

# Empirical convergence analysis

## Replay threshold

The plots below show how the value of the root state changes with the EVB threshold  $\xi$  at multiple horizons, inverse temperatures, and for various prior beliefs. Note that the root  $Q$ -values (i.e, the current MF  $Q$ -values) were set to 0 everywhere, and the  $Q$ -values at all other belief states were set to the expected immediate reward.

The prioritisation was determined by:

$$\text{EVB}(b, a) = p(b_\rho \rightarrow b, \pi) \times \mathbb{E}_{b' \sim p(b'|b, a, \pi)} [v(b') - v(b)]$$

where  $b_\rho$  is the root belief and  $\pi$  is a softmax policy.

A couple of interesting observations. First, it looks like higher inverse temperature ( $\beta$ ) needs fewer replays to approach the Bayes-optimal value at the root (it doesn't completely get there though because of the softmax). Second, there is an interesting horizon -  $\xi$  trade-off – higher horizons need fewer replays to get a better root value estimate than lower horizons; this means that one can (potentially) increase the horizon at the expense of setting a stricter (higher)  $\xi$  threshold. As we have seen before, there is also a somewhat "emergent" depth-breadth trade-off which depends on uncertainty – i.e., for the more uncertain arm the algorithm tends to backup values "in breadth", and for the more certain arm it goes more "in depth".

