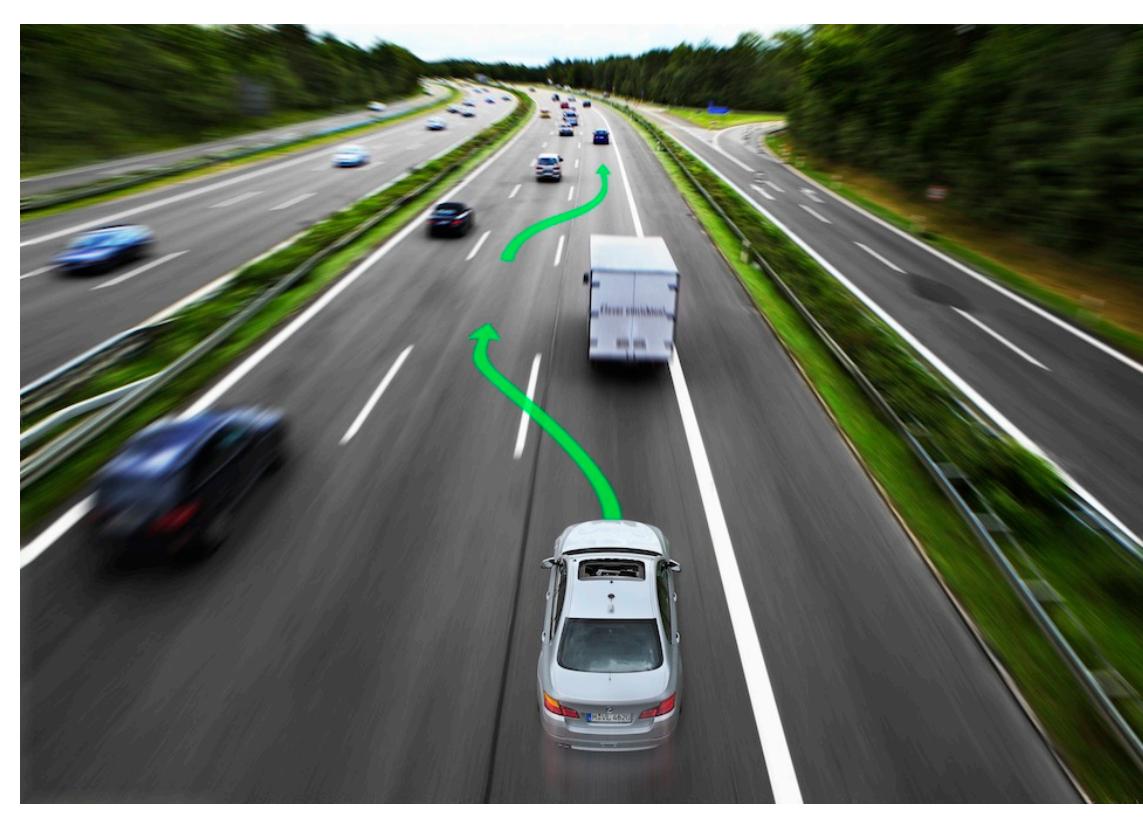
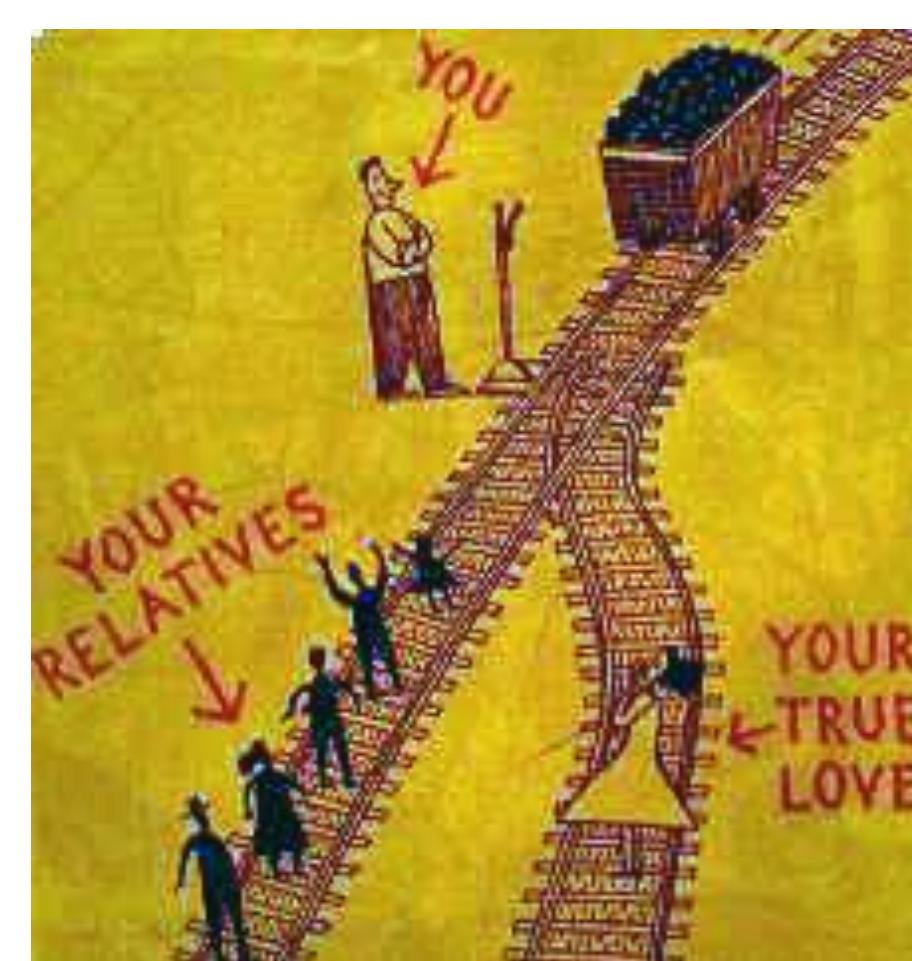


