoulli setting). Of note, the maximum is found using a dichotomy method using a precision of  $10^{-5}$  and a maximum number of iterations of 50.

For KL-UCB++ [Ménard and Garivier, 2017], the function  $F_i$  is expressed as follows

$$F_i = \max \left\{ \theta \in \Theta : N_i(t) \text{KL} \left( \frac{r_i(t)}{N_i(t)}, \theta \right) \leq \ln_+ \left( \frac{T}{KN_i(t)} \ln_+^2 \left( \frac{T}{KN_i(t)} \right) + 1 \right) \right\}, \quad (89)$$

where  $\ln_+(x) = \max(\ln(x), 0)$ , and  $T$  is the stopping time of the bandit game. Therefore, KLUCB++ is not an anytime algorithm, but it still underperforms when compared to AIM and Thompson sampling.

#### D.4.3 Thompson sampling

At each step, Thompson sampling [Thompson, 1933, Kaufmann et al., 2012a,b] selects an arm at random, based on the posterior probability that maximizes the expected reward. In practice, it draws  $K$  random values according to each arm’s mean posterior distribution and selects the arm with the highest sampled value as:

$$a_t = \operatorname{argmax}_{i=1..K} \left( Z_i(\hat{\mu}_i(t), N_i(t)) \right), \quad (90)$$

where  $Z_i(t)$  is drawn according to the posterior distribution of the  $i$ th arm’s mean. Here, we used a uniform prior on  $[0, 1]$  for Bernoulli rewards and a uniform prior on  $\mathbb{R}$  for Gaussian rewards to provide a direct comparison with AIM.

Finally, for Thompson sampling plus [Jin et al., 2022], denoted TS+, each sampled value for the comparison is drawn according to  $Z_i(\hat{\mu}_i(t), N_i(t))$  with a probability  $1/K$  or taken equal to  $\hat{\mu}_i(t)$ , otherwise.

#### D.4.4 MED

At each step, the minimal empirical divergence (MED) algorithm [Honda and Takemura, 2011], selects an arm at random, based on a tailored distribution building on the Kullback-Leibler distance to the better empirical arm. In practice, the arm  $a_t = i$  is drawn with a probability:

$$p_i = \frac{\exp \left[ -N_i(t) \text{KL} \left( \frac{r_i(t)}{N_i(t)}, \bar{\mu}_{M_t} \right) \right]}{\sum_{j=0}^K \exp \left[ -N_j(t) \text{KL} \left( \frac{r_j(t)}{N_j(t)}, \bar{\mu}_{M_t} \right) \right]}. \quad (91)$$

### D.5 Additional experiments

#### D.5.1 Approximate information maximization for Gaussian rewards and close arms

For completeness, we provide in Figure 4 below regret performances in which the arms’ mean values are close ( $\Delta\mu = 0.01$ ), and are thus difficult to distinguish, for Gaussian reward distributions. Here, AIM shows state-of-the-art performance comparable to Thompson sampling, even outperforming it at longer times.

#### D.5.2 Approximate information maximization for Bernoulli rewards and close arms

For completeness, we provide in Figure 5 below regret performances in which the arms mean value are close ( $\Delta\mu = 0.01$ ) for Bernoulli reward distributions. As for Gaussian rewards, AIM shows state-of-the-art performance comparable to Thompson sampling even when arms mean rewards are difficult to distinguish.

#### D.5.3 Investigating approximate information maximization for large simulation volumes

To refine the numerical investigation of AIM’s regret performance, we replicate the experiments of the main text using a larger volume of simulations (but on shorter timescales). For the 2-armed bandit with Gaussian and Bernoulli rewards, the regret performance under a uniform prior is averaged over more than  $10^5$  runs. Similarly, for the 50-armed bandit with Bernoulli rewards, the regret is averaged over  $4 \times 10^4$ . This leads to the results shown in Figures 6 to 8, confirming the results of Figures 2 to 3.

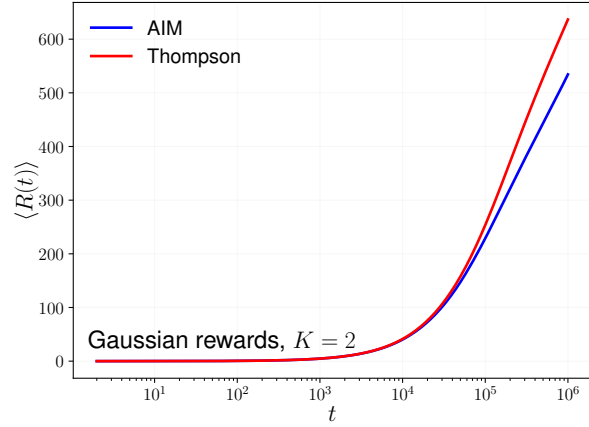


Figure 4: Temporal evolution of the regret for 2-armed bandit with Gaussian rewards ( $\sigma = 1$ ) for close mean parameters. In blue AIM, in red Thompson sampling. Arm mean reward values are fixed with  $\mu_1 = 0.8$  and  $\mu_2 = 0.79$ , the regret is obtained by averaging over  $10^5$  realizations.

