1 We would like to thank all the reviewers for their time and helpful feedback.

2 All reviewers mention that the current version of the paper is a little dense. We will work on improving clarity and
3 accessibility and try to convey more intuition in the main text. We would like to thank reviewers 1 and 2 for their helpful
4 suggestions on how to make the paper more accessible.

5 We also plan to add more simulation results in the final version. To reviewer 1, we will add simulation results for
6 the conjecture. To reviewer 2, we completely agree that it would be very interesting to see simulated examples of
7 information gain and regret bounds. For the linear-Gaussian bandit, the regret bound derived in the paper is comparable
8 to the best bounds in the literature up to logarithmic factors. It will be interesting to see the simulated total information
9 gain and plot the corresponding regret bounds. For MDPs, we suspect that the upper bounds on total information gain
10 are loose and the actual information gain might be much smaller.

11 Several reviewers ask about how the paper might lead to new algorithms. This is definitely an important direction for
12 future research. As a starting point, it may be interesting to experiment with UCB algorithms that use information
13 gain in their upper confidence functions. They come with regret guarantees for problems discussed in the paper, and it
14 would be interesting to see how they perform in practice. UCB algorithms that are studied previously usually construct
15 confidence sets around the model parameter based on past observations, and they do not consider how much information
16 would be revealed about the model when we see a new observation. By explicitly taking into account information
17 gain, these new UCB algorithms tend to give a higher bonus to actions that reveal more information about the system
18 compared to existing UCB algorithms. It would be interesting to see whether these new algorithms would work better
19 for some classes of problems.

20 Finally, we will work on clarifying the notations and proofs.

21 Below are additional responses to individual reviewers.

22 **Reviewer 1**

23 – *Line 26: please add references for "Most analyses.."*
24 Thanks for pointing it out. We will add the references.

25 – *Line 69: should $a$ be $A_\ell$ here?*
26 We use $A_\ell$ to denote a possibly random action selected by the algorithm, while $a$ is used to denote some fixed action.
27 We will clarify this in the final version.

28 – *It seems the $\theta$ variable is overloaded at different places. For example, Line 111 $\theta$ refers to the true model but*
29 *elsewhere it refers to the random variable.*
30 Under our Bayesian framework, the true model is a random variable. We will emphasize this in the final version.

31 **Reviewer 2**

32 – *The paper offers novel regret bounds for Thompson sampling and UCB algorithms but the key idea is mostly a*
33 *straightforward extension of Russo and Van Roy (JMLR'16)...*
34 It is true that our main idea bears similarity to Russo and Van Roy (JMLR'16), but a major difficulty of applying
35 their approach to complex environments is that it is unclear how to analyze the mutual information between optimal
36 actions and observations, as optimal actions themselves are fairly complicated objects in complex environments. Our
37 approach considers a more natural information gain between the model and observations, which sidesteps their issue
38 and allows us to obtain information-theoretic bounds for more complex problems like MDPs and factored MDPs that
39 are not achieved in their paper.

40 – *There are some problems with the notation in the paper...*
41 We will clarify our notation on probability measures. We use $T$ to denote the total number of time periods for the
42 linear bandit problem as it is more consistent with other bandit literatures, but I see how it can cause confusions and
43 we will make it clear in the final version.

44 – *...presenting less material (e.g., leave the factored MDPs for the appendix)...*
45 We include factored MDPs in the main text as an example to demonstrate that our approach is able to handle complex
46 structured environments that could not be handled using tools from previous work such as Russo and Van Roy
47 (JMLR'16).

48 **Reviewer 3**

49 – *...more exhaustive descriptions of the proof will help the readers and the reviewer to understand the paper...*
50 We will polish the proofs for the final version.