Kinesthetic teaching difficult for high-DOF robots



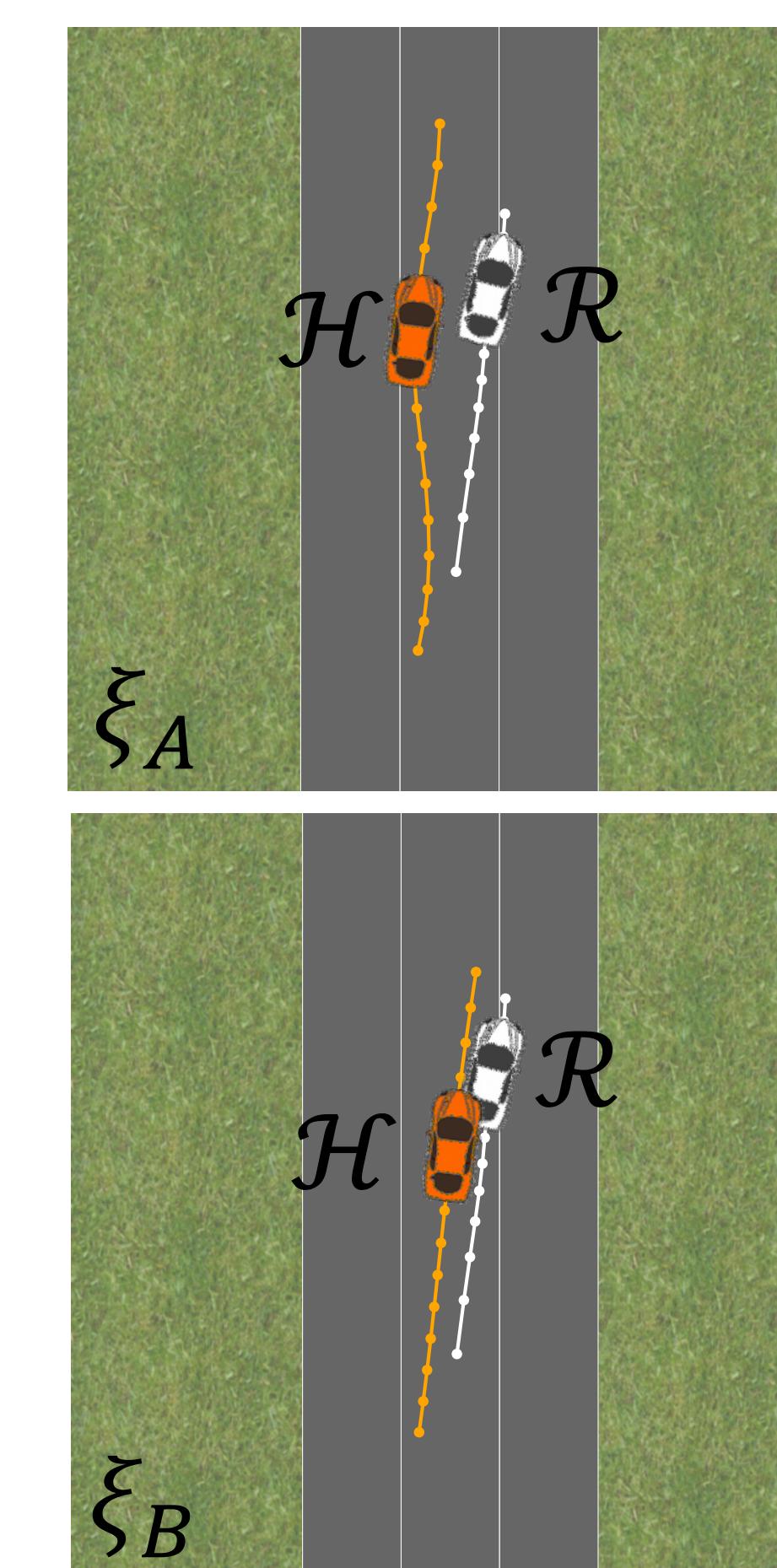
Want cars to drive differently than we drive



