

Supplemental Materials

S1 Calculating $\nabla_{\theta} I_{\text{state}}(t)$

We want to evaluate:

$$\nabla_{\theta} I_{\text{state}} = \sum_g \rho_G(g) \sum_s (\nabla_{\theta} p(s | g)) \log \frac{p(s | g)}{p(s)} \quad (\text{S1})$$

$$+ \sum_g \rho_G(g) \sum_s p(s | g) \frac{\nabla_{\theta} p(s | g)}{p(s | g)} \quad (\text{S2})$$

$$- \sum_g \rho_G(g) \sum_s p(s | g) \frac{\nabla_{\theta} p(s)}{p(s)} \quad (\text{S3})$$

$$\equiv T_1 + T_2 - T_3, \quad (\text{S4})$$

where we denote the three terms by T_1 , T_2 , and T_3 . The effect of T_1 follows from the policy gradient theorem and amounts to adding the following to the reward return:

$$\sum_{t'=t}^T \log \frac{p_{\text{emp}}(s_{t'} | g)}{p_{\text{emp}}(s_{t'})}. \quad (\text{S5})$$

By the same argument, $T_2 = \sum_g p(g) \sum_s \nabla_{\theta} p(s | g)$ simply results in the addition of 1 to the info return at each time step.

Finally, we have the third term:

$$T_3 = \sum_g \rho_G(g) \sum_{s_t} \frac{p(s_t | g)}{p(s_t)} \nabla_{\theta} p(s_t) \quad (\text{S6})$$

$$= \sum_g \rho_G(g) \sum_{s_t} \frac{p(s_t | g)}{p(s_t)} \nabla_{\theta} \sum_{g'} \rho_G(g') \rho_S(s_0) \prod_{t'=0}^t \pi_{g'}(a_{t'} | s_{t'}) P(s_{t'+1} | s_{t'}, a_{t'}) \quad (\text{S7})$$

$$= \sum_g \rho_G(g) \sum_{s_t} \frac{p(s_t | g)}{p(s_t)} \sum_{g'} \rho_G(g') \rho_S(s_0) \prod_{t'=0}^t (\nabla_{\theta} \pi_{g'}(a_{t'} | s_{t'})) P(s_{t'+1} | s_{t'}, a_{t'}) \quad (\text{S8})$$

$$= \sum_{g, s_t} \rho_G(g) \rho_S(s_0) \prod_{t'=0}^t \pi_g(a_{t'} | s_{t'}) P(s_{t'+1} | s_{t'}, a_{t'}) \times \quad (\text{S9})$$

$$\sum_{g'} \rho_G(g') \frac{\pi_{g'}(a_{t'} | s_{t'})}{\pi_g(a_{t'} | s_{t'})} \frac{\nabla_{\theta} \pi_{g'}(a_{t'} | s_{t'})}{\pi_{g'}(a_{t'} | s_{t'})} \frac{p(s_t | g)}{p(s_t)} \quad (\text{S10})$$

$$\equiv \mathbb{E}_{\tau} \left[\sum_{g'} \rho_G(g') \prod_{t'=0}^t \frac{\pi_{g'}(a_{t'} | s_{t'})}{\pi_g(a_{t'} | s_{t'})} (\nabla_{\theta} \log \pi_{g'}(a_{t'} | s_{t'})) \frac{p(s_t | g)}{p(s_t)} \right] \quad (\text{S11})$$

where in the fourth line we multiply and divide by the policy under both g and g' in order to employ the log derivative trick and to express the equation as an expectation under the present goal. The end result is the update in equation 11.

S2 Experimental parameters and details

S2.1 Simple spatial navigation

In order to allow Bob to integrate information about the goal over time and remember it to guide future actions, we endow Bob with a recurrent neural network (RNN) to process Alice's state-action

pairs. We used a gated recurrent unit (GRU) Cho et al. [2014] to which Alice’s state-action pairs are fed as a one-hot vector. We chose to use a scalar core state for the GRU since it was simply tasked with tracking Bob’s belief about one of two goals, and could thus assign each goal to a sign of the GRU core state/output, which is what Bob chose to do in practice. The GRU output $z_t = \text{RNN}(s_t^{\text{alice}}, a_t^{\text{alice}})$ was then concatenated with a one-hot representation of Bob’s own state s_t^{bob} and fed into a fully-connected, feed-forward layer of 128 units with two readout heads: a policy head (a linear layer with $|\mathcal{A}|$ units followed by a softmax, yielding $a_t^{\text{bob}} \sim \pi^{\text{bob}}(s_t^{\text{bob}}, z_t)$) and a value head (a single linear readout node, yielding $v_t = V^{\text{bob}}(s_t^{\text{bob}}, z_t)$).

	Alice	Bob
training time, in steps	100k	200k
max episode length, in steps	100	100
entropy bonus (logarithmically annealed from/to)	.5, .005	.5, .01
learning rate (Adam)	2.5×10^{-2}	5×10^{-5}
weight on value function regression term	.5	.5
discount γ	.8	.8

Table 1: **Training parameters.**

S2.2 Key game

The only difference from the previous set of training parameters is that Alice now trains longer (250k instead of 100k steps).

	Alice	Bob
training time, in steps	250k	200k
max episode length, in steps	100	100
entropy bonus (logarithmically annealed from/to)	.5, .005	.5, .01
learning rate (Adam)	2.5×10^{-2}	5×10^{-5}
weight on value function regression term	.5	.5
discount γ	.8	.8

Table 2: **Training parameters.**