$\alpha_0 = 5, \beta_0 = 1, \alpha_1 = 1, \beta_1 = 1$

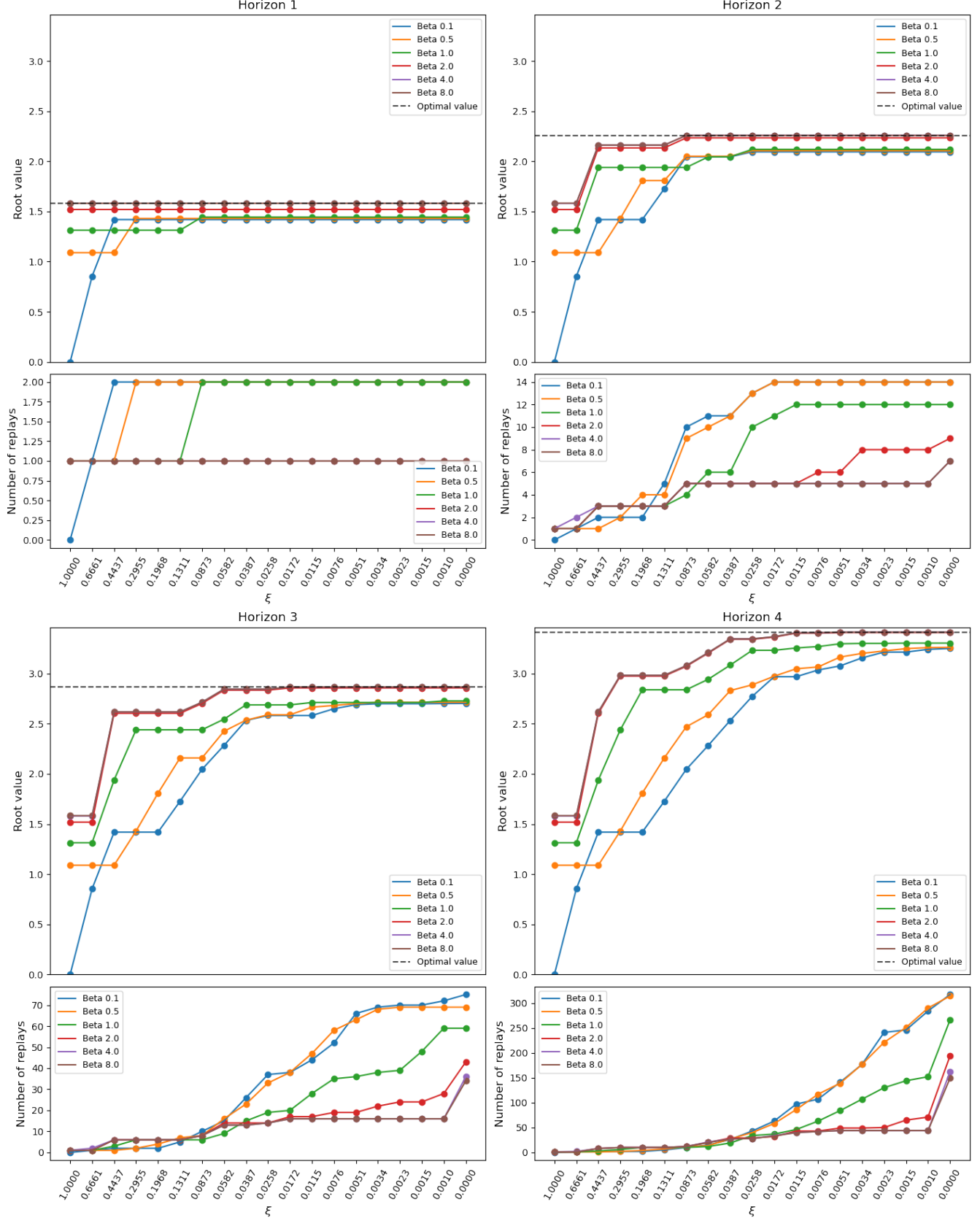