Hard to collect data in extreme scenarios

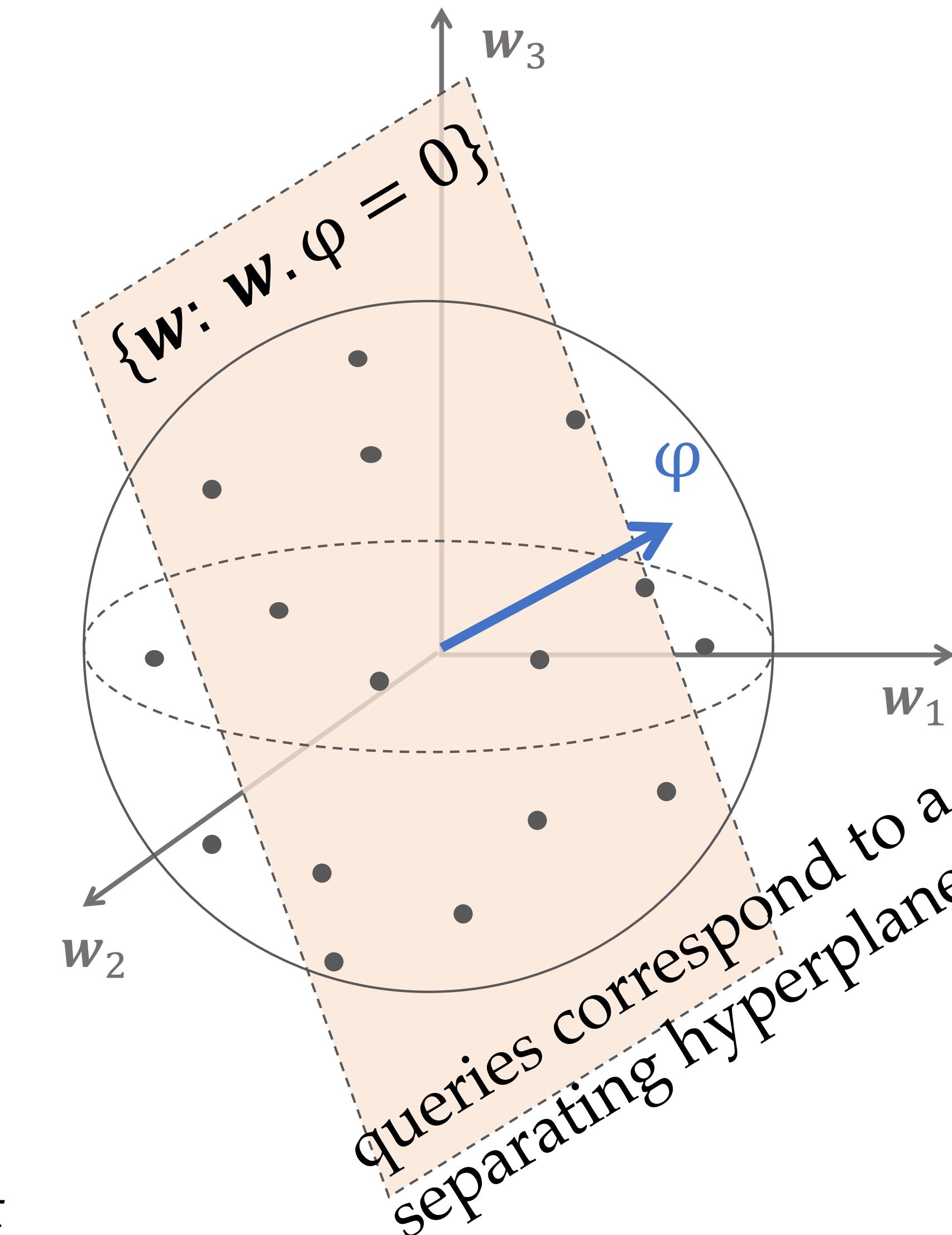
People can't always provide demonstrations of what they actually want a robot to do. We instead leverage **comparisons** as useful observations about the desired robot reward function.

## I. Query the human



$\xi_A$  or  $\xi_B$  ?  $\rightarrow I_t$

## II. Update belief over $\mathbf{w}$



## III. Actively synthesize queries

*expected volume removed*

$$\max_{\varphi} \min\{\mathbb{E}[1 - f_{\varphi}(\mathbf{w})], \mathbb{E}[1 - f_{-\varphi}(\mathbf{w})]\}$$

