1 Thank you to all reviewers! The response to reviewer 2 may be most important to how the paper is interpreted:

2 **Reviewer 2:** *while the principle is certainly common in the literature, this is the first paper to demonstrate frequentist*
3 *regret guarantees for perturbation induced exploration. Proof techniques do not appear to be terribly different than the*
4 *prior arm, with the key differences appearing in Lemmas 4-5.. . . . . . However, I do not know if the fact that randomized*
5 *value functions do in fact ensure frequentist regret guarantees is terribly surprising.*

6 Managing to give such a clean, seemingly straightforward, frequentist analysis of RLSVI seems to be a major
7 contribution over prior art. The journal paper Osband et al. [2017] develops a theory of recursive stochastic-dominance
8 relations to study the algorithm, hence requiring very different techniques than the rest of the RL literature. The paper on
9 frequentist analysis of posterior sampling by Agrawal and Jia [2017] builds on those stochastic-dominance techniques,
10 is immensely technical, requires modifying Thompson sampling to get the proof to work, and contains a critical flaw in
11 the proof currently posted online. I've worked very hard to uncover a new proof that hopefully makes it easy for future
12 researchers to transfer results known for optimistic algorithms over to randomized value function approaches.

13 On whether the results are surprising: Strong theory sometimes takes years to develop and in the meantime people
14 can start to get used to the main ideas. This paper tries to provide some backing for the claim that "Training a
15 value function estimation scheme on noise-perturbed data generates a highly sophisticated form of exploration that is
16 fundamentally quite different from what is generated by employing stochastic policies." It has been just three years
17 since the first paper making such a claim was published [Osband et al., 2016]. To my understanding, this claim was
18 often met with skepticism, especially because it lacked a frequentist regret bound to back it up. Such a bound has
19 been elusive since then. Things can seem quite clear in hindsight, but I think the claims in this paper would have been
20 shocking 7 years ago (before any analysis of Thompson sampling in bandits even exited.) That's noteworthy, since
21 we're studying extremely old questions in sequential decision making.

22 **Responses to Reviewer 1:** [Paraphrasing] *(1) Would a similar analysis yield a high-probability regret bound? . . . (2)*
23 *While the proof for Lemma 4 & 5 is described very well in the main text, it would be helpful to have a short explanation*
24 *how this is used to achieve Lemma 6.. . . (3) How does this setting for $\beta$ affect the empirical performance of the algorithm.*
25 *. . . (4) The authors chose the setting with time-dependent dynamics which is a little less common than the default setting*
26 *where dynamics and rewards are identically distributed across time steps within the episode. . .*

27 (1) Thanks so much! I now believe the high probability bounds work out. Effectively, the techniques in this paper bound,
28 with high probability, the conditional expected regret $\mathbb{E}[V(M, \pi^*) - V(M, \pi^k) \mid \mathcal{H}^k] \le B_k$ by some simple terms $B_k$
29 whose sum we know how to control. Rather than take an expectation, the Azuma-Hoeffdin'g inequality should bound
30 the sum of martingale differences: $\sum_{k=1}^{K} \left( \left(V(M, \pi^*) - V(M, \pi^k)\right) - \mathbb{E}[V(M, \pi^*) - V(M, \pi^k) \mid \mathcal{H}_{k-1}] \right)$ The main
31 challenge seems to be that $B_k$ involves terms that are not uniformly bounded, since we add Gaussian noise. The
32 standard ways in which researchers bounds things like $B_1 + \cdots + B_K$ could get unusually messy as a result.

33 (2-4) In the revision, I'll explain more about Lemma 6 and clearly highlight open questions regrading points (3) and (4).
34 I think setting $\beta$ too large is similar to setting overly large optimism bonuses for optimistic algorithms. For practical
35 applications, I suspect the noise variances one adds should also be adaptive to the data. We're essentially (recursively)
36 applying linear-regression to minimize the Bellman residuals. One can write down variants of RLSVI that calibrate the
37 noise they add to the variance of observed residuals, but I'd like to do careful empirical evaluation. I'll try to work
38 through the time-homogeneous analysis. The challenge is that mis-estimation of a single state could lead to error in
39 every Bellman update. One needs to be careful to avoid introducing even more factors of $H$.

40 **Responses to Reviewer 3:** *In response to the comments on bootstrapping and the correct form of prior randomness .*

41 Thanks, this really gets to the crux of what makes boostrap-like methods for exploration so different from the treatment
42 statistics books. You're right that in initial periods, RLSVI is really no better than uniform exploration, though perhaps
43 the same comment applies to optimistic algorithms. But my understanding is that, no matter what, any algorithm will
44 visit certain states/actions many times. These become well understood and the variance of the injected noise begins to
45 vanish at those states/actions. Some other parts of the state/action space are still poorly understood. Because we add
46 lots of noise to the estimated rewards at those states, there is a significant chance they appear to be even better than they
47 really are in any given episode, in which case the algorithm will deftly navigate through the well understood part of
48 the state space trying to reach them. Once the algorithm starts actively trying to reach poorly understood states, it's
49 performing the kind of multi-period exploration that's essential for efficient RL.

50 To inject prior randomness, there is a natural counterpart (to sampling $Q$) for linear models, where you regularize the
51 parameter vector to a prior sample. This proecdure actually corresponds to Algorithm 1 if specialize it to the tabular
52 case (think $\theta = Q$). Things are more subtle with neural networks but [Osband et al., 2018] offers one approach. You're
53 right, that algorithm otherwise applies directly to settings with function approximation.