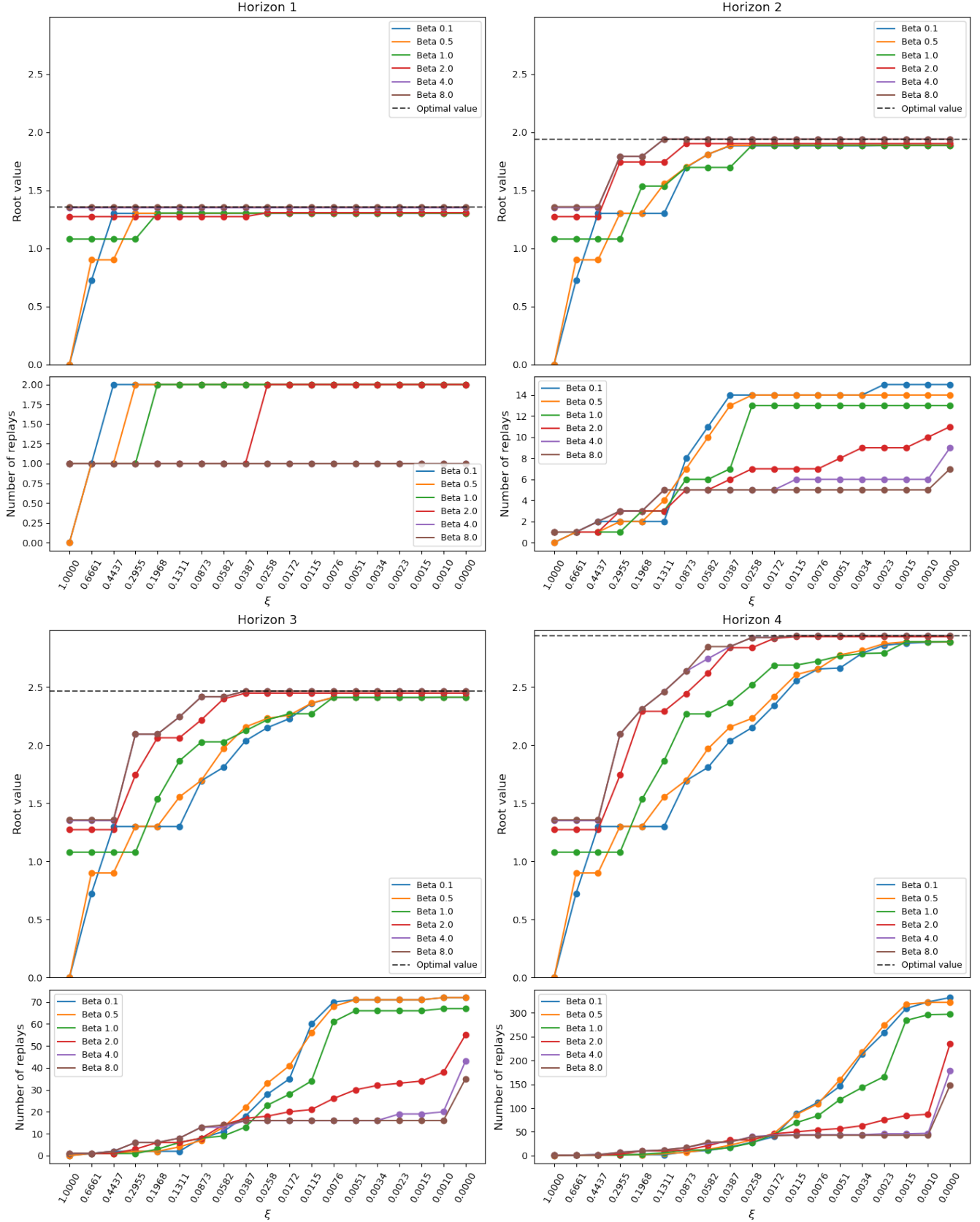
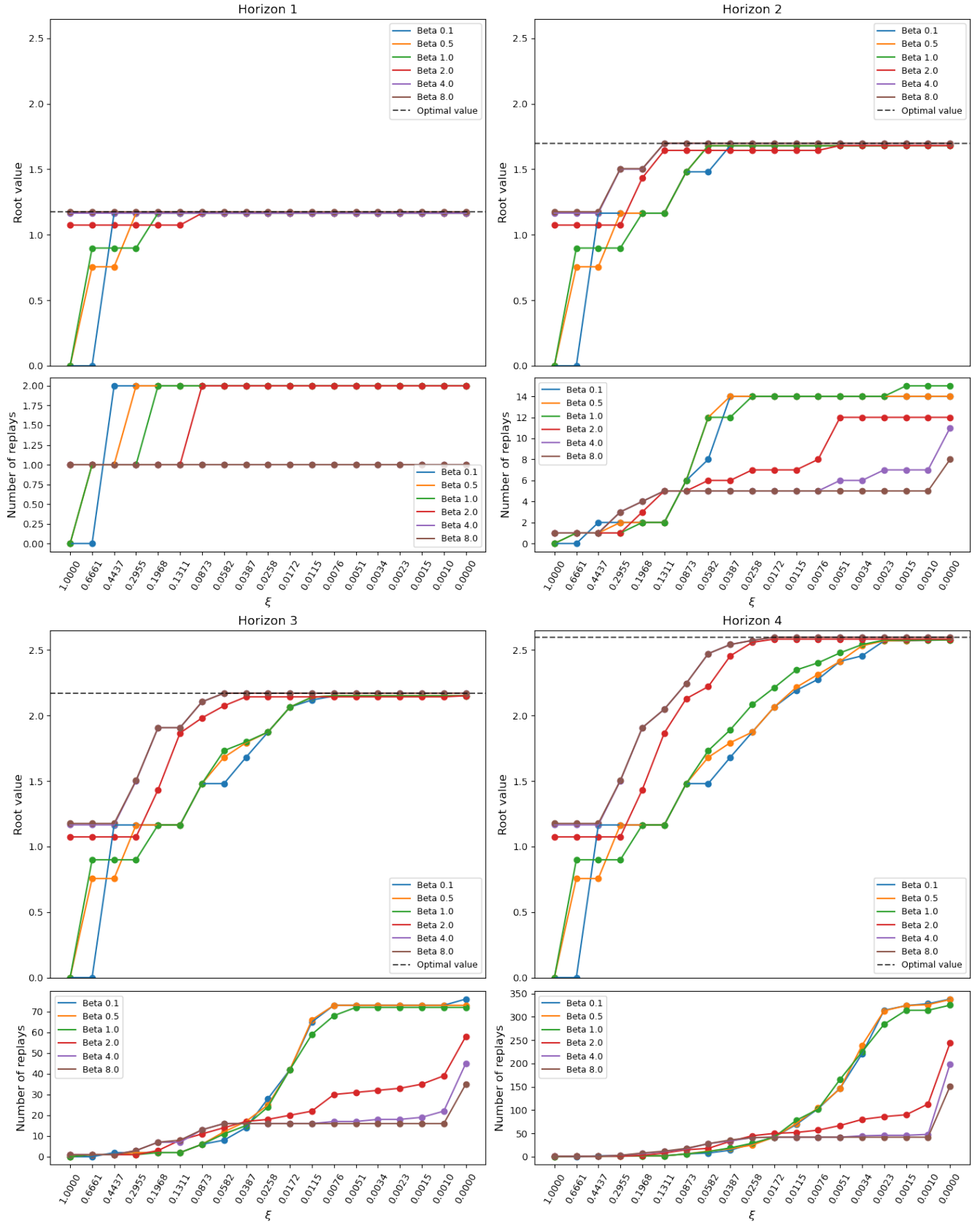


