

1 We thank the reviewers for all their remarks and comments. We address these remarks separately below.

2 **Reviewer 1:** Our bounds indeed require the horizon T to be bigger than the number of arms K . We study an
 3 intermediate confidence regime which is non-asymptotic but in which the algorithm still gets to access several samples
 4 per arm. The setting where the algorithm should not even pull each arm once (due to a strong structure and low
 5 confidence level) is a challenge that we do not address.

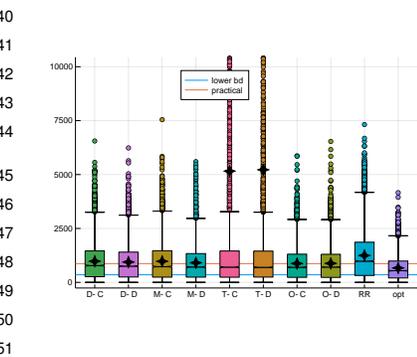
6 You mention that the lower bounds we use are difficult to compare to the bounds obtained in other works and we agree
 7 that their non-explicit nature can have that effect. For best-arm with Gaussian arms with variance σ^2 , it is shown in [13]
 8 that the lower bound T^* such that $\liminf_{\delta \rightarrow 0} \frac{\mathbb{E}[\tau_\delta]}{\log(1/\delta)} \geq T^*$ verifies $\sum_{k \neq *} \frac{2\sigma^2}{\Delta_k^2} \leq T^* \leq 2 \sum_{k \neq *} \frac{2\sigma^2}{\Delta_k^2}$. On the upper
 9 bound side, we can divide the existing work in two categories: papers that get only asymptotic optimality, like [13,
 10 20, 24]; and papers that get finite confidence guarantees but not asymptotically optimal (with the right multiplicative
 11 factors), like [5, 16]. In the second case, more effort is spent towards making the other terms of the bound small. For
 12 the Gaussian top-k arms problem, the difference between the lower bound we use and the gap-based one found for
 13 example in [Chen et al., Nearly Instance Optimal Sample Complexity Bounds for Top-k Arm Selection, 2017] is also a
 14 multiplicative factor. For non-Gaussian arms, the difference can be much larger. See [13, section 2.2] for a discussion.

15 **Reviewer 2:** We disagree with your assessment of the contributions of the paper, expressed in that sentence: “This
 16 paper adds a link between the principle of optimism in the face of uncertainty (widely used in regret minimization
 17 setting) and pure exploration problems, and proposes algorithms based on this principle.” Our design exploits a game
 18 point of view through a combination of two adversarial algorithms playing against each other, one of which deals with
 19 uncertainty by using optimism. The optimism is more akin to a trick than the main feature. The main contribution is
 20 to show how the interaction of the two players iteratively solve the lower bound problem to tend towards the optimal
 21 sampling behaviour (and complexity). No other stochastic bandit paper uses a similar design. The reference [9] notes
 22 that the lower bound can be seen as a game and uses that observation to derive lower bounds, but that paper does not
 23 use this hindsight to form algorithms (their algorithm is a slightly modified track-and-stop).

24 For the experimental part, we are inclined towards keeping it fairly light in the main part of the paper due to space
 25 constraints. We can however add comparisons to more algorithms in the appendix. As our algorithm can deal with
 26 any exploration problem on which one of the required oracles can be implemented, and not only best-arm problems,
 27 we compared to similar general algorithms. We will include the algorithm of <https://arxiv.org/abs/1602.08448> in the
 28 experiments of appendix G.

29 We agree that the computational complexity claims would be better supported with data, which we will include. To give
 30 a quick idea: on best-arm identification (Figure 1a, section 4), taking the time per iteration of uniform sampling as
 31 1, the algorithms D, M, T and O have iteration times of 6.4, 3.8, 121 and 526 respectively. Our new algorithm D has
 32 same order of computational complexity as M [24] up to an overhead due to the computation of the optimism terms.
 33 Track-and-Stop (T) is slower and our new algorithm O is again slower, as predicted (it uses a more complicated oracle).
 34 On the thresholding task, all oracles are closed form and all algorithms have similar iteration times (6.1, 5.1, 7.1, 5.2).

35 **Reviewer 3:** We doubt that a meaningful moderate confidence bound can be derived from [13], for reasonable
 36 confidence δ . Their continuity based argument introduces a quantity $\varepsilon > 0$ and the forced exploration leads to a bound
 37 that depends on ε and the parameters of the problem through polynomial terms with high exponents (for example $1/\varepsilon^4$).
 38 ε must then be chosen as a function of δ to obtain a moderate confidence bound. Our best efforts lead to a bound with a
 39 lower order term proportional to $(\log(1/\delta))^{11/12}$, which is $o(\log(1/\delta))$ but potentially still big for moderate confidence.



40 We evaluated track-and-stop (TaS) experimentally in the paper, and the exper-
 41 iments (in particular Figure 1b and Figures 5) show that while TaS is asymptotically
 42 optimal, the empirical performance may be poor, even for $\delta = e^{-20}$.
 43 On the left, we reproduced Figure 1b with larger markers for the mean (crosses).
 44 TaS (called T-C and T-D) performs much worse than even uniform sampling.

45 About [Chen et al., 2017]: we indeed should have discussed it since it addresses
 46 the same problem. We now discuss it and will include this comparison in the
 47 paper. Their algorithm has several features in common with TaS: a phase of
 48 forced exploration restricts the candidates answers to only 1 (through the use
 49 of a confidence region) and forms estimates of the parameters; the algorithm
 50 then verifies that it is indeed the correct answer up to the required confidence
 51 level by using a plug-in estimate of the solution to the lower bound problem.

52 Their definition of “optimal” is not the same as ours. Our algorithm verifies $\lim_{\delta \rightarrow 0} \mathbb{E}[\tau_\delta] / \log(1/\delta) = T^*$, where T^*
 53 is the lower bound complexity. They show a $256T^*$ bound. Their algorithm’s complexity is optimal in the sense that it
 54 is proportional to T^* . Ours matches the lower bound with the right multiplicative constant.