1 We thank reviewers R1, R2, R3, and R4 for their constructive and helpful feedback.

2 **Scope and significance (R4).** Our work derives the normative solution to the problem of how to dynamically allocate
3 noisy and limited memory resources for reinforcement learning. This work could have implications for machine learning
4 in the long run but our intended audience is currently neuroscientists, since memory access is inherently noisy in neural
5 circuits. We aim to use this work to generate testable predictions as to what should be observed in neural circuits
6 when learning complex tasks that require memory access. One possibility is that more neurons are devoted to specific
7 state-action pairs (in parietal cortex or in the basal ganglia, where q-values are putatively encoded) or memory of the
8 q-values might be sampled for a longer duration when higher precision is warranted, thus modulating the speed-accuracy
9 trade-off characteristic of decision making with sequential sampling. We aim to explore these ideas in future work.

10 **Disambiguating DRA from approaches in RL (R4).** R4 is correct in pointing out that other groups have proposed
11 alternative approaches to deal with memory limitations in RL, such as using regularization (SAC [A1]), or using neural
12 networks for representing policy and value functions, and even compressing state representations with graph Laplacian
13 [A2]. Our work is meant to complement these previous studies. SAC, for instance, directly penalizes the policy entropy
14 while maximizing reward to encourage exploration. In DRA, we penalize precise representations of (q-)values instead.
15 The use of Laplacian in RL [A2], on the other hand, hints at yet another problem involving efficient use of memory –
16 compact representation of states (*e.g.*, *chunking*) – which is something we very much look forward to addressing in
17 future work. We modified the *Related work* section to discuss and point out how these proposals complement our work.

18 **Justifying our assumptions (R4).** In our work, we assume that agents have limited capacity to store (q-)values in
19 their memory, but that they can observe and store states and actions perfectly. This choice is deliberate. It allowed us to
20 develop a normative solution to the problem of allocating resources to items in memory. As far as we know, we are the
21 first ones to provide such a solution. We agree with R4 that it would be important in the future to combine our approach
22 with alternatives that focus instead on regularizing the policy or compressing states, including heuristic approaches such
23 as truncation of planning trees, a strategy suggested Huys et al. [A3]. However, we hope that R4 will agree with us and
24 other reviewers that the existence of other approaches does not take anything away from our contribution: We propose a
25 normative solution to an important problem that "*is of high interest to the community (R2)*", and which will eventually
26 make concrete experimental predictions. "*This is not a paper searching for state of the art results, and it should not be
27 treated as such; rather, it is an exposition of a particular idea, and it did well to explore it (R1)*".

28 **Convergence and baseline (R1).** Following a suggestion from R1, we found with new analyses that the asymptotic
29 values of memory precision, $\sigma_*$, are largely independent of the choices of initial value $\sigma_0$. We also found that
30 convergence speed does not depend on the difference between the initial and optimal asymptotic values. Also, we can
31 confirm that letting $\lambda \to 0$ reduces DRA to SARSA. As suggested by R1, we pick $\lambda = 0$ as the baseline to compare
32 convergence speed and found that DRA is $1.5\times$ and $1.4\times$ slower in the grid-world and mountain car tasks respectively.

33 **Comparison with equal resource allocation (R1, R3).** We tested DRA against a model that allocates equal resources
34 to all memories for the grid-world and mountain car tasks (ref. Fig. 3d), and report a $2\times$ and $1.3\times$ improvement in the
35 objective (Eq.1) respectively with flexible resource allocation (DRA). We added these results to the revised manuscript.

36 **Sampling trajectories (R3).** Conceptually, the sampling of memories is similar to Dyna, with the difference being
37 that instead of randomly sampling individual memories, DRA samples entire trajectories *on-policy* (though they could
38 also be drawn from stored episodic memories). In DRA, replays are entirely forward. In our current implementations,
39 we sample trajectories at the end of each trial from the starting state for that trial to termination. In principle, we could
40 also do so multiple times during the trial, *e.g.* at each state, and not necessarily until termination.

41 **Generalizing DRA (R2, R4).** We aim to address non-independent memories and continuous state spaces in future by
42 considering a Gaussian process prior over the q-values and using GPSARSA [A4] instead of SARSA to update the mean
43 q-values, and extend to non-tabular settings by incorporating compact state representations, *e.g.*, [A2].

44 **Response to remaining comments.** **R1:** We have now clarified the meaning of "more" or "less" resources in the
45 text, but we insist that our arguments apply to all systems that are restricted to sample from value distributions but
46 cannot access its mean and precision directly. We depicted the normalized entropy in Fig. 1 simply to ease visualization,
47 but performed appropriate checks as mentioned earlier. **R2:** The mountain car task was included to demonstrate general
48 applicability of DRA to arbitrary problems. We now dedicate a section in the main text to *Related work* and discuss [A1,
49 A2, A5]. **R3:** Figs. 1b & 1c come from different simulations. We have rewritten the caption for Fig. 3 (also suggested
50 by **R4**) clarifying the y-axis confusion in 3f. We also mention that $a_{\text{dec}}$ & $t_{\text{non-dec}}$ come from simulations. **R4:** In all
51 tasks, we systematically varied the values of $\lambda$ & $\sigma_{\text{base}}$ and report that the qualitative results hold for arbitrary values of
52 these parameters with $\sigma_{\text{base}}$ having a slightly stronger effect than $\lambda$.

53 **Code release** We intend to release the code on GitHub as soon as the submissions are no longer anonymous.

54 **References.** [A1] Haarnoja et al. *ArXiV* 2018. [A2] Wu et al., *ArXiV* 2018. [A3] Huys et al. *PNAS* 2015. [A4] Engel
55 et al. *ICML* 2005. [A5] Mattar & Daw. *Nat. Neuro.* 2019.