Figure 1: Note that the same inverse temperature parameter  $\beta$  was used within the tree and at the root (which is the plotted root value).

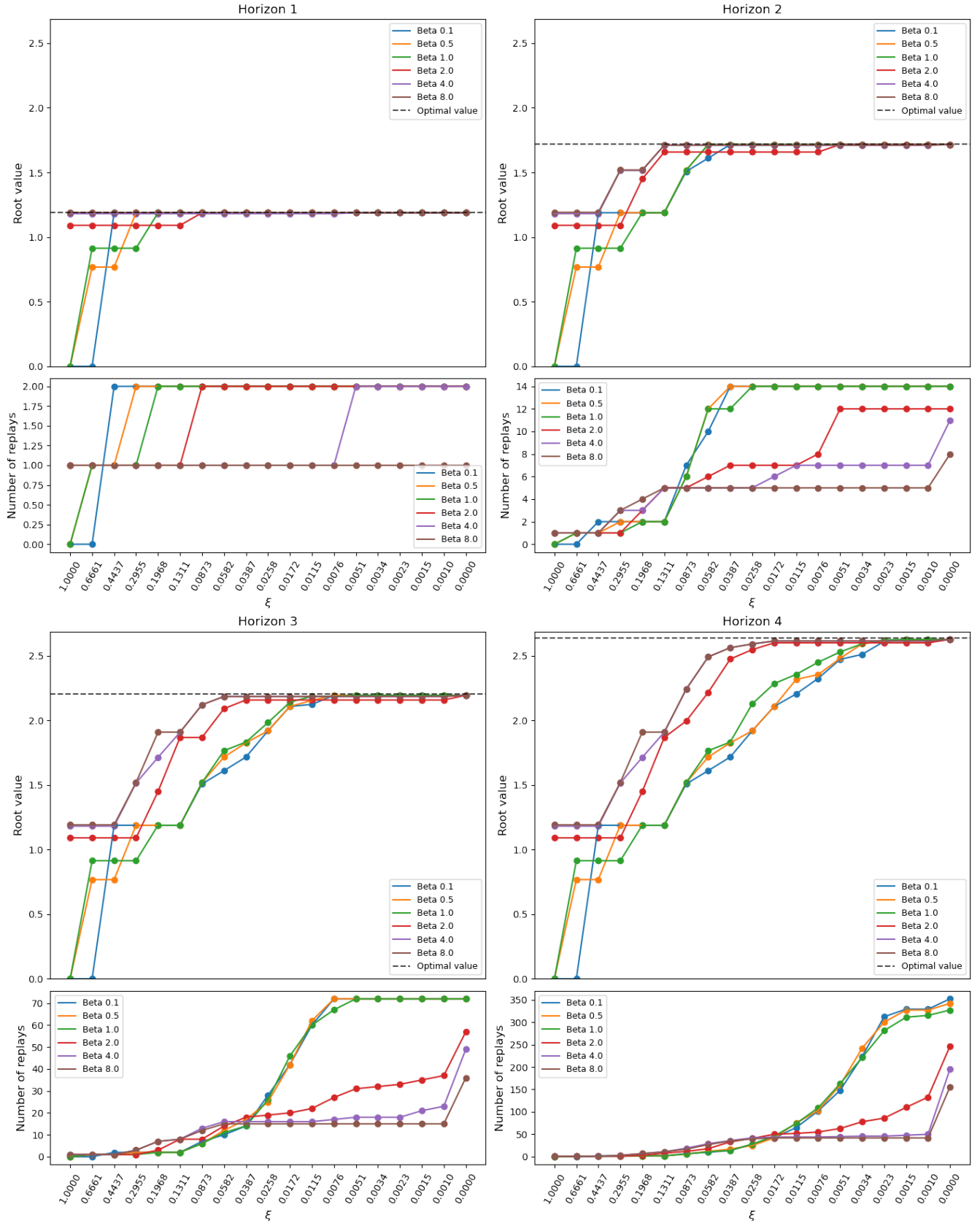
$\alpha_0 = 5, \beta_0 = 2, \alpha_1 = 3, \beta_1 = 2$



$\alpha_0 = 9, \beta_0 = 6, \alpha_1 = 3, \beta_1 = 2$



$\alpha_0 = 10, \beta_0 = 6, \alpha_1 = 3, \beta_1 = 2$



## Maze simulations

The prioritisation of experiences at **distal** information states was determined by:

$$\text{EVB}(\langle s, b \rangle, a) = p(\langle s_\rho, b_\rho \rangle \rightarrow \langle s, b \rangle, \pi) \times \mathbb{E}_{\langle s', b' \rangle \sim p(\langle s', b' \rangle | \langle s, b \rangle, a, \pi)} [v(\langle s', b' \rangle) - v(\langle s, b \rangle)]$$

where  $s_\rho$  is the root physical state and  $b_\rho$  is the root belief.

