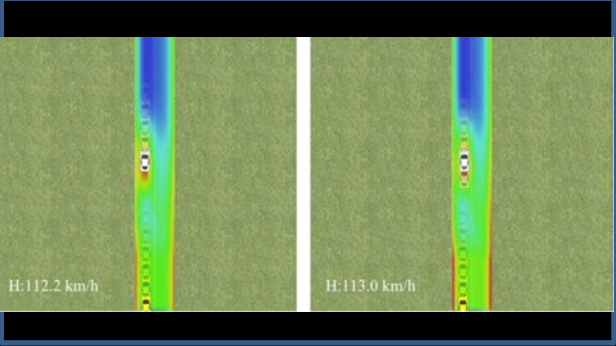
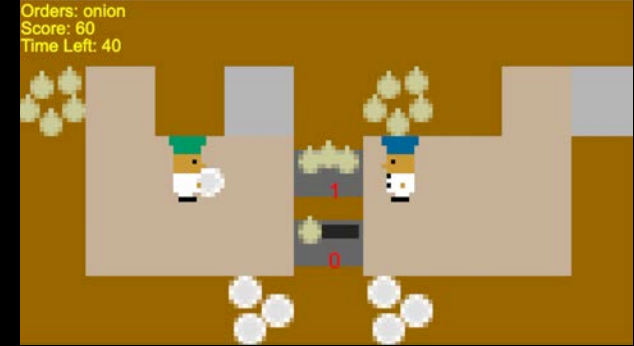
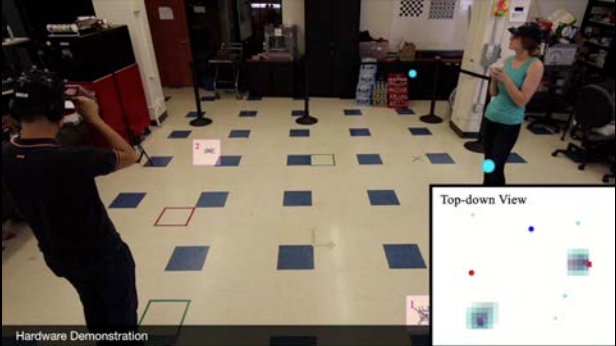


Optimizing Intended Reward Functions

Anca Dragan





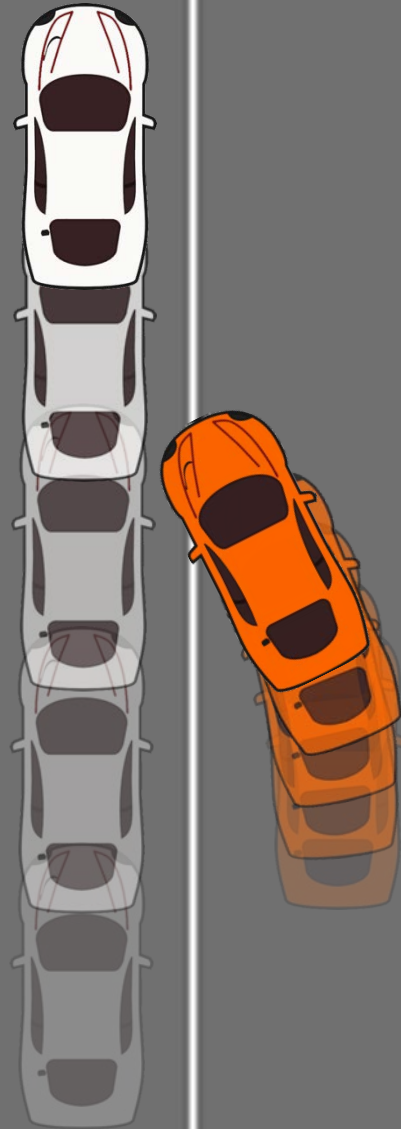
Example: Capture influence on human action