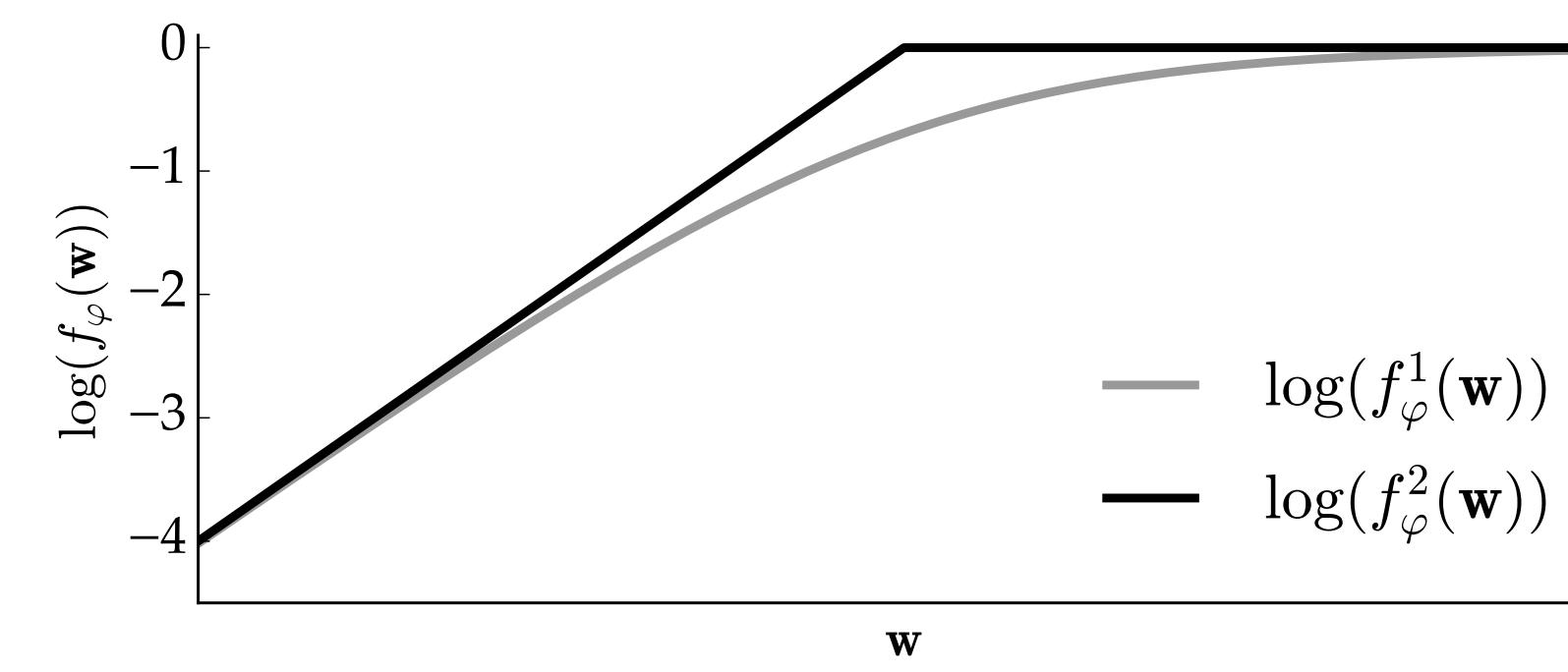
subject to  $\varphi \in \mathbb{F}$

$$\mathbb{F} = \{\varphi: \varphi = \Phi(\xi_A) - \Phi(\xi_B), \xi_A, \xi_B \in \Xi\}$$