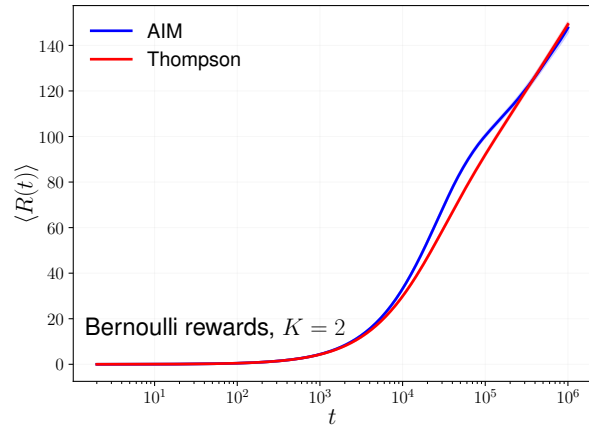


Figure 5: Temporal evolution of the regret for 2-armed bandit with Bernoulli rewards for close mean parameters. In blue AIM, in red Thompson sampling. Arm mean reward values are fixed with  $\mu_1 = 0.8$  and  $\mu_2 = 0.79$ , the regret is obtained by averaging over  $10^5$  realizations. Confidence intervals shows the standard deviation.

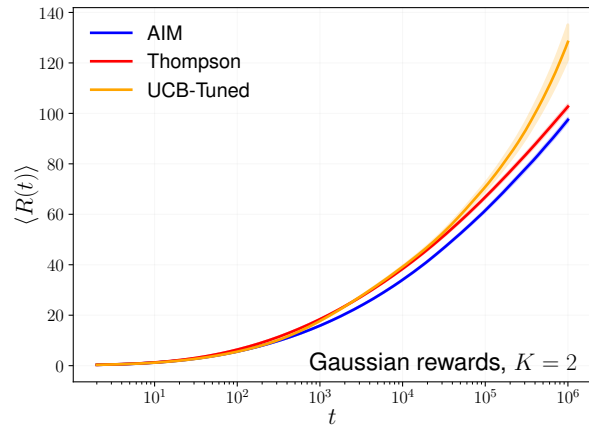


Figure 6: Evolution of the Bayesian regret for 2-armed bandit with Gaussian rewards under a uniform mean prior. The regret is averaged over more than  $10^5$  runs. Confidence intervals shows the standard deviation. Confidence intervals show the standard deviation.

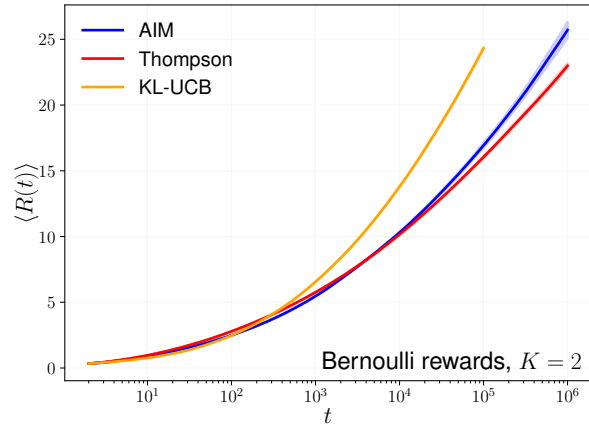


Figure 7: Evolution of the Bayesian regret for 2-armed bandit with Bernoulli rewards under a uniform mean prior. The regret is averaged over more than  $10^5$  runs. Confidence intervals show the standard deviation.

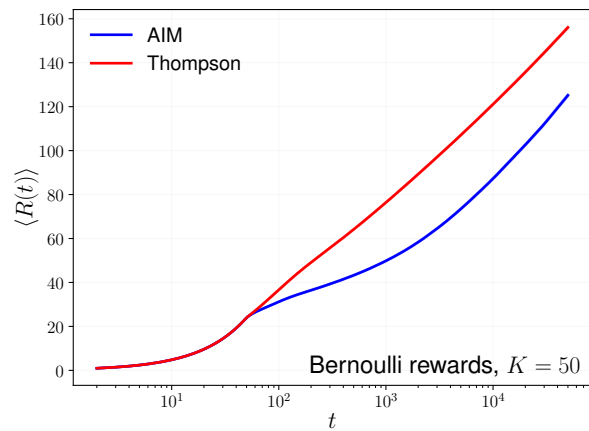


Figure 8: Evolution of the Bayesian regret for 50-armed bandit with Bernoulli rewards under a uniform mean prior. The regret is averaged over  $4 \times 10^4$  runs. Confidence intervals show the standard deviation.