efficiency



Before: not the most efficient

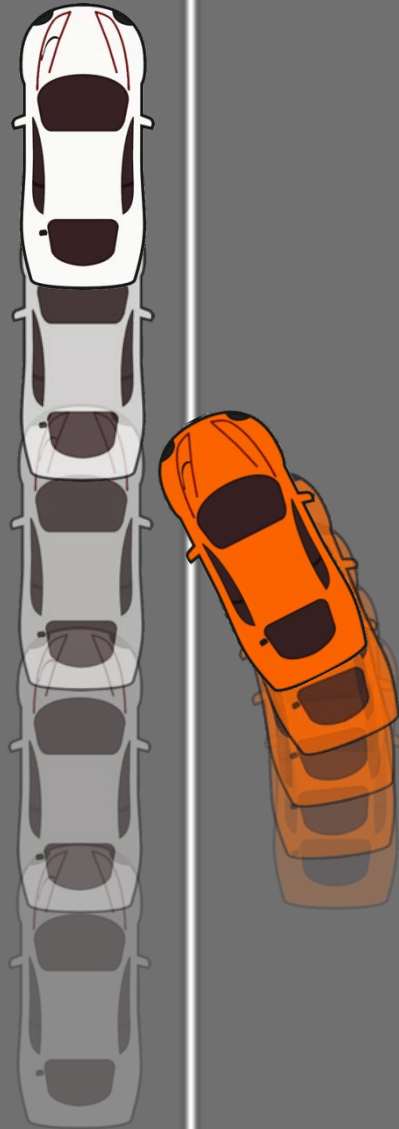


Before: not the most efficient



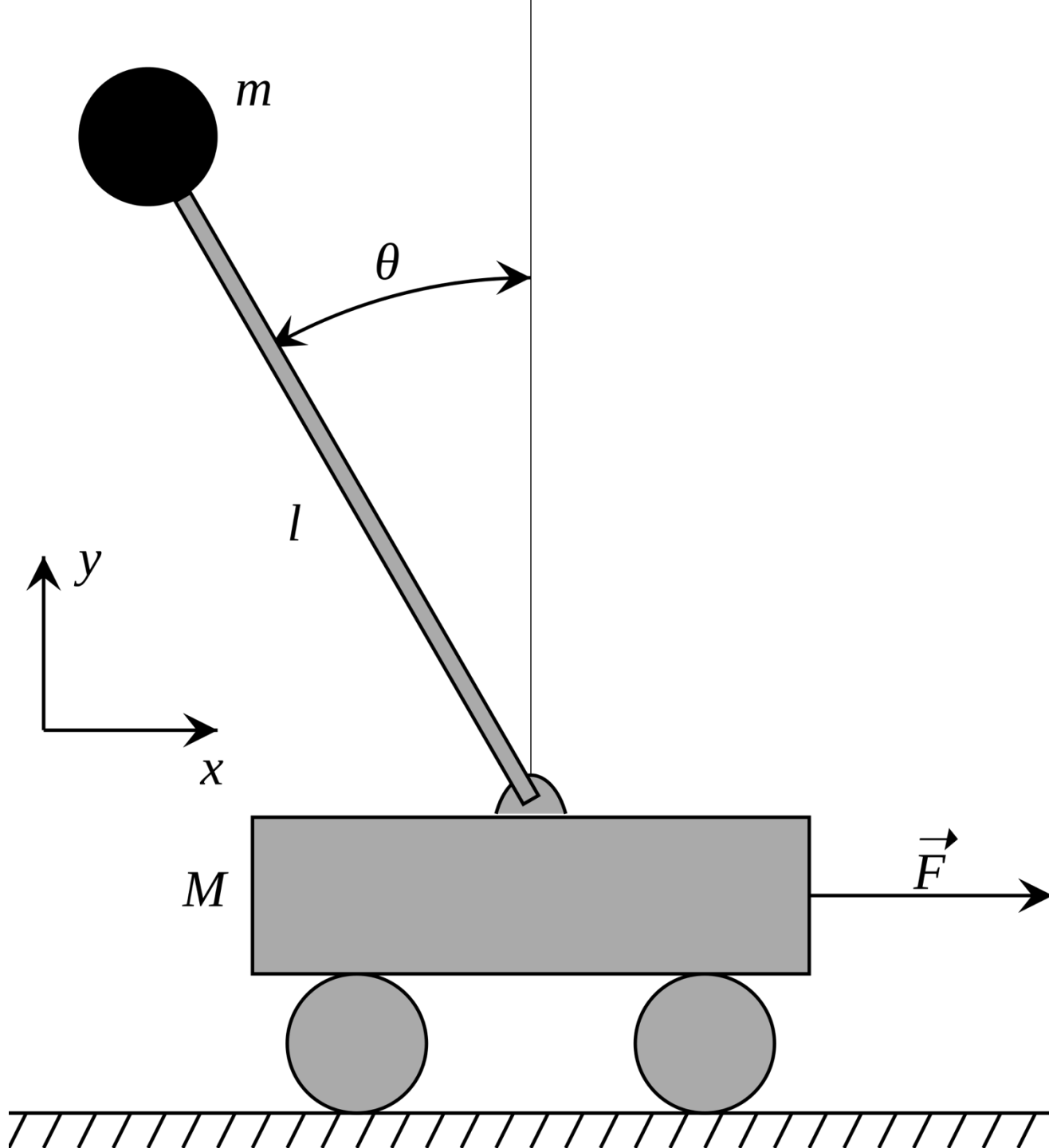


Before: not the most efficient

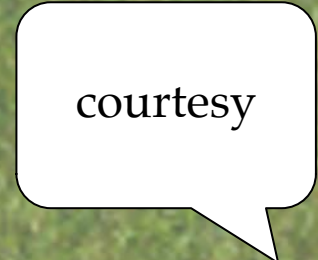
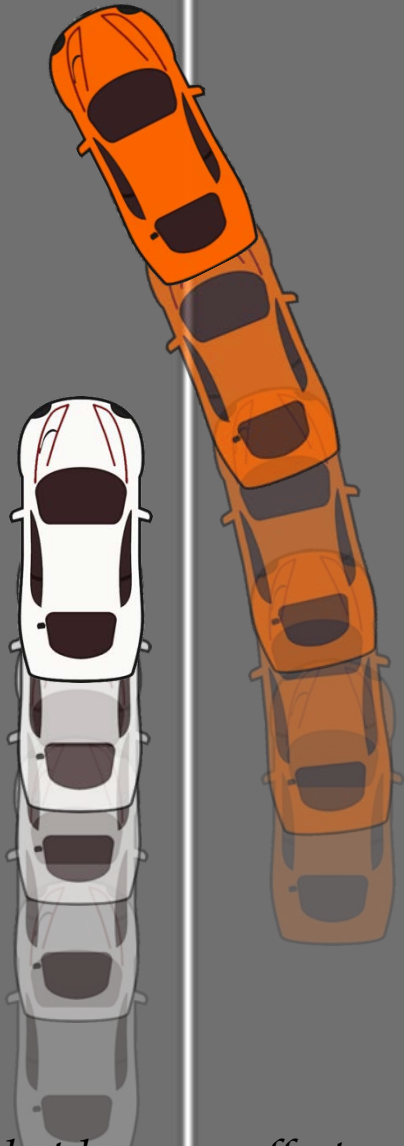


influence!





Now: sometimes "too" efficient



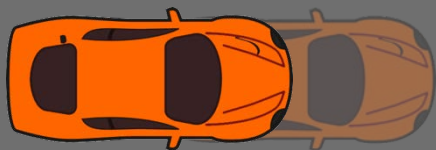
Planning for autonomous cars that leverage effect on human actions [RSS'16, with Sadigh, Sastry, Seshia]

Add courtesy..



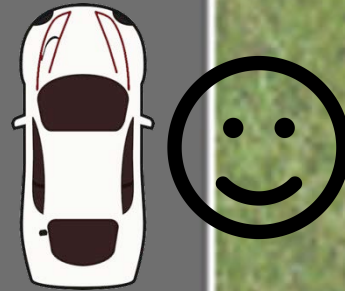
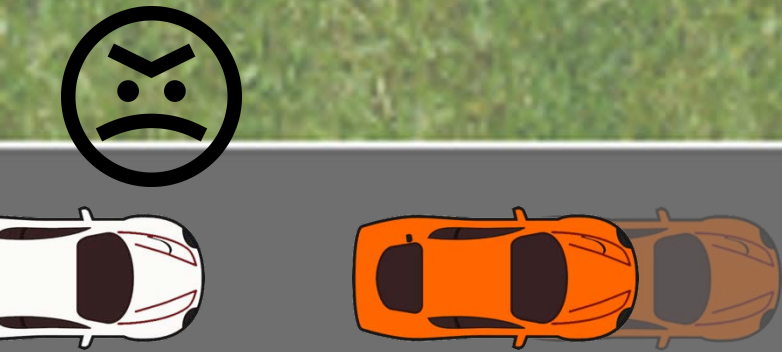
Planning for autonomous cars that leverage effect on human actions [RSS'16, with Sadigh, Sastry, Seshia]

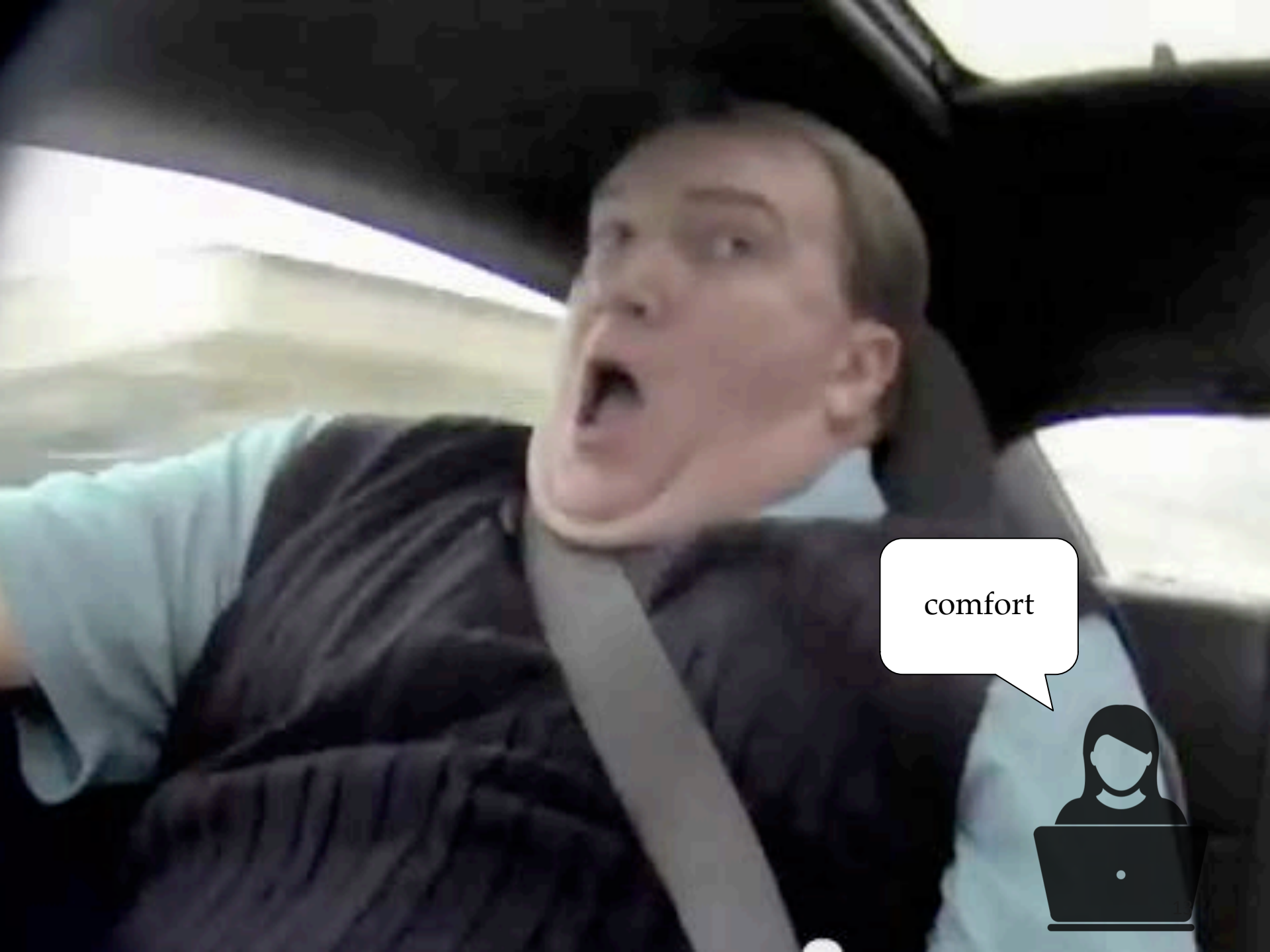
But now, car inches backwards to get you to go!