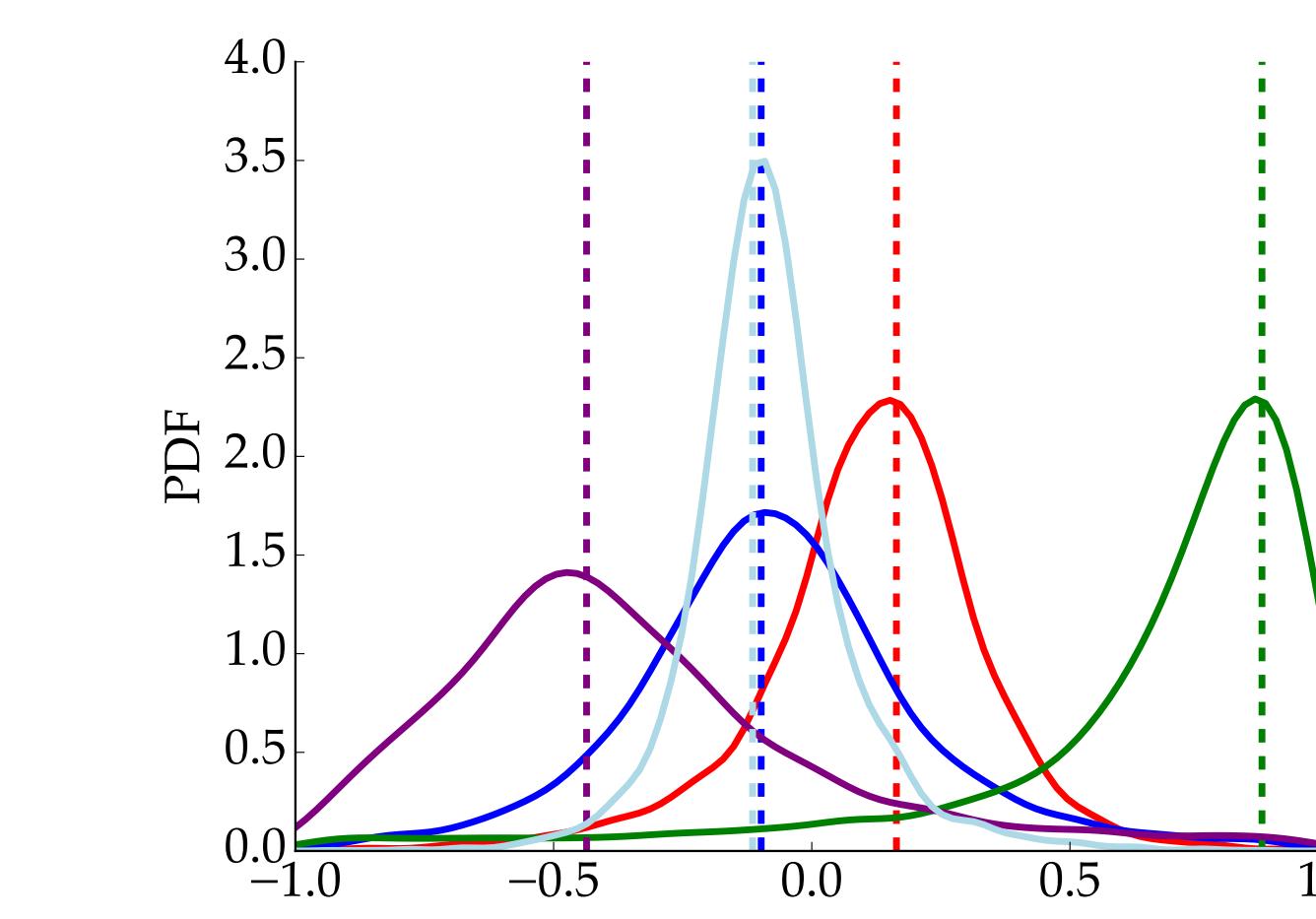
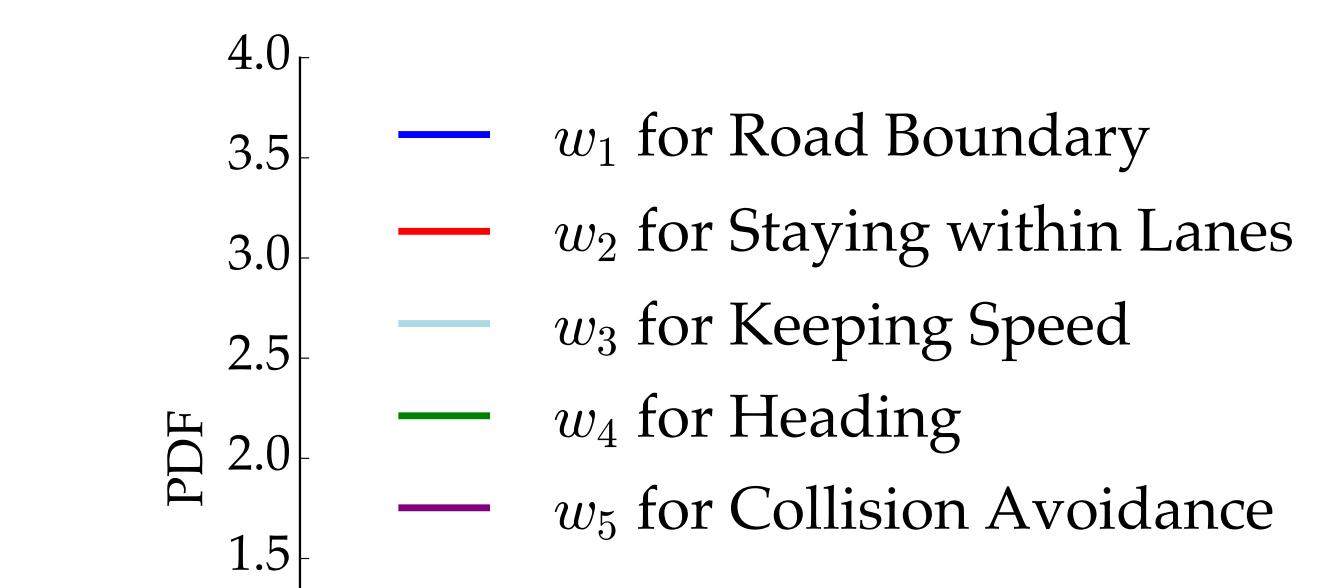
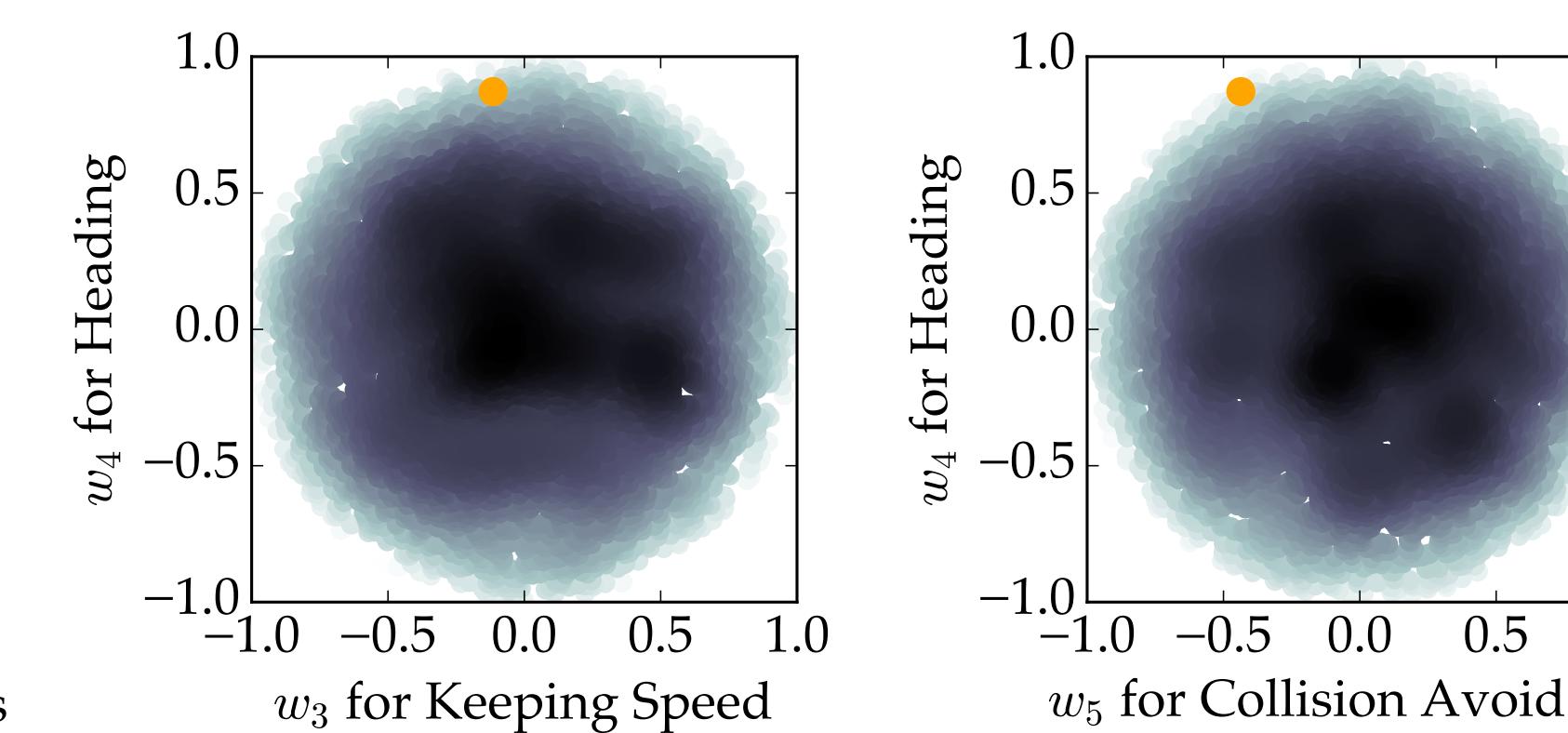
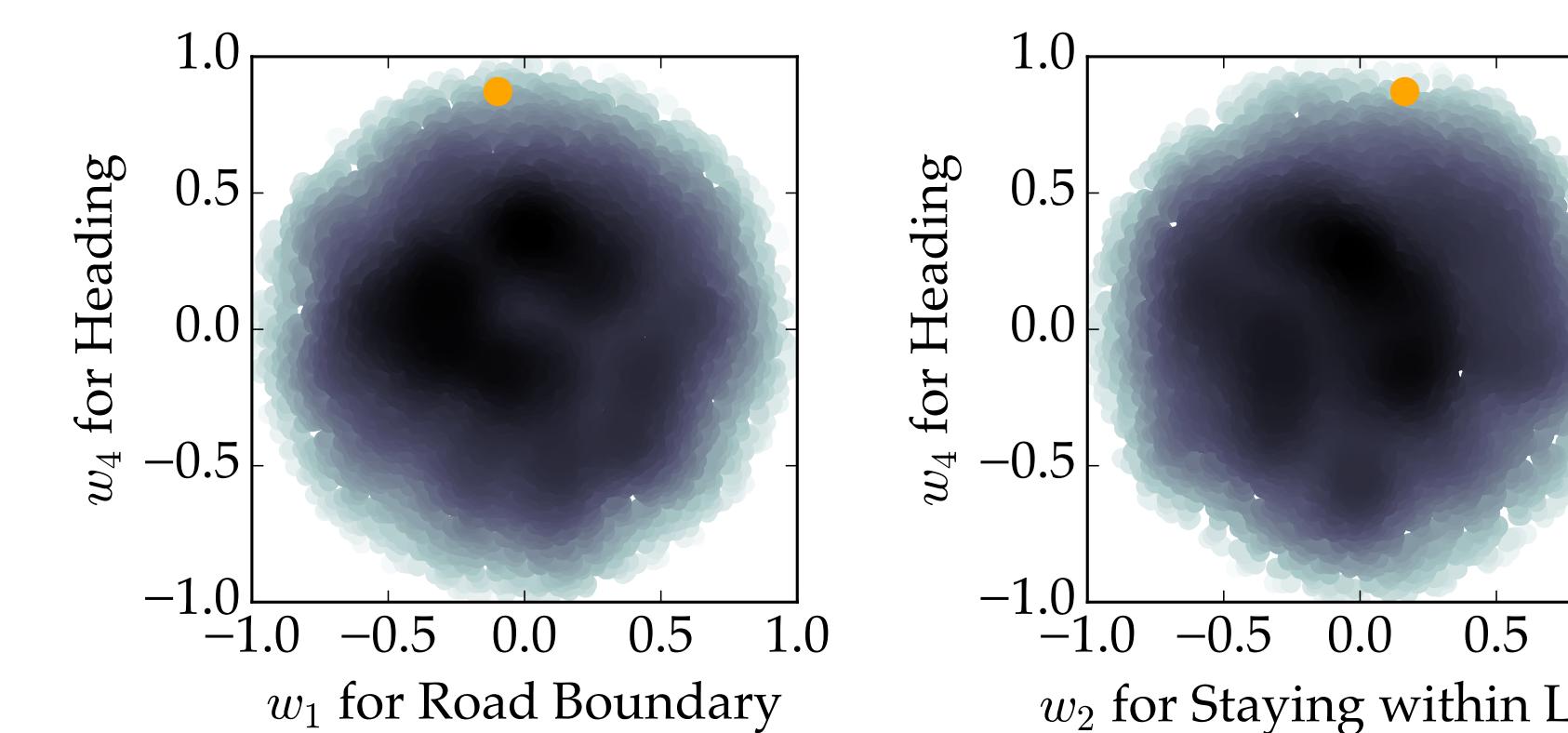
