

1 We thank the reviewers for their insightful comments. The numbered citations refer to references in the submitted paper.
2 The remaining additional citations are listed at the end of the page.

3 **Reviewers 3 and 4: motivation.** We admit the comment of **Reviewer 4** about the lack of motivation on *the combined*
4 *scenario of nonstochastic bandits with corruption*. The isolated motivation about bandits with corruption is mentioned
5 in 2nd paragraph of the introduction (L25-35) with extensive recent results [8,9,11,14,15,17,18,25]. The motivation
6 about nonstochastic bandits is given in the 3rd paragraph (L36-43). There exist application scenarios where either
7 ‘nonstochastic bandits’ or ‘bandits with corruption’ in isolation fail to fully characterize the underlying model. To
8 response *the combined motivation* and strengthen the motivation, we will add one such example, described below, to
9 the paper. This will also address **Reviewer 3**’s comments on *the real example of targeted attacks* and **Reviewer 4**’s
10 comments on attacks to the *actual vs. observation of the reward*.

11 A concrete example of nonstochastic bandits with corruptions is the Online Shortest Path Routing (OSPR) problem
12 under the denial of service (DoS) attacks. OSPR is a classic example of MAB problems [20]. And there is also
13 extensive research on routing under DoS attacks, including the recent work [Zhou et al., 2019] focusing on bandit
14 modeling of this scenario. OSPR could be reasonably modeled as nonstochastic bandits when the delays on the links
15 change dynamically [György et al. 2007], or once is it difficult to characterize the combined distribution of a path
16 including multiple links [20]. In this nonstochastic scenario, the DoS attack could be modeled by our bandit with
17 targeted corruptions. Specifically, the DoS attacker can be aware of the selected paths by detecting the transmitted
18 packets over the path and manipulate the latency of the selected path by flooding the path with dummy packets. Also,
19 the budget of the attacker is simply the available resources for the DoS attacker to keep her undetectable. Arguably,
20 none of ‘nonstochastic bandits’ and ‘bandits with corruption’ alone would suffice to fully characterize the underlying
21 model here. We believe presenting this example in the introduction would prove useful to connect the dots between
22 *nonstochastic bandits* and *bandits with corruption*.

23 **Reviewer 2: experiments.** We will add the following experiment to the supplementary material.
24 The goal here is to compare ExpRb with Exp3. We constructed a simple scenario where
25 the attacker follows a simple policy that attacks the optimal arms (see L172-181 of the
26 paper) with 1 high-reward and $K - 1$ low-reward arms. In Fig. 1, we report the average
27 regret of 100 independent runs, with $\Phi = O(\sqrt{T})$. The results show that the regret of
28 Exp3 is largely degraded with the attack, while ExpRb achieves a sublinear regret. This
29 experiment is not meant to be exhaustive, rather it is intended to validate the theoretical
30 results and illustrate the potential of our approach.

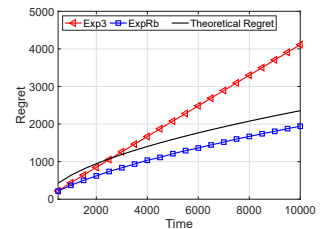


Figure 1: ExpRb vs. Exp3

31 **Reviewers 3 and 4: lower bound. R3/Q1:** Different bandit algorithms tolerate
32 different levels of corruption. Hence, finding a more refined bound for the attacker’s
33 budget is highly algorithm-specific and we were not able to generalize our result in
34 Corollary 2 for algorithms beyond Exp3. **R3/Q2:** Without substantial change in the proofs, the lower bound in
35 Theorem 1 can be applied to the algorithms that are aware of the existence of corruptions, but do not know the budget.
36 We will add a remark after Theorem 1 to involve this case. **R4 (the impossibility Theorem 1):** The impossibility result
37 is a generic result for any attack-agnostic algorithms; one can have more refined versions for specific algorithms such
38 as our result for Exp3 (see **R1/Q1**). **R4: the algorithm needs to know the attacker’s budget):** Yes, it is a negative
39 result on the attack-agnostic model. ExpRb algorithm is parameterized by a robustness parameter γ . Our further results
40 which will appear in the future version can characterize the regret as a function demonstrating the regret reduction with
41 improper γ . As future work, it is promising to dig out more interesting results on the attack-agnostic case.

42 **Additional comments: Reviewer 1:** Both questions are valid and due to a mistake when defining \mathcal{T}_i , $t_i(n)$ and N_i in
43 the proof of Lemma 8. Thanks for your careful reading. \mathcal{T}_i should be referred to (Do you simply mean: " \mathcal{T}_i denotes")
44 the set of time slots that the i -th arm is selected and the selection probability for arm i is lower than the previous
45 one (only in this case, $\delta(t)$ will be larger than 0). And accordingly, $t_i(n)$, $n = 1, 2, \dots, N_i$ are the indices for those
46 time slots. By redefining those variables, we hope we can clarify the reviewer’s concern as follows. **R1/Q1:** We only
47 consider the time slots that the selection probability for the i -th arm is smaller than the previous, so the algorithm falls
48 into the case in Line (5) of algorithm 1, but not case (1) in line 522. By checking the conditions for Equation (8), we
49 have $\tilde{p}_i(t_i(n)) = p_i(t_i(n))$. **R1/Q2:** Yes, we do require that condition. This inequality holds, since we only consider
50 the time slots that the selection probability for the i -th arm is smaller than the previous one.

51 Additional References

52 [Zhou et al., 2019] Zhou, P., Xu, J., Wang, W., et al. (2019). Toward Optimal Adaptive Online Shortest Path Routing
53 With Acceleration Under Jamming Attack. *IEEE/ACM Transactions on Networking*, 27(5), 1815-1829.

54 [György et al. 2007] György, A., Linder, T., Lugosi, G., and Ottucsák, G. (2007). The on-line shortest path problem
55 under partial monitoring. *Journal of Machine Learning Research*, 8(Oct), 2369-2403.