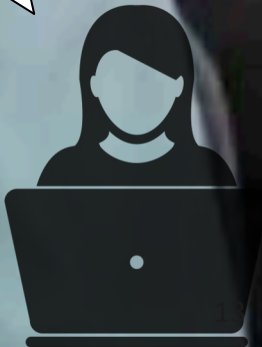
Planning for autonomous cars that leverage effect on human actions [RSS'16, with Sadigh, Sastry, Seshia]

But now, car inches backwards to get you to go!





comfort



*optimization, search, constraint
satisfaction, satisficing, RL...*

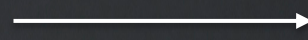
task specification \longrightarrow behavior

cost, reward, goal, loss, constraints,...

?????

*optimization, search, constraint
satisfaction, satisficing, RL...*

task **specification**



behavior

cost, reward, goal, loss, constraints,...



SCORE

0

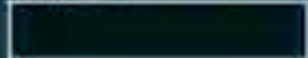
LAPS

-/3

TIME

0:01

TURBO



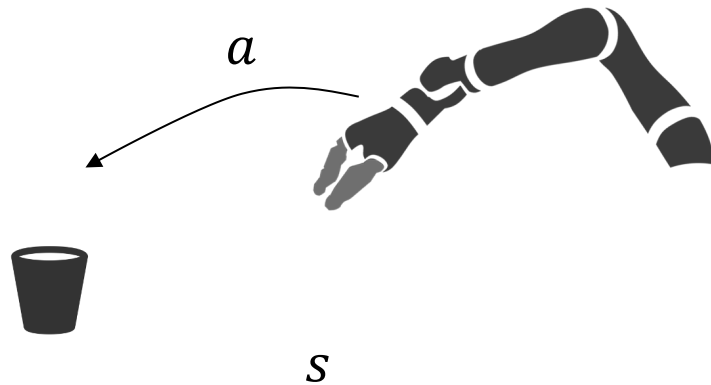
[MORE GAMES](#)



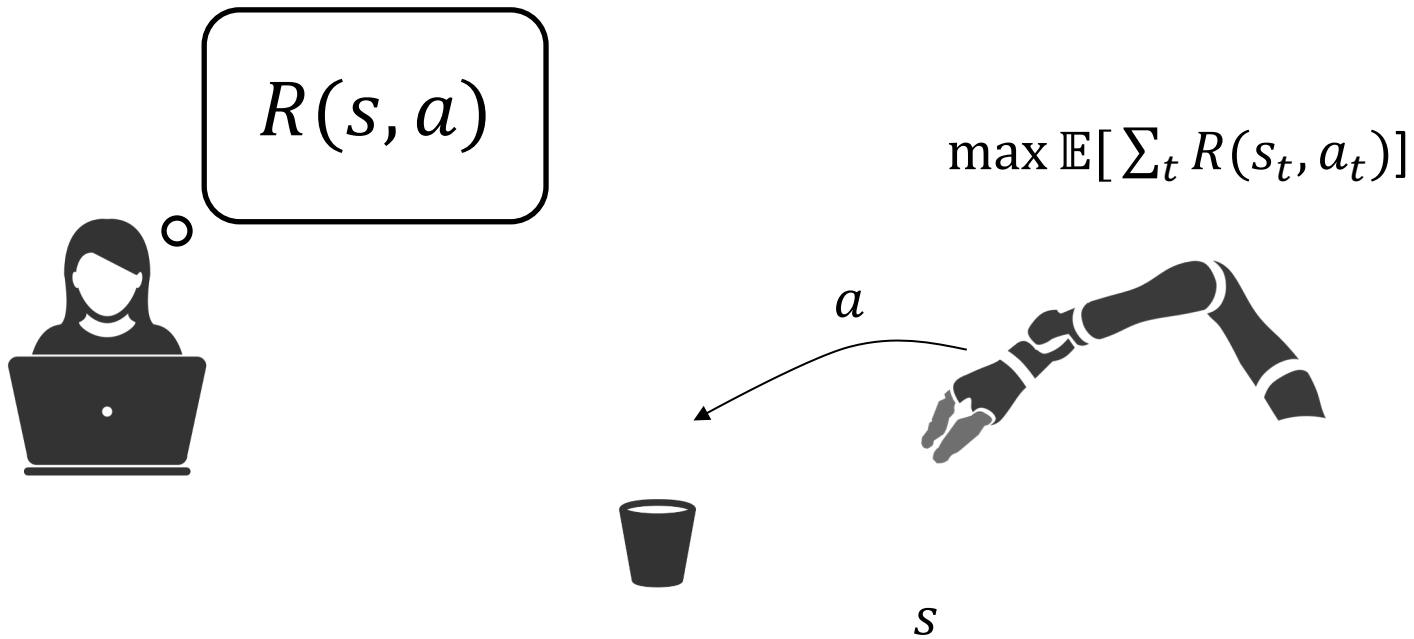
What we pretend AI is:

$$R(s, a)$$

$$\max \mathbb{E}[\sum_t R(s_t, a_t)]$$



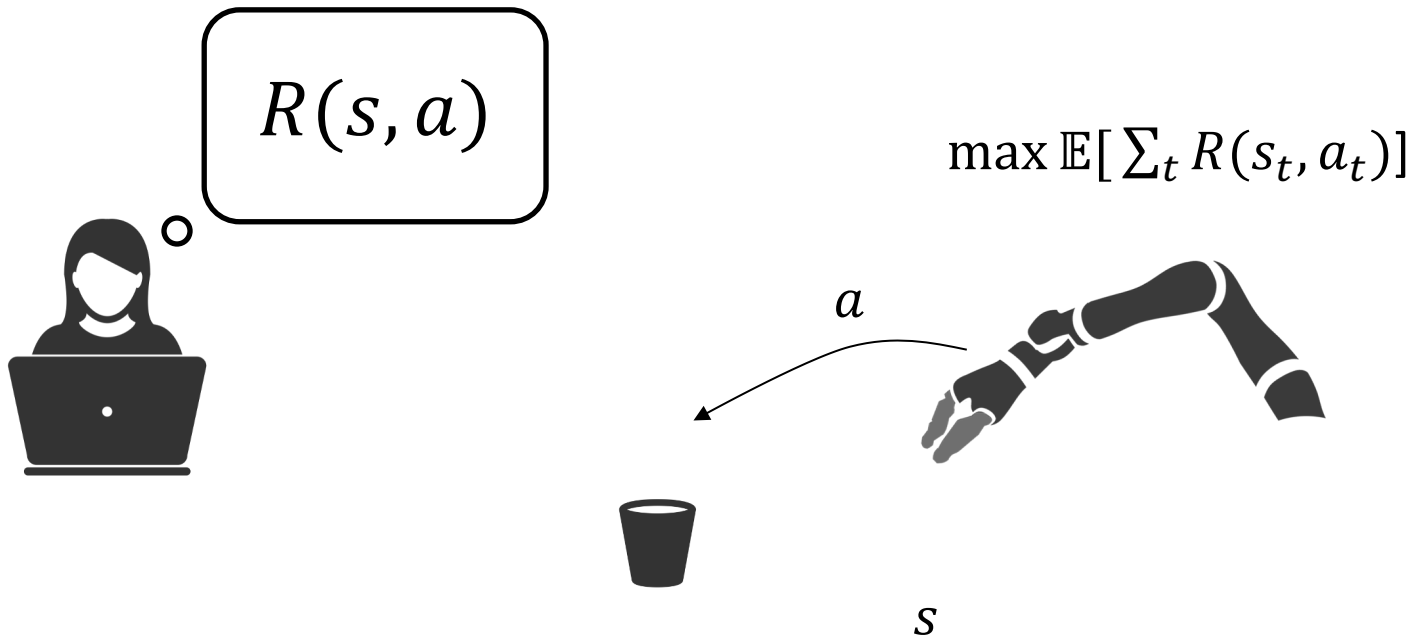
What AI actually is:



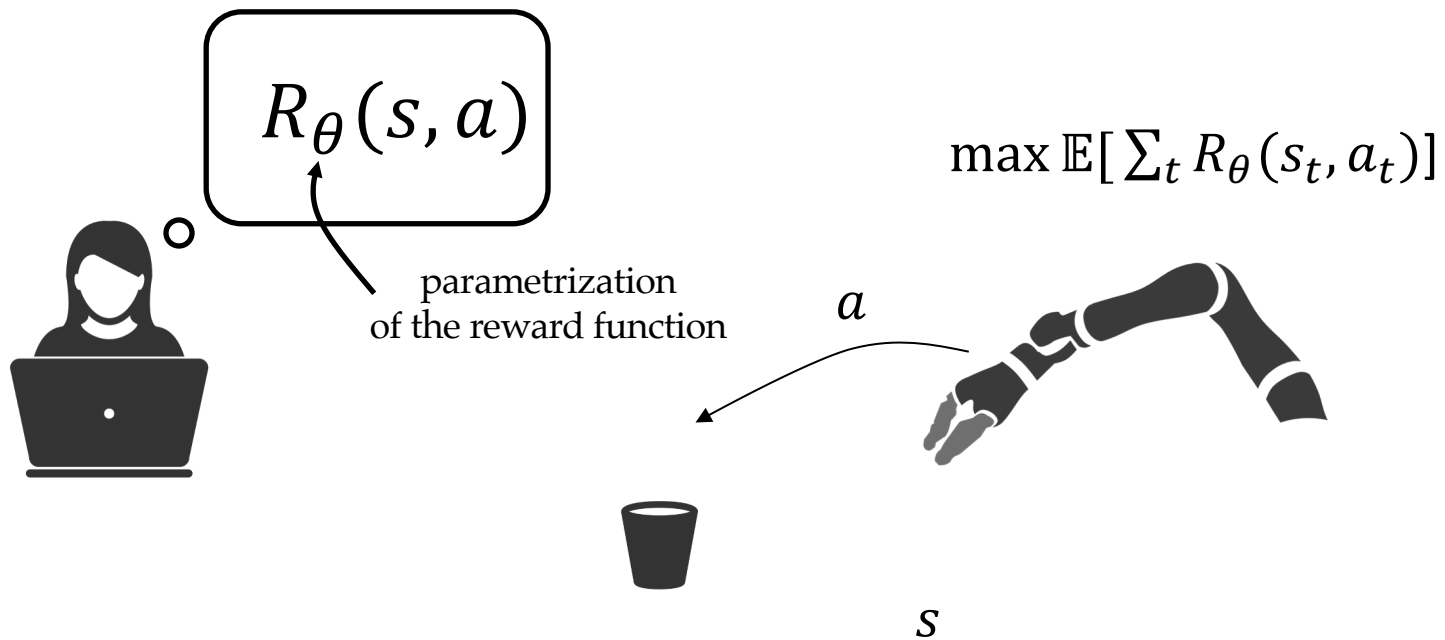
AI \neq optimize specified reward

AI = optimize intended reward

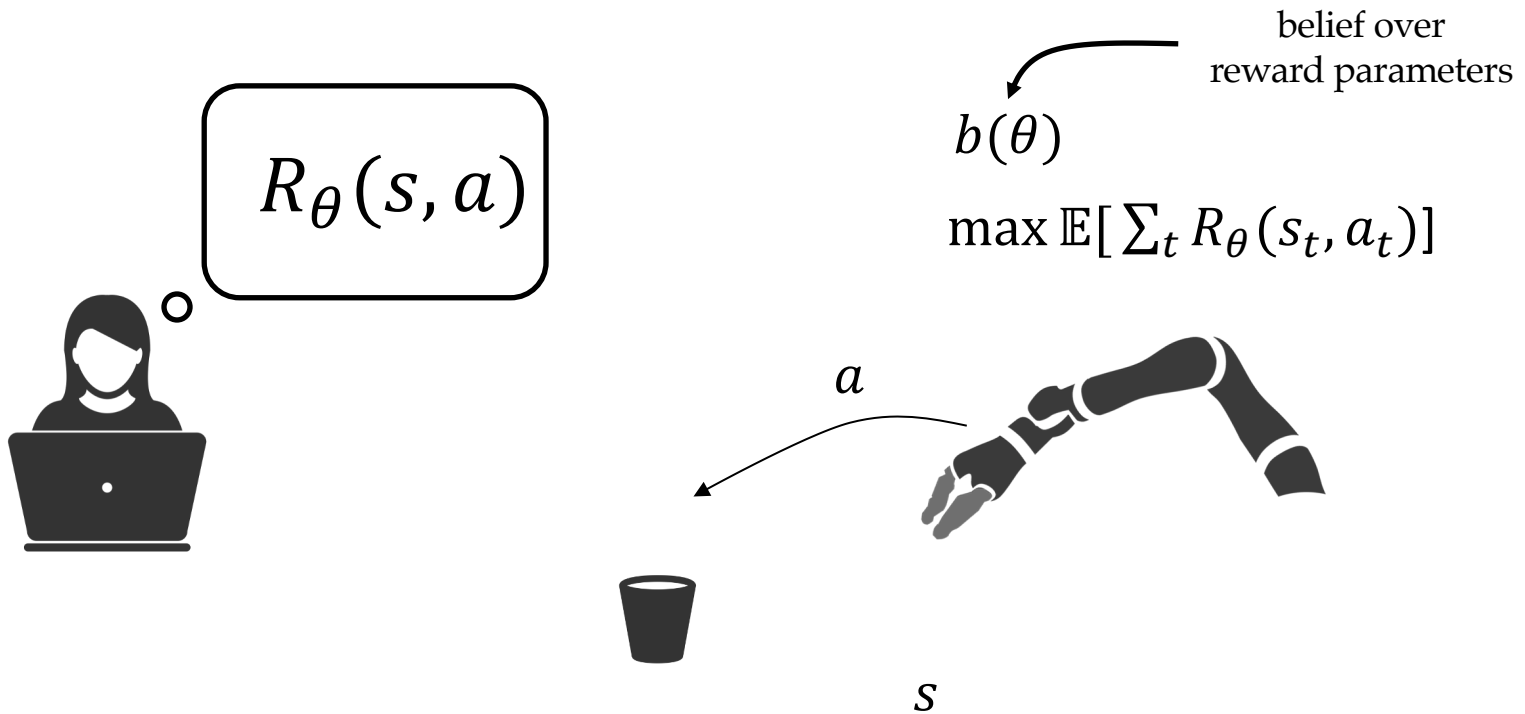
Optimize intended reward



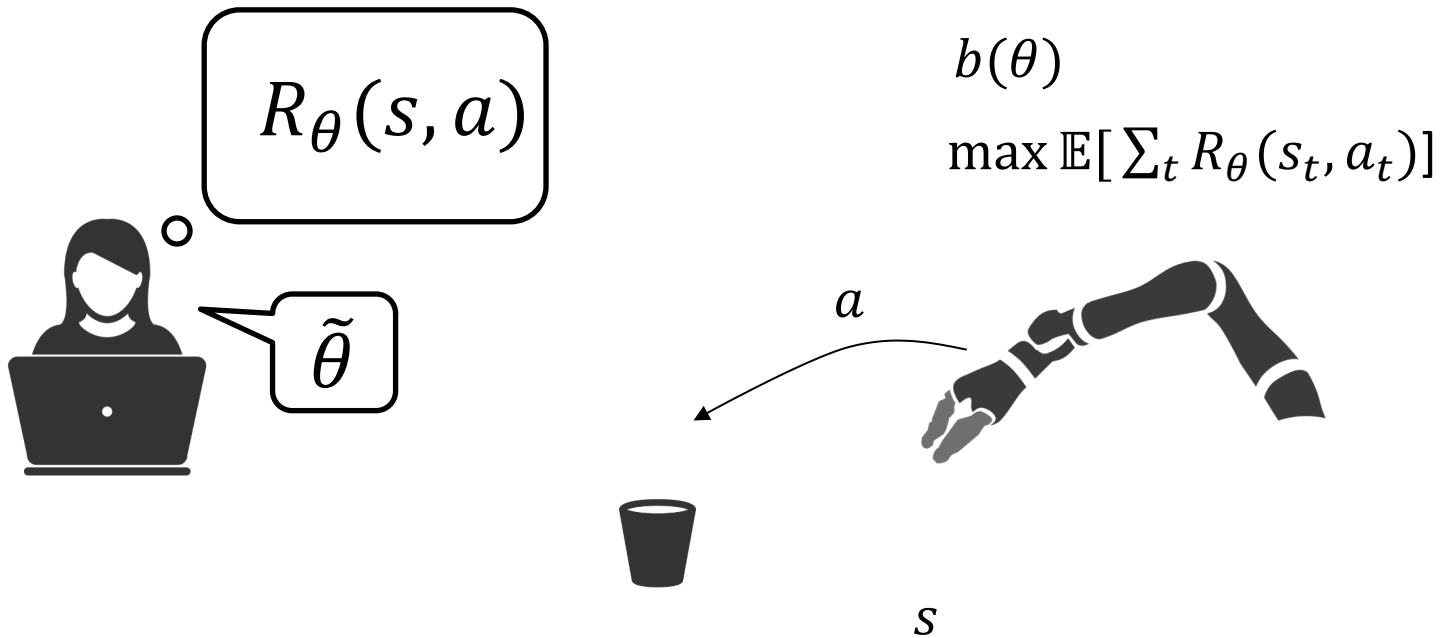
Optimize intended reward



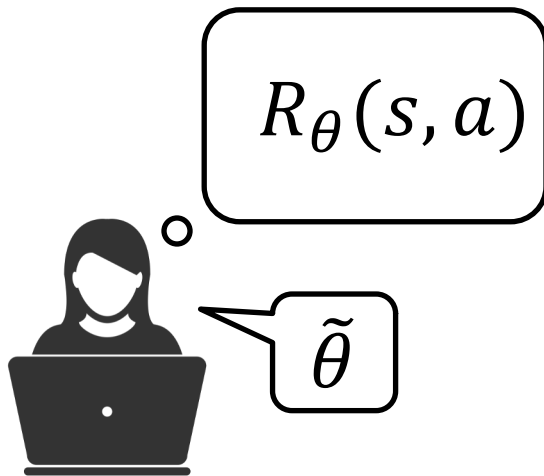
Optimize intended reward



Why treat specified rewards as definition?

