

1 **Reviewer 1:** Thank you for your constructive feedback. We hope our response below addresses your concerns:

2 • Indeed the constant  $\alpha$  in the CR bound could exponentially increase as  $p$  increases, especially for marginally-stable  
3 systems and highly unstable systems. But we believe it is mainly from the system’s property, instead of our algorithm’s  
4 limitation. In other words, marginally stable or highly unstable  $A$  matrices are inherently harder to achieve small CR.  
5 As you mentioned, any online policy will suffer  $\alpha^2$  in the competitive ratio (Appendix E). It is interesting to understand  
6 which system will make online control harder (we partially discussed that in Appendix B, focusing on 1-D systems),  
7 and we will add more discussion about this point in the main paper.

8 • In online learning, it is common to have  $q_{\min}$  and  $q_{\max}$  in the CR or regret bounds (e.g., [23, 31]). We also want to  
9 point out that even if  $q$  is fixed, bounding competitive ratio is non-trivial, and often it is unbounded [7].

10 • We agree our control setting is a subset of general LQR. We emphasize that our goal in the control part is to show the  
11 possibility for a policy to be competitive to the true, dynamic optimal offline cost, something not shown before. The  
12 only prior result is [23], which uses invertible  $B$  and known  $w_t$  at time  $t$ . This goal is important because all other prior  
13 works focus on achieving small regret compared to the best linear policy. Given our example in Appendix B, one may  
14 wonder if it is even possible to match the offline optimal. We give a positive answer. Though our setting is not fully  
15 general, it strictly generalized the prior art. We see this as a significant step towards competitive control.

16 • The most significant difference between our approach and the classical robust control (e.g.,  $H_\infty$ ) is that  $H_\infty$  is neither  
17 online nor adaptive, i.e., the policy will not change even if some  $w_t$  is not adversarial. Our framework is naturally  
18 adaptive from the estimation set  $W_t$  (i.e., better estimation leads to more aggressive policy). We will add this discussion.

19 • Being optimistic is a common and powerful heuristic in general online learning, especially in regret minimization  
20 (e.g., UCB in multi-armed bandit, efficient Q-learning). But our optimistic idea is quite novel in SOCO, naturally  
21 because previous settings focus on precise information cases and then there is no need for being optimistic.

22 **Reviewer 2:** Thank you for your constructive feedback. We hope our response below addresses your concerns:

23 • We feel our assumption on strongly convexity and smoothness is not a significant weakness because: (1) Assuming  
24 strongly convexity and smoothness are very common in the online learning and optimization community [30, 31]; (2)  
25 Strong smoothness is not needed in [24] because the agent has the perfect prediction of the next hitting cost function, but  
26 it is critical in our setting where the prediction is imperfect, because the competitive ratio can be unbounded otherwise  
27 (e.g., consider the case when the hitting cost is an indicator function).

28 • Our setting strictly generalizes [23], where  $B$  is invertible and  $w_t$  is perfectly known at step  $t$ . To address the clarity  
29 issues, we will add more background and make notations easier to follow (e.g., add a notation list).

30 • For technical questions: (1) In line 207, we need to make an additional assumption that  $\alpha$  is large (strictly speaking,  
31 larger than a constant  $c$  such that  $c > 1$ ). In this case,  $\lambda$  must be in the order of  $O(m)$ . (2) In line 153, when  $m$  tends to  
32 zero, the competitive ratio is of order  $O(1/\sqrt{m})$  when  $\alpha = 1$ , and  $O(1/m)$  if  $\alpha > 1$ . Hence we report  $O(1/m)$  which  
33 is more conservative. (3) In line 150, it is not necessary to assume  $\alpha^2 < m + 1$ , because all assumptions in line 150  
34 hold if we let  $\lambda_2 = 0$  and  $\lambda_1 = 2m/(m + \alpha^2 - 1 + \sqrt{(m + \alpha^2 - 1)^2 + 4m})$ . We will discuss them in the revision.

35 **Reviewer 3:** Thank you for your constructive feedback. We hope our response below addresses your concerns:

36 • Previous competitive ratio results in SOCO focus on the setting when both the geometry and minimizer of  $f_t$  are  
37 revealed before the agent picks  $y_t$  [16, 23, 24]. In contrast, the minimizer is unknown when picking  $y_t$  in our setting.  
38 Our setting generalizes the previous ones, and is practical in many cases. For example, in power systems the geometry  
39 is governed by network topology (usually revealed before decision making) and minimizer is decided by users (which  
40 could be revealed afterwards). When reduced to control, the geometry is from cost functions and the minimizer is from  
41 adversarial disturbance. It is why we need  $p$  steps of future costs, but don’t need any future disturbances. We want to  
42 emphasize (1) the access to future cost functions is common in control if the focus is on disturbances in dynamics (e.g.,  
43 in LQ tracking problem the cost functions are pre-given) and (2) the only existing competitive policy [23] needs both  
44 future costs and disturbances, and we show a competitive policy exists even if disturbance is unknown in advance.

45 • Note that the cost bound in our Theorem 2 is  $C_1 \cdot \text{cost}(\text{opt}) + (a + b - d) \sum_t \|v_t - \tilde{v}_t\|^2$ , and we can get a pure  
46 competitive ratio  $C_2 \cdot \text{cost}(\text{opt})$  because we have a “ $-d$ ” term before the path length  $\sum_t \|v_t - \tilde{v}_t\|^2$ . We name it  
47 “beyond worst case” because the first bound will be tighter if the estimation is more precise. But note that we can  
48 always get a constant-competitive result even if the estimation is totally off, and  $C_2$  does not depend on  $\sum_t \|v_t - \tilde{v}_t\|^2$ .  
49 Our numerical examples in Appendix C provide the example: if  $w_t$  is smooth (i.e.,  $\|w_{t+1} - w_t\|$  is small), then the  
50 path length is small and the first bound gets tighter (Fig. 1(b,d)). However, if the estimation set is large, the path length  
51 may get large and then the second bound would be better.

52 • For technical questions, we think the definition of  $g_t$  on page 14 is correct because the switching cost has the  
53 coefficient of  $1/2$  by definition (between line 109 and 110). In Eq. (7), we should change  $t$  to  $i$  and  $q$  to  $p$ .

54 **Reviewer 4:** Thank you for the feedback on paper organization and presentation. We would move our examples in the  
55 main body (we will have one more page in the revision), add a notion table and restate all claims before the proofs.