The implementation is very similar to that of the bandit case – during planning (replay), the agent generates a planning tree from every state and for every action. The updates can then happen either in distal parts of the trees (only 1 update is allowed at each iteration), or at the root. The value of each tree node is initialised to the agent’s current MF  $Q$ -values.

For the **root** action updates, the EVB is computed in the standard M&D fashion:

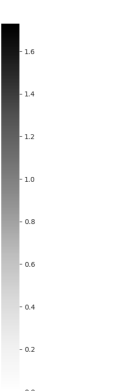
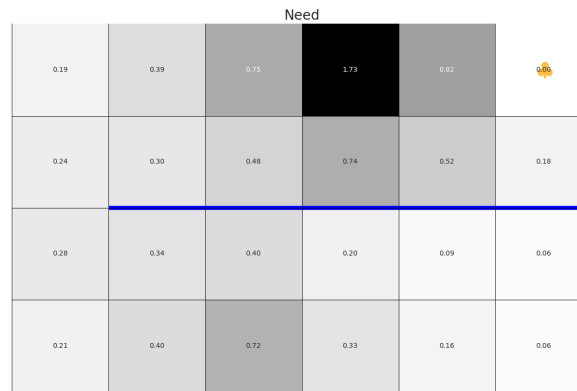
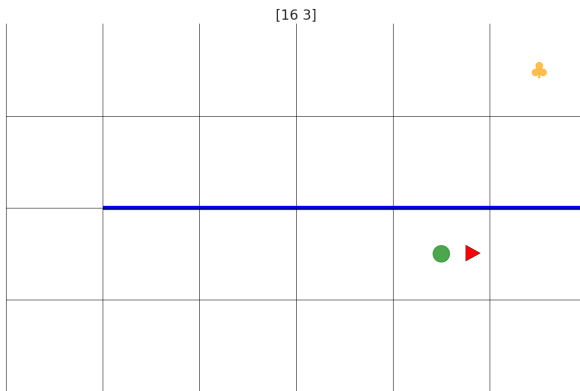
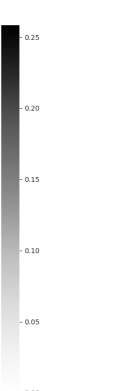
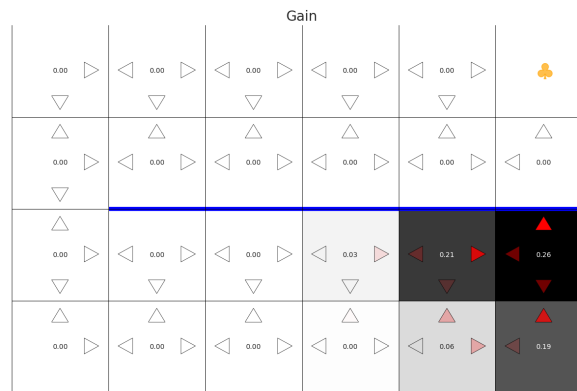
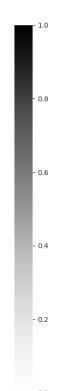
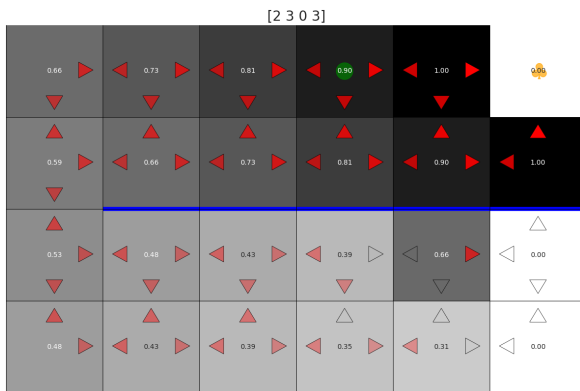
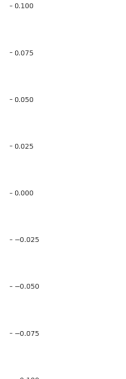
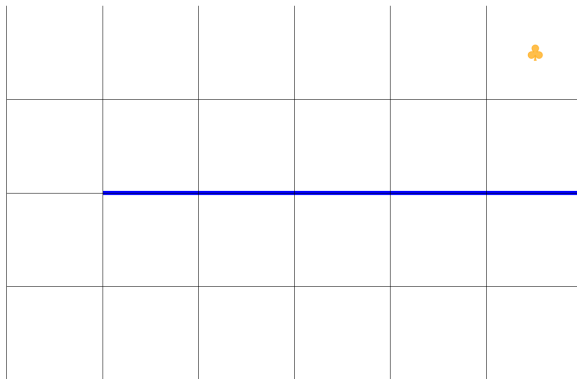
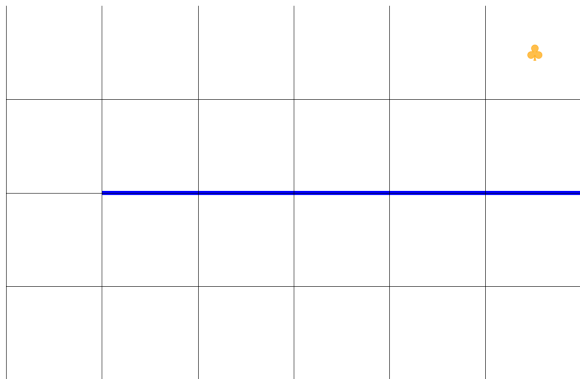
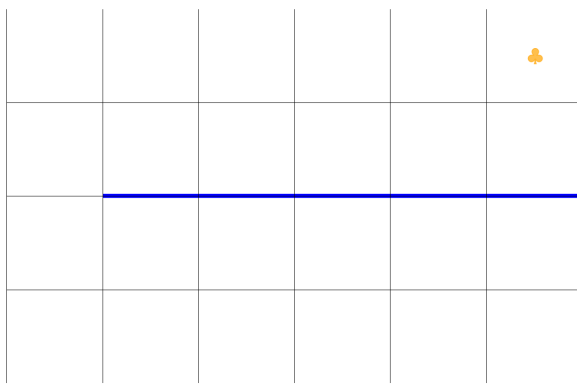
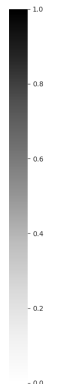
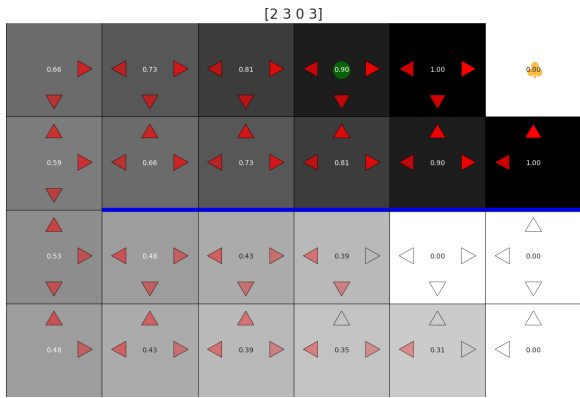
$$\text{EVB}(\langle s_\rho, b_\rho \rangle, a) = \text{Need}(\langle x, b_\rho \rangle \rightarrow \langle s_\rho, b_\rho \rangle, \pi) \times \text{Gain}(\langle s_\rho, b_\rho \rangle, a)$$