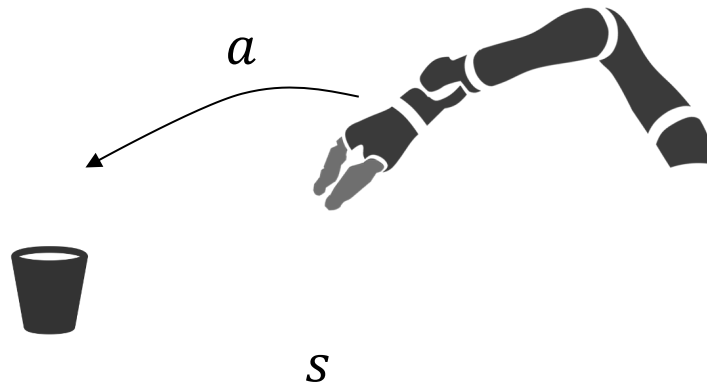


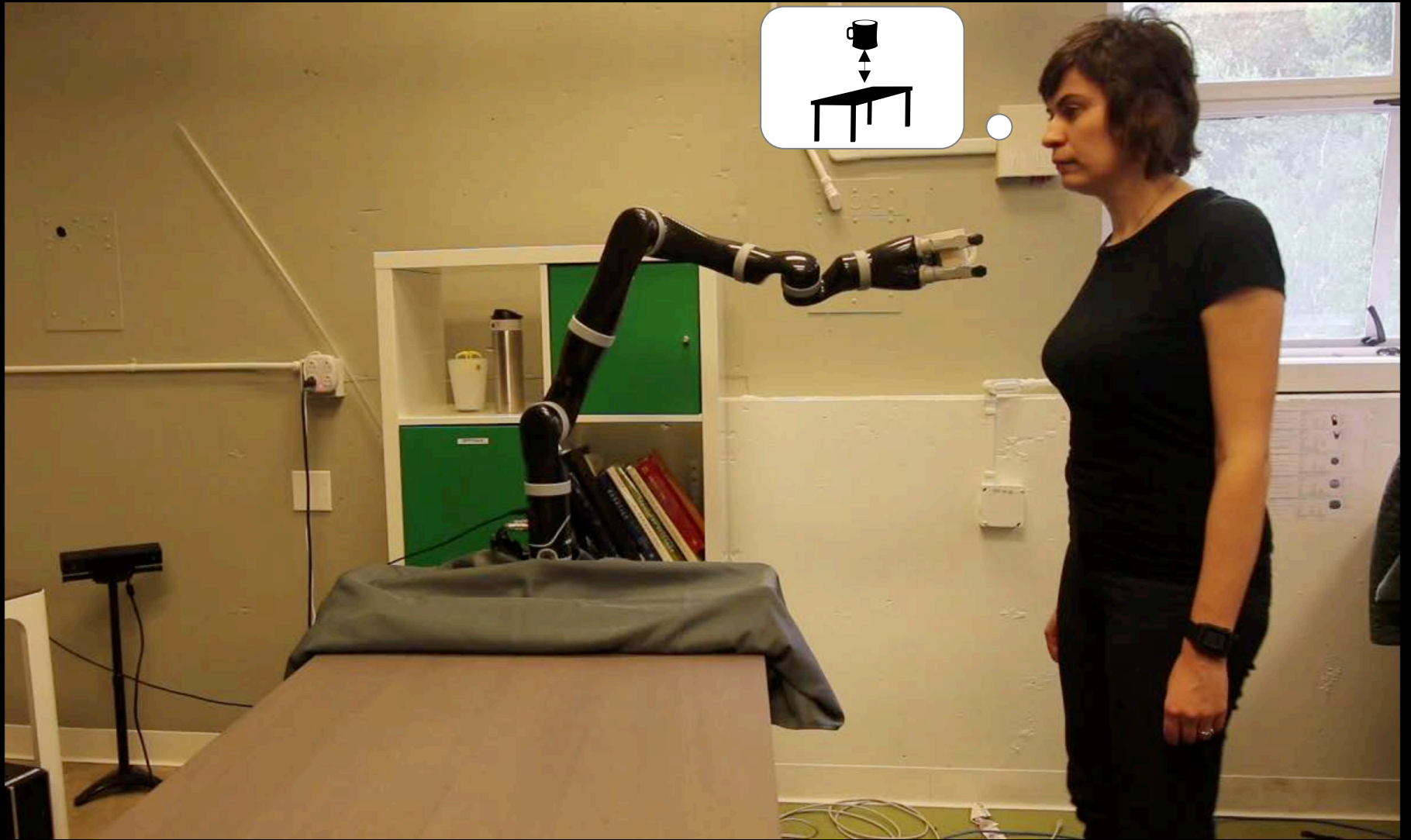
Why treat specified rewards as definition?

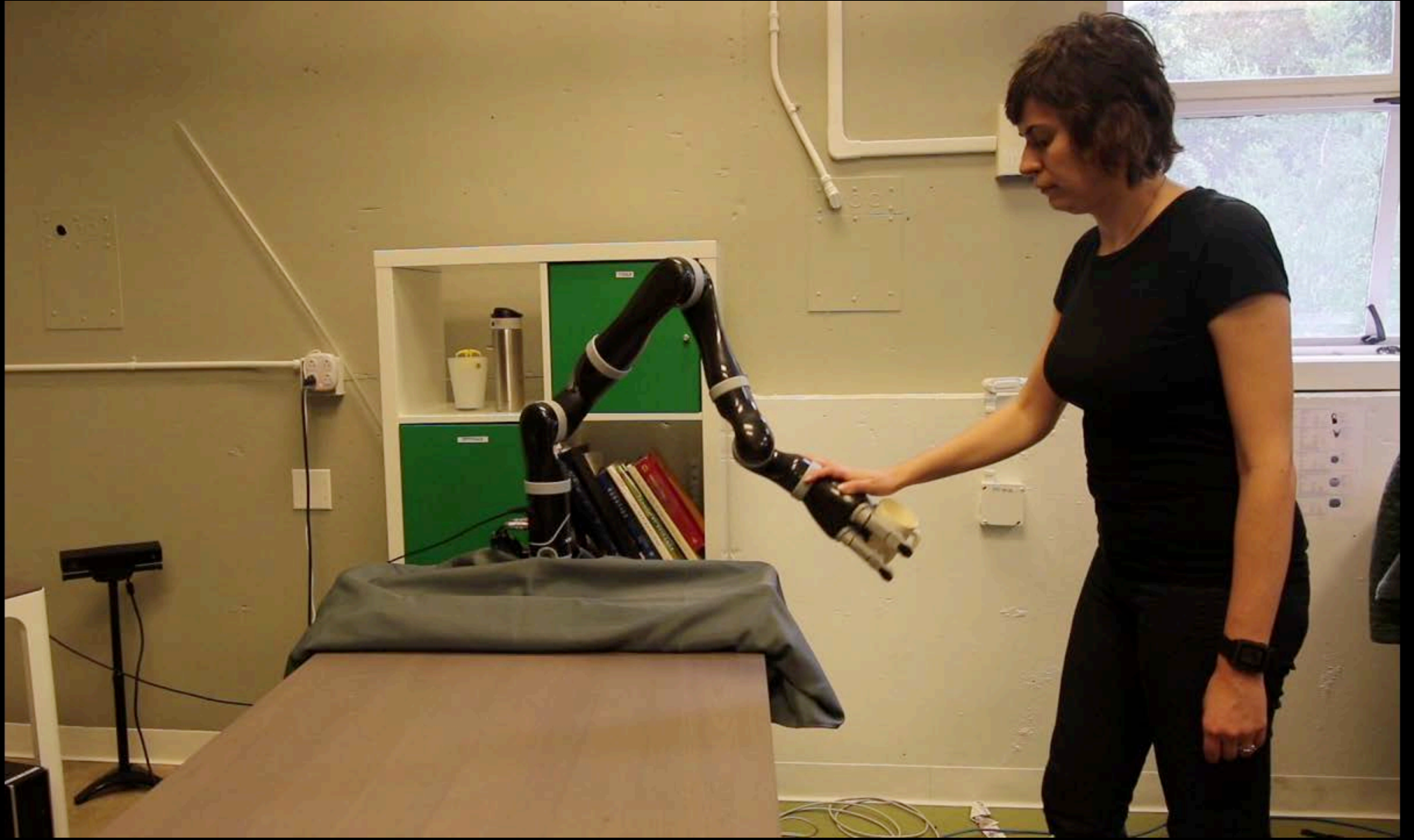


$$b(\theta) = \begin{cases} 1, & \theta = \tilde{\theta} \\ 0, & \text{else.} \end{cases}$$

$$\max \mathbb{E}[\sum_t R_{\theta}(s_t, a_t)]$$



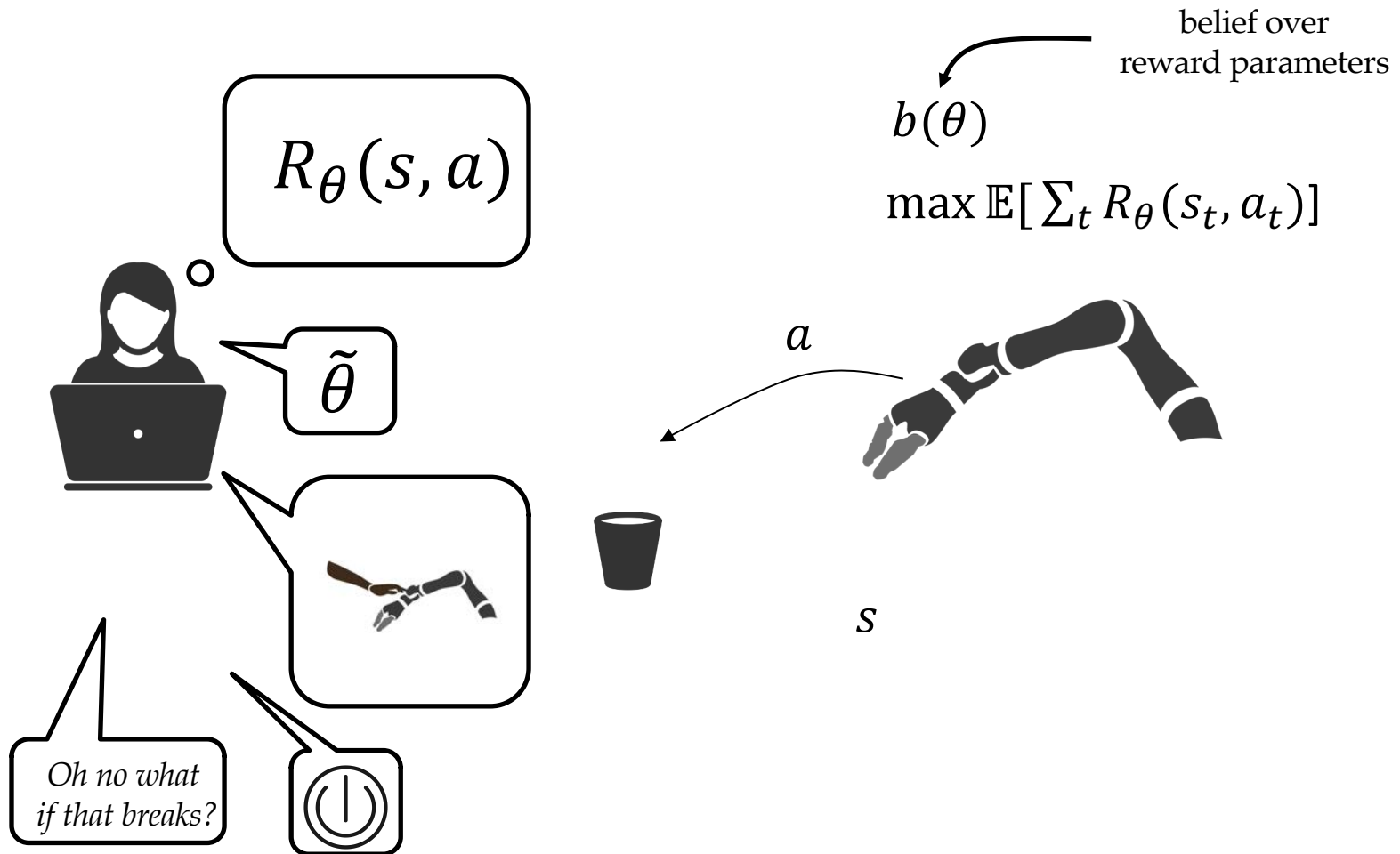




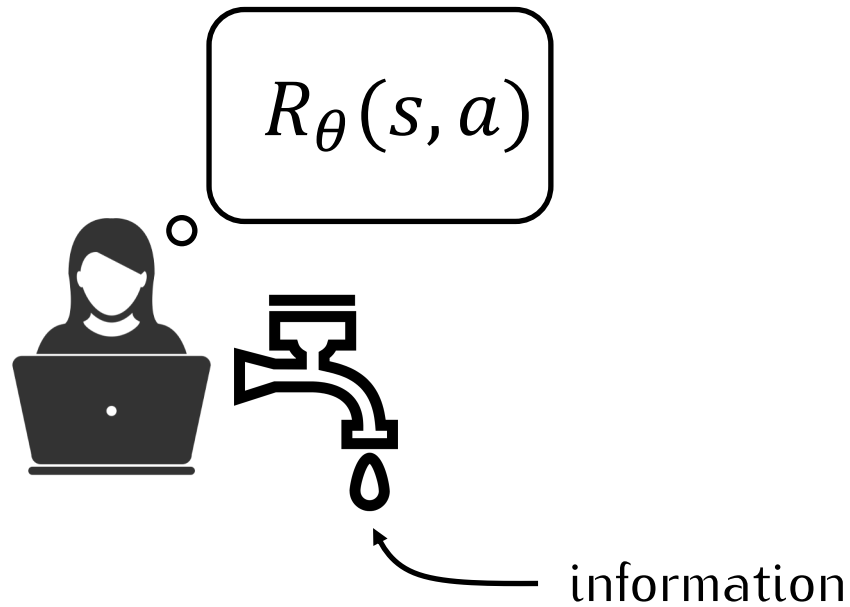


Agents overlearn from specified rewards,
but underlearn from other sources.