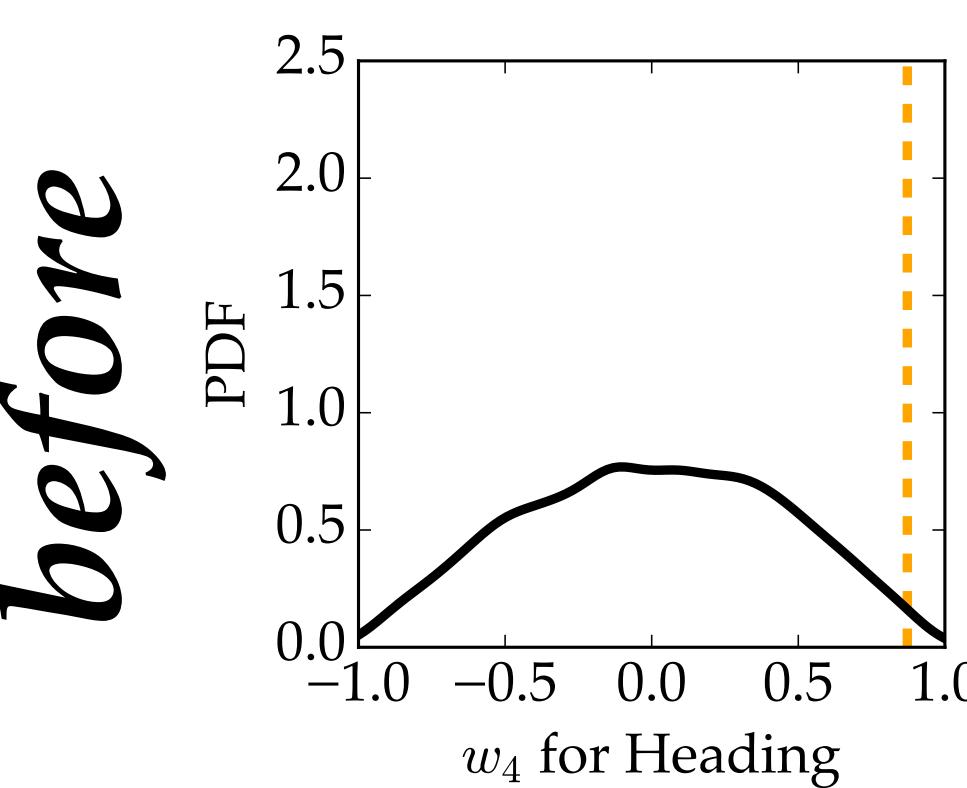
$$f_{\varphi}^1(\mathbf{w}) = p(I_t | \mathbf{w}) = \frac{1}{1 + \exp(-I_t \mathbf{w}^\top \varphi)}$$

$$f_{\varphi}^2(\mathbf{w}) = \min(1, \exp(I_t \mathbf{w}^\top \varphi))$$

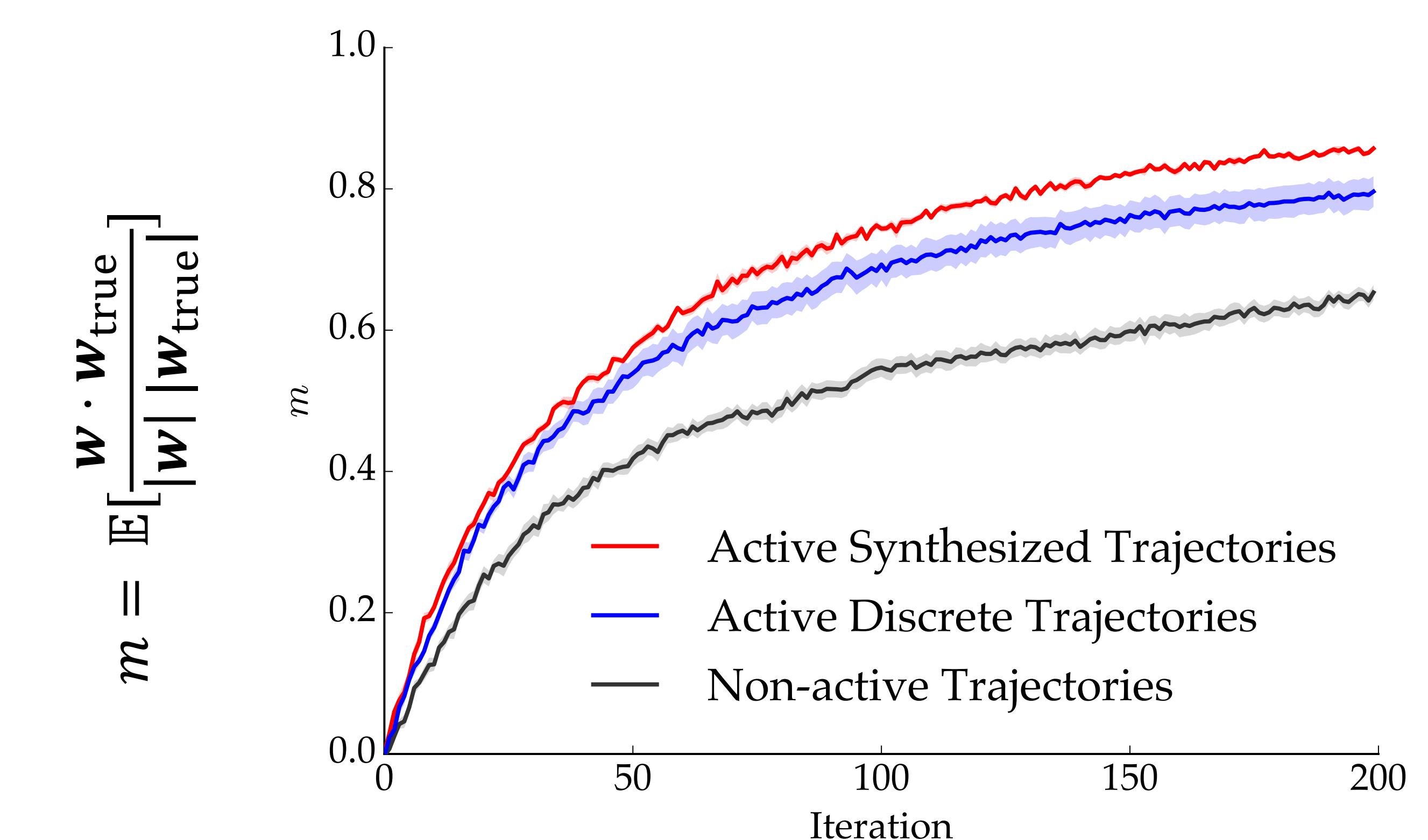
Log-concave update function for  $\mathbf{w}$  enables bound on the number of iterations to converge