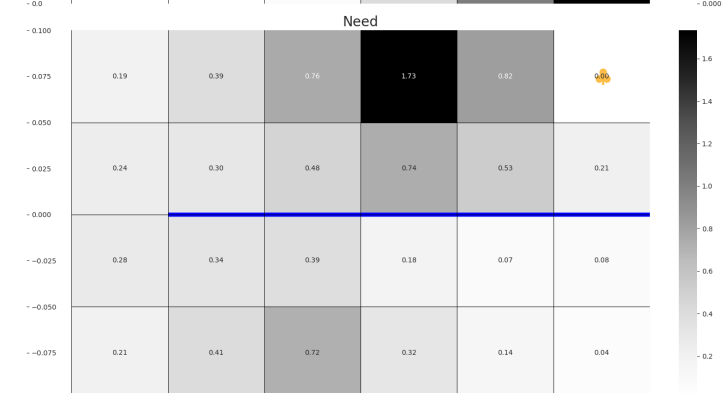
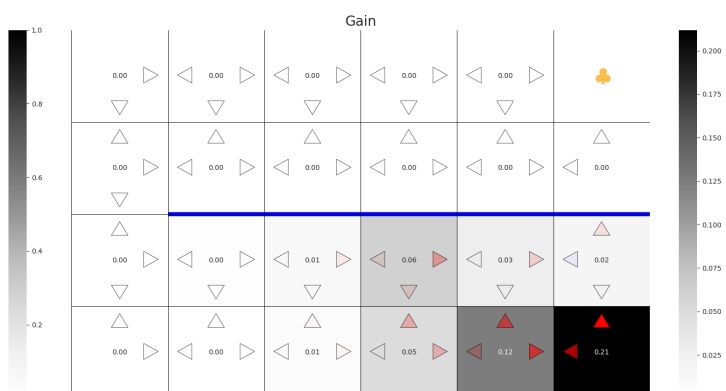
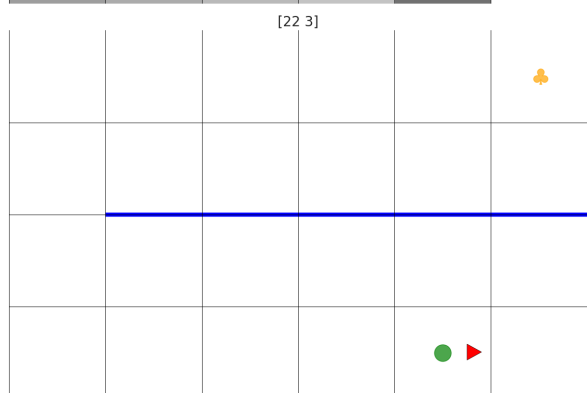
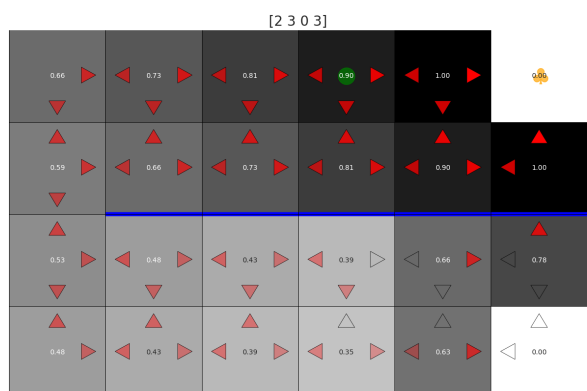
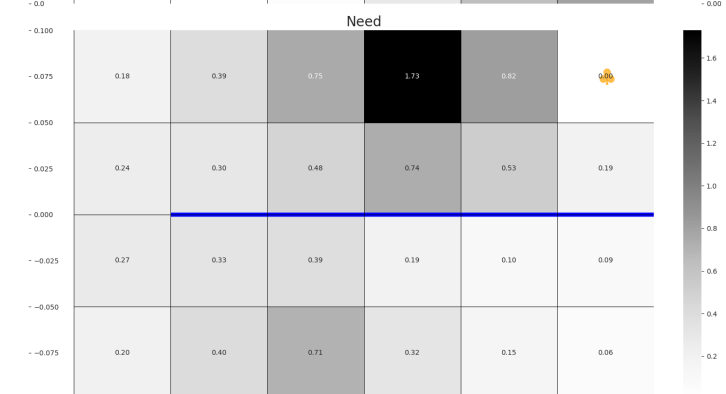
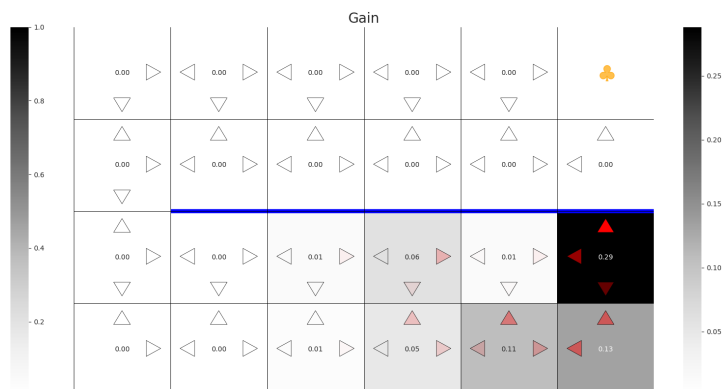
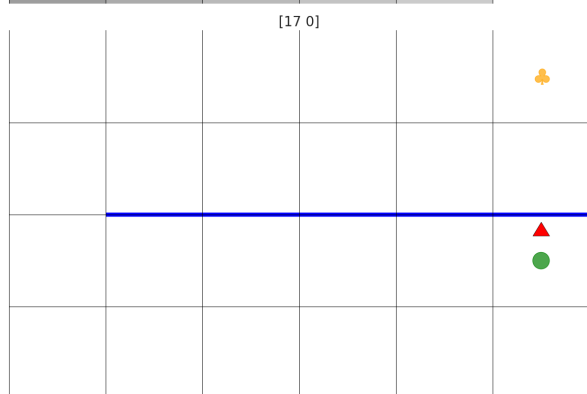
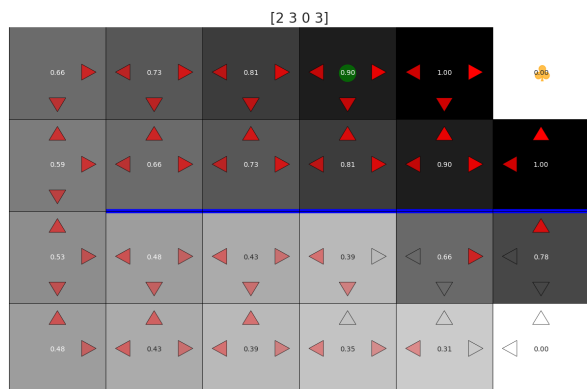
This is a bit tricky to write down in a less cryptic way. Essentially, the assumption is that state  $s_\rho$  can be reached from the agent’s current location  $x$  without belief dynamics – that is, Need in this case is exactly the same as in M&D.

