We thank the reviewers for their effort! Below we answer some of their questions.

General: we want to clarify that our results are purely theoretical, resolve an important question in exploration for control, and introduce techniques that we hope other researchers will build upon for practical applications.

**(R2) (necessity of geometric exploration)**
The value of the geometric exploration is that it achieves *much* better dependence in the dimension (not a constant factor!). Suppose $d=d_x=d_u$. The poly(d) in Theorem 2 (bandit case) is actually $d^{30}$ (!), while in Theorem 1 (geometric exploration) we have just $d^3$.
This is a major improvement over the bandit algorithm.
Since this was the *only* weakness presented in your review (other than the complaint on the organization of the Appendix), please consider updating your score.

**(R3) (relevance)**
1) We believe that our paper is highly relevant, since it tackles an important problem at the intersection of control theory and online learning, and follows up on numerous other works. As evidence, all the other reviewers said that this work has definite relevance.
2) If you cannot assess the value of a theoretical paper, please consider not reviewing it. Science is a long-term endeavor. It took millennia for number theory to evolve from a leisurely activity of the ancient Greeks to the basis of modern cryptography.

**(R1) (extensions of the result: best LDC, more general stochastic noise, partial observability)**
Yes, we can achieve all these extensions. We did not mention them in the paper, because we wanted to highlight our novel theoretical ideas. We will add one more section where we explain how all these extensions can be achieved.
1) the disturbance-based policies can express all stabilizing LDCs with internal state, so our regret bound actually holds wrt this richer policy class.
2) We can deal with any stochastic bounded disturbance (we need boundness for Lipschitz concentration to hold). The only place where the assumption of gaussian disturbances really helps is that given some policy M and matrices A, B, we can compute offline (to a very good approximation) the stationary cost C(M|A,B), because we know the disturbance distribution. However, even when we don't, we can still use the estimated disturbances (\hat{w}) as samples to approximate this expectation (i.e., the stationary cost).
3) The extension to partial observation is tedious but straightforward and uses the idea of "nature's y's", exactly as in [20].

**(R1) (poly time)**
Yes, we can prove polynomial time. We considered the implementation of the linear optimization oracle via the ellipsoid method to be folklore. We will add a formal proof at the Appendix.

**(R4) (dimension dependence running time)**
The overall dimension dependence in the runtime is $(d_x*d_u)^7$. We did not put effort into making this algorithm practical.

**(R1) (need for stabilizing controller)**
This is a standard assumption in online control. Without this assumption, the regret has an additive term which is exponential in the dimension (see the recent paper "Black box control for LDS" by Chen and Hazan).

**(R4) (Barycentric spanners, d vs d+1 elements)**
We use *affine* barycentric spanners and an affine basis has d+1 elements. We need this because the state is an affine function of the policies (not just linear).

Finally, we would like to thank reviewer 4 for their suggestions on the background section, which we will incorporate in the paper.