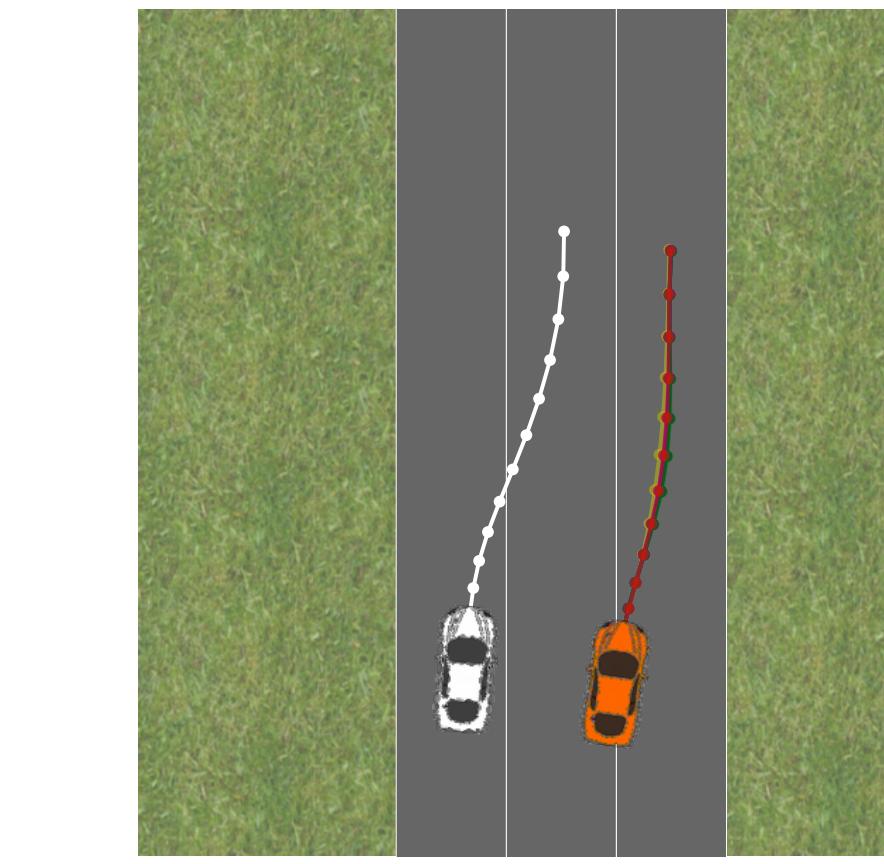
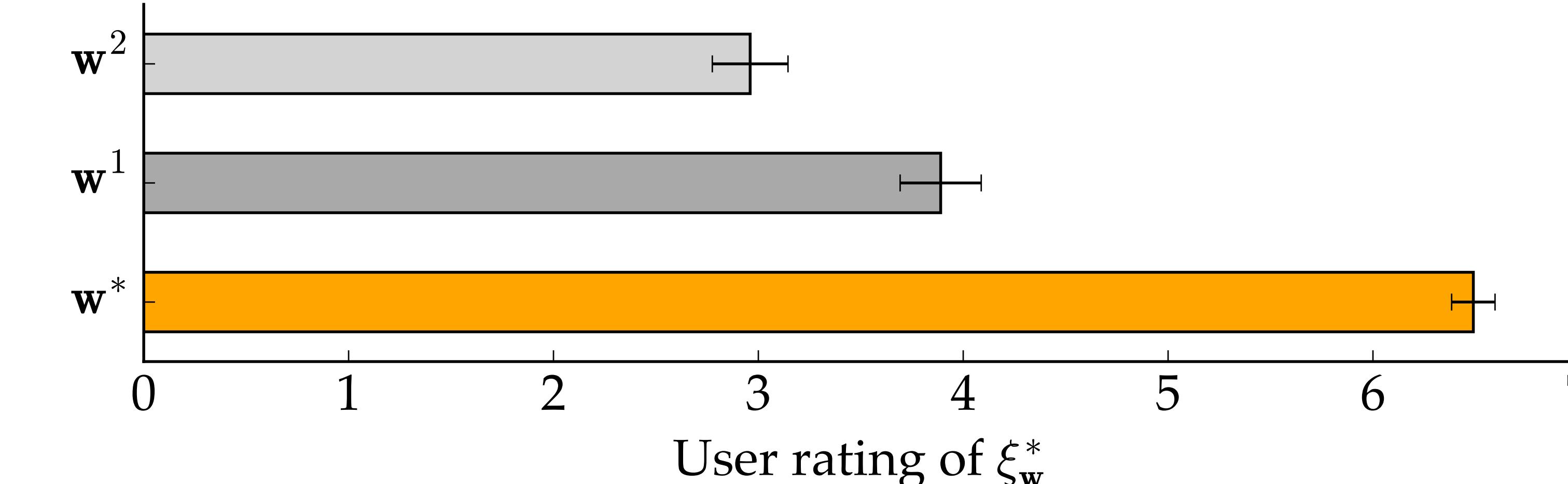
## Converging to ground-truth reward



## Active synthesis helps



... when queries = real trajectories



Learned  $\mathbf{w}^*$