The crucial difference is the Gain term. Our Bayes-adaptive prioritised sweeping procedure effectively ”chooses” which computation to perform at the root. Since it approximates a Bayes-optimal value (yet with a heuristic at the distal belief in the form of the current MF  $Q$ -values), the future information gain is subsumed by M&D’s Gain at the root.

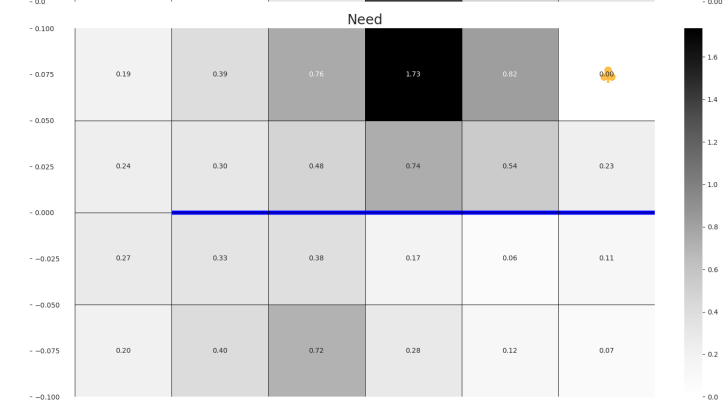
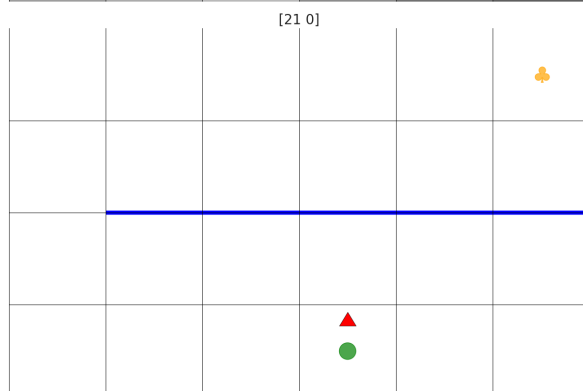
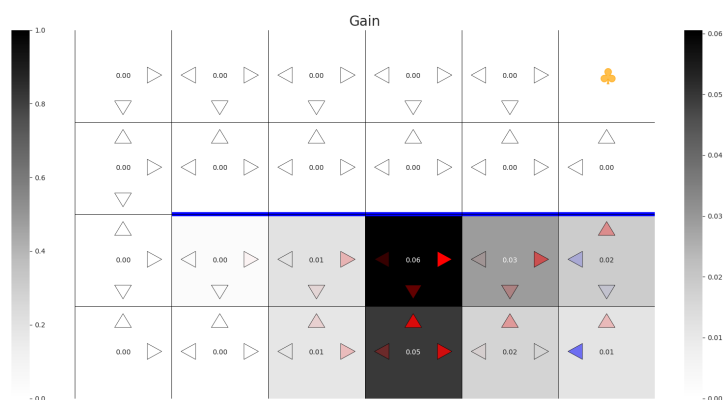
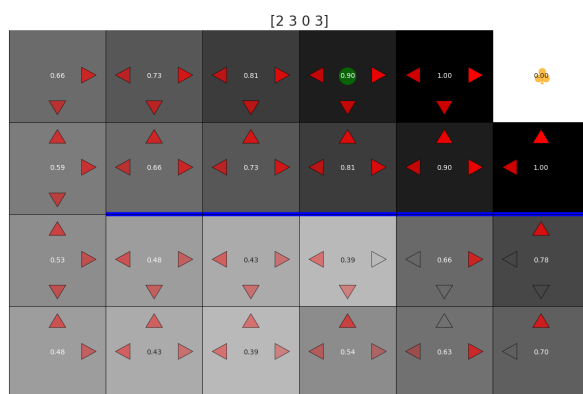
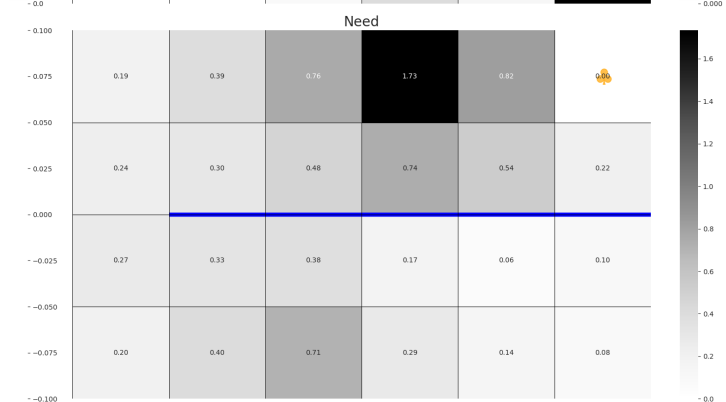
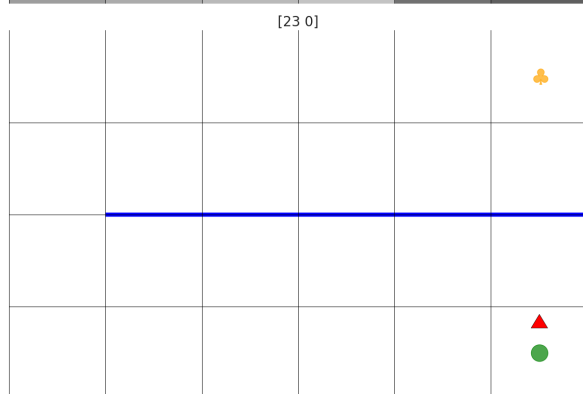
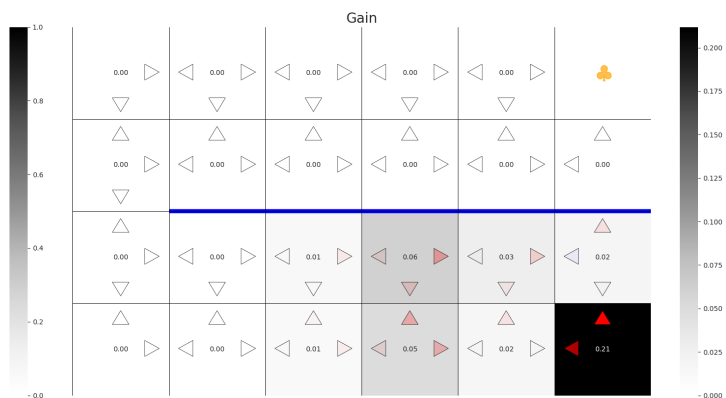
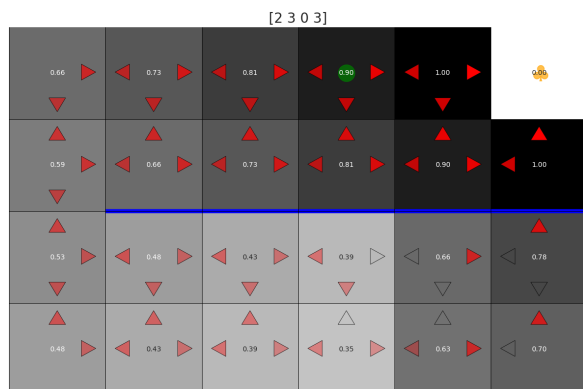
Below I show an example exploratory replay from a toy simulation. First, the agent was simulated in the environment for 3000 moves without being allowed to replay (so only MF learning) – and luckily, with this seed, it never got to experience some of the bottom rightmost states (the ones where MF  $Q$ -values are all 0). The agent is completely certain about all transitions; however, at move 3000 I set its Beta posterior to (1, 1) for the action ’up’ below the rightmost barrier – hence making it expect that the probability of that barrier being removed is 0.5. I then let it replay with horizon 2. The start state is the square which is the 3rd from the left at the very bottom.

