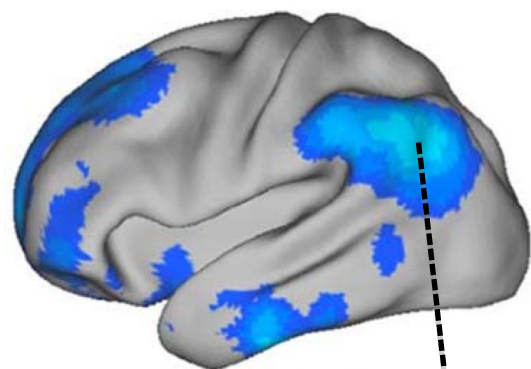


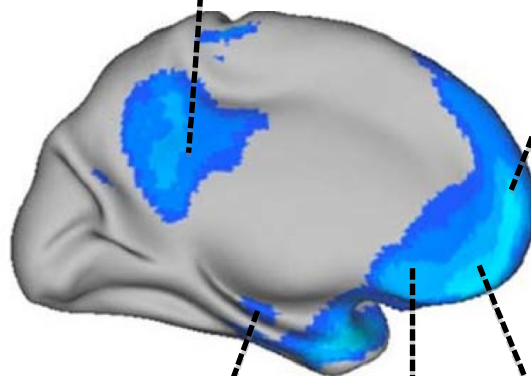
$$\nabla_{\theta_{k+1}^Q} \mathcal{L}(\theta_{k+1}^Q) = \mathbb{E}_{(s,a,r,s') \sim \mathcal{U}(\mathcal{D})} \left[\underbrace{(Q(s,a|\theta_{k+1}^Q) - y)}_{\text{regret}} \underbrace{\nabla_{\theta_{k+1}^Q} Q(s,a|\theta_{k+1}^Q)}_{\text{aversiveness}} \right]$$

parameter update

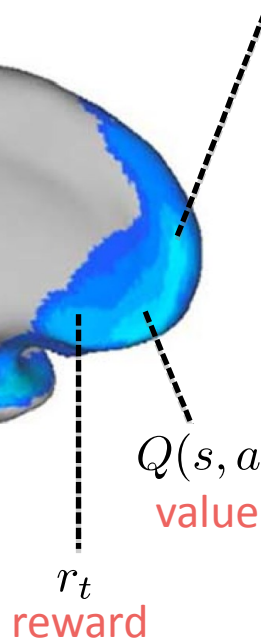
$$Q^\pi(s,a) = \mathbb{E}_{s' \sim p(s'|s,a)} [r(s,a) + \gamma \max_{a' \in \mathcal{A}} Q^\pi(s',a')] \\ \text{Bellman equation}$$



s, a
world knowledge

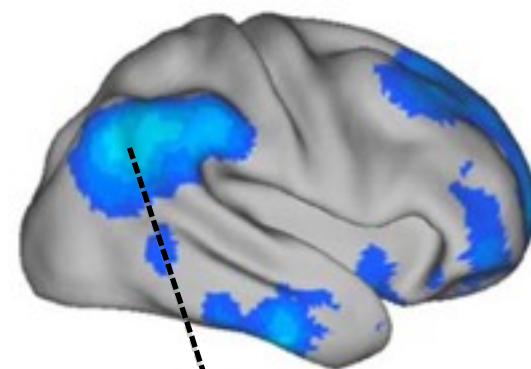


$(s,a,r,s') \sim U(\mathcal{D})$
experience sampling



$Q(s,a|\theta_k^Q)$
value function

r_t
reward



$Q(s,a|\theta_{k+1}^Q) - y$
Prediction error signaling