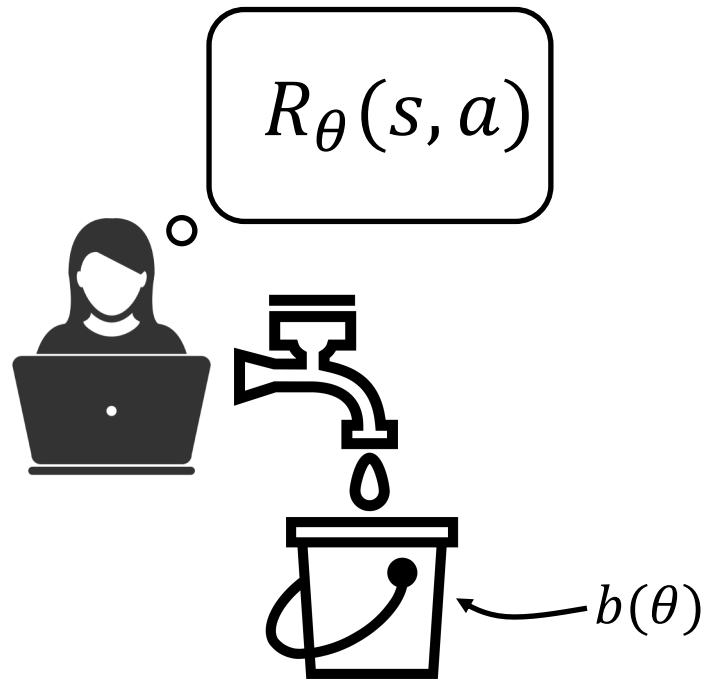
Optimize intended reward





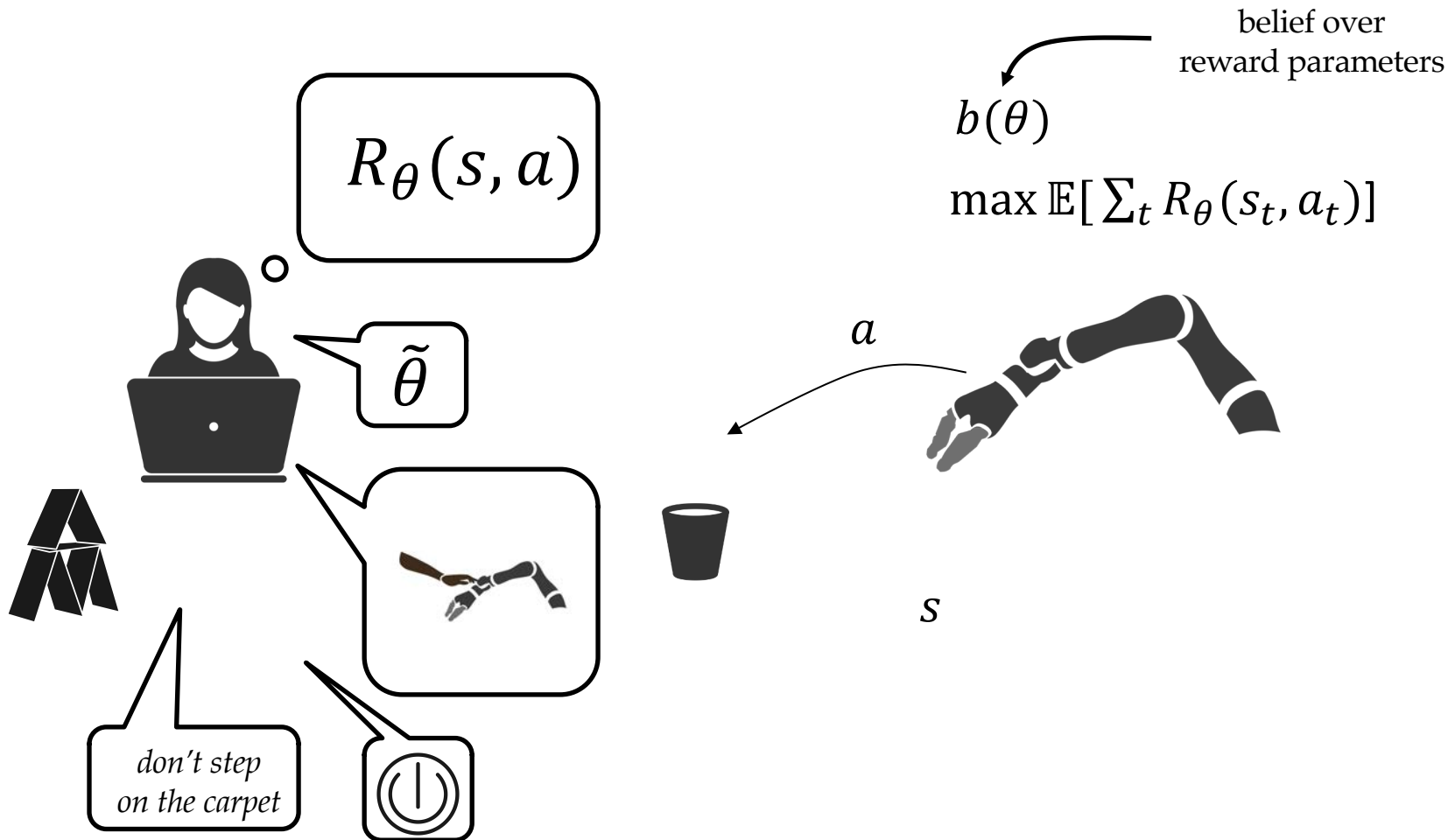


Humans leak information about the reward.

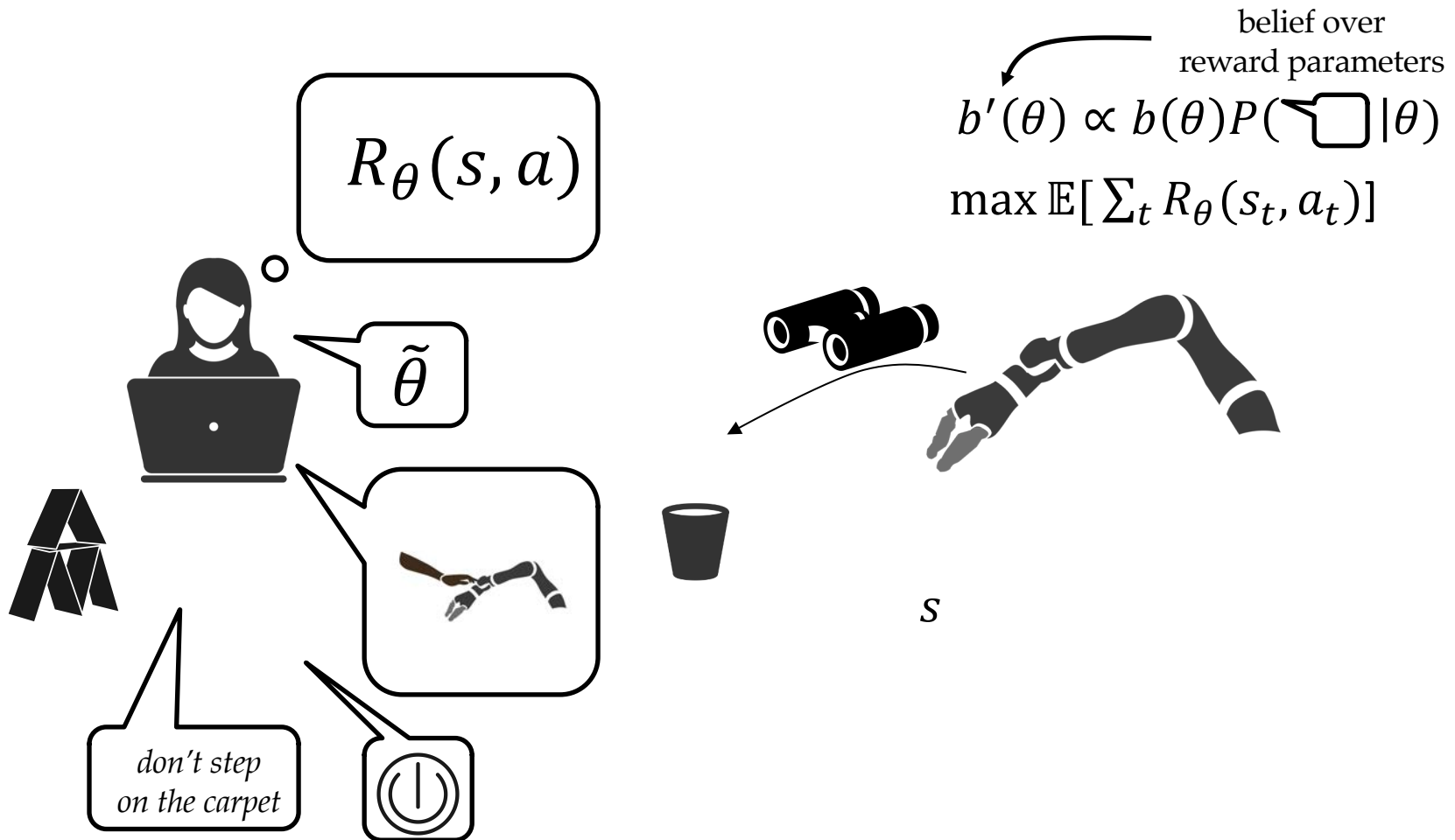


How should the robot extract it into an updated belief?

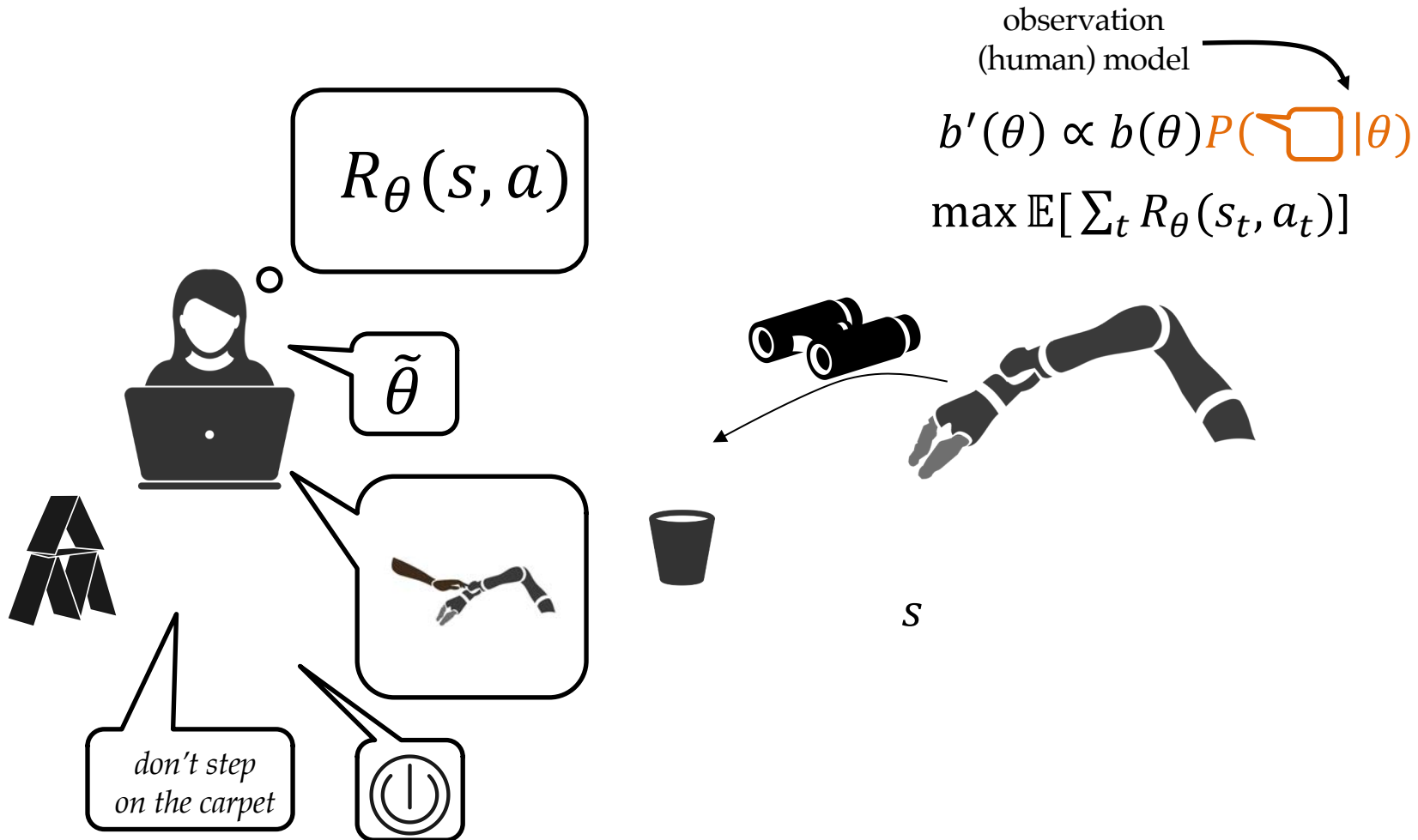
Human feedback, from specifying a reward to turning the robot off, is evidence about the intended reward.



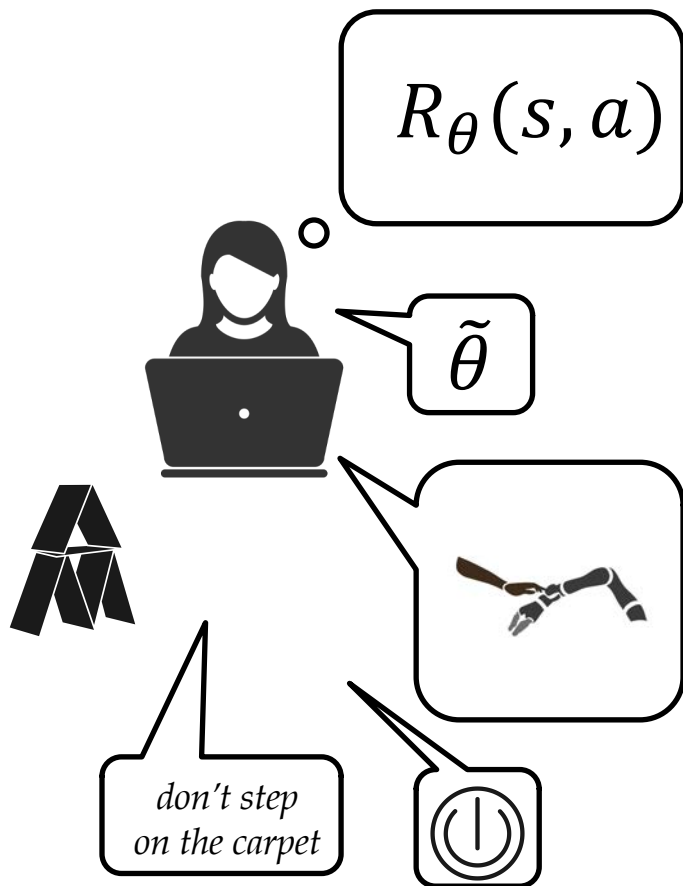
Human feedback, from specifying a reward to turning the robot off, is evidence about the intended reward.



Human feedback, from specifying a reward to turning the robot off, is evidence about the intended reward.



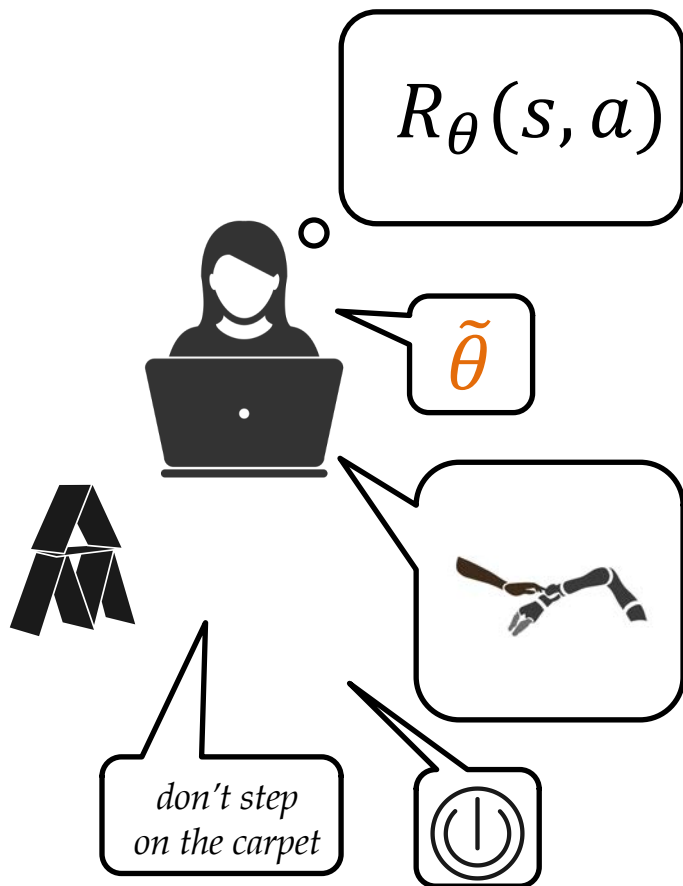
What is a human model that can be used to make sense of all these types of human feedback?



observation
(human) model

$$b'(\theta) \propto b(\theta) P(\text{observation} | \theta)$$

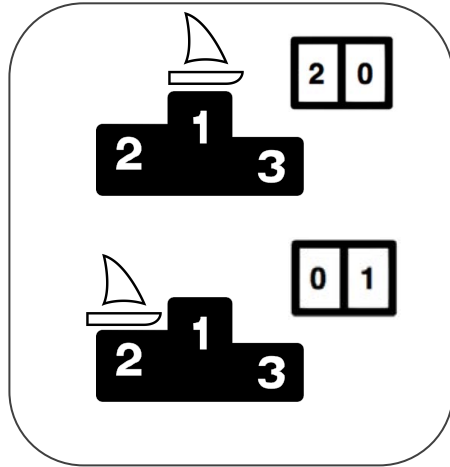
How can we model reward design/specification as a noisy and suboptimal process?



observation
(human) model

$$b'(\theta) \propto b(\theta)P(\tilde{\theta} | \theta)$$