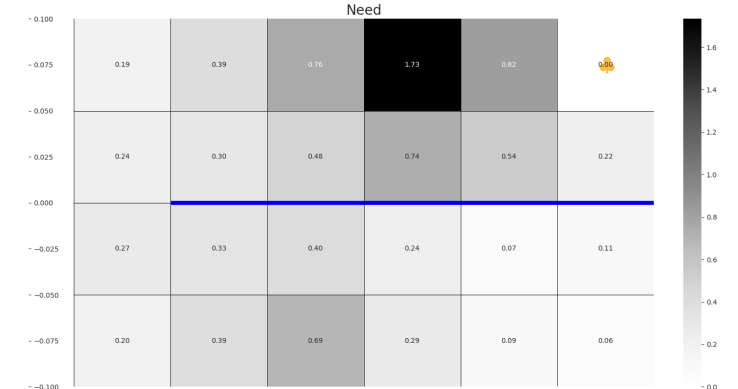
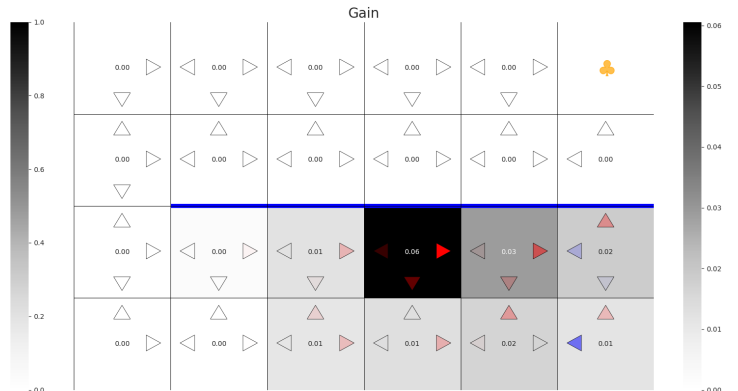
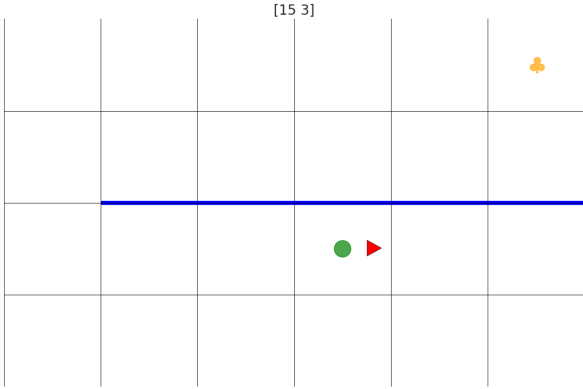
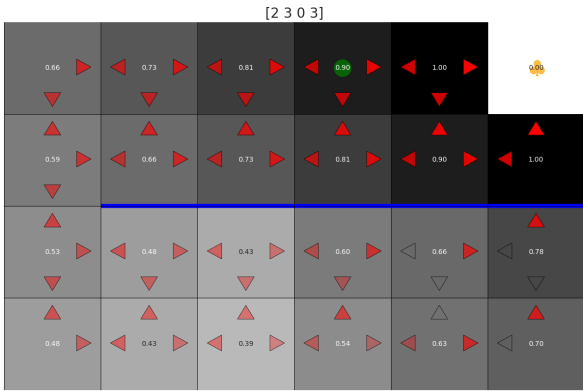
The top left plots show MF  $Q$ -values and the agent’s current location (green dot); the bottom left plots show the replay chosen by the agent; the top right plots show Gain; and the bottom right plots show Need. Note that only root updates are shown.











I think this is very interesting and there are quite a few important considerations here. First, the horizon really matters – since it helps overcome the very low Need at (currenty) unexpected states. In the example shown above, for instance, if the horizon was set to 1, there would be no exploratory replays at all. This is because Need at the state adjacent to the shortcut is very small. However, for a larger horizon (2), the agent could plan crossing the barrier from a state which had a slightly higher Need – just enough so that the expected value of that replay crossed the threshold, even though Gain was slightly lower due to  $\gamma$  (which was set to 0.9 in this simulation). That replay, in turn, increased the Need for the subsequent state (the one adjacent to the barrier), and hence that replay followed next. This is the reason why the model predicts such exploratory replays to proceed in a forward manner.

One thing which feels a bit uncomfortable is that Gain turns out to be much more powerful than Need (because Gain can dramatically change Need). This, of course, depends on the horizon and learning rates (replay learning rate was set to 1 in this simulation), as well as the discount factor.

For clarity, I also wrote down a very high-level pseudocode of the algorithm on the next page.

## Replay pseudocode

### Initialise

$x \leftarrow$  agent's current physical state  
 $b_\rho \leftarrow$  agent's current belief  
 $H \leftarrow$  horizon  
 $\alpha_r \leftarrow$  replay learning rate  
 $\gamma \leftarrow$  discount factor  
 $\xi \leftarrow$  EVB replay threshold

### for each $s \in \mathcal{S}$ do

build belief tree  $t_s$  with root information state  $z_\rho = \langle s_\rho = s, b_\rho \rangle$  up to horizon  $H$   
 initialise  $Q_{t_s(h)}(z = \langle k, b \rangle, a) \leftarrow Q^{MF}(\langle k, b_\rho \rangle, a) \quad \forall h \in H, z \in t_s(h)$

### done

### while $\max_{z, a, t_s, h} \text{EVB}_{t_s(h)}(z, a) > \xi$ do

#### for each $h$ in $(H - 1)$ do

#### for each information state $z = \langle k, b \rangle \in t_s(h)$ do

#### for each action $a \in \mathcal{A}$ do

#### if $(h > 1)$ do

$\text{EVB}_{t_s(h)}(z, a) \leftarrow p(z_\rho \rightarrow z, \pi) \times \mathbb{E}_{z' \sim p(z'|z, a, \pi)} [v(z') - v(z)]$

#### else if $(h == 1)$ do

$\text{EVB}_{t_s(h)}(z_\rho, a) \leftarrow \text{Need}(\langle x, b_\rho \rangle \rightarrow z_\rho, \pi) \times \text{Gain}(z_\rho, a)$

#### done

#### done

#### done

$z, a, t_s, h \leftarrow \text{argmax}_{z, a, t_s, h} \text{EVB}_{t_s(h)}(z, a)$   
 if  $(h == 1)$  do

$Q^{MF}(z, a) \leftarrow Q_{t_s(h)}(z, a) + \alpha_r \mathbb{E}_{z' \sim p(z'|z, a, \pi)} [R(z') + \gamma \max_{a' \in \mathcal{A}} Q_{t_s(h+1)}(z', a') - Q(z, a)]$

#### for each $s \in \mathcal{S}$ do

build belief tree  $t_s$  with root information state  $z_\rho = \langle s_\rho = s, b_\rho \rangle$  up to horizon  $H$

initialise  $Q_{t_s(h)}(z = \langle k, b \rangle, a) \leftarrow Q^{MF}(\langle k, b_\rho \rangle, a) \quad \forall h \in H, z \in t_s(h)$

#### done

#### else if $(h! = 1)$ do

$Q_{t_s(h)}(z, a) \leftarrow Q_{t_s(h)}(z, a) + \alpha_r \mathbb{E}_{z' \sim p(z'|z, a, \pi)} [R(z') + \gamma \max_{a' \in \mathcal{A}} Q_{t_s(h+1)}(z', a') - Q(z, a)]$

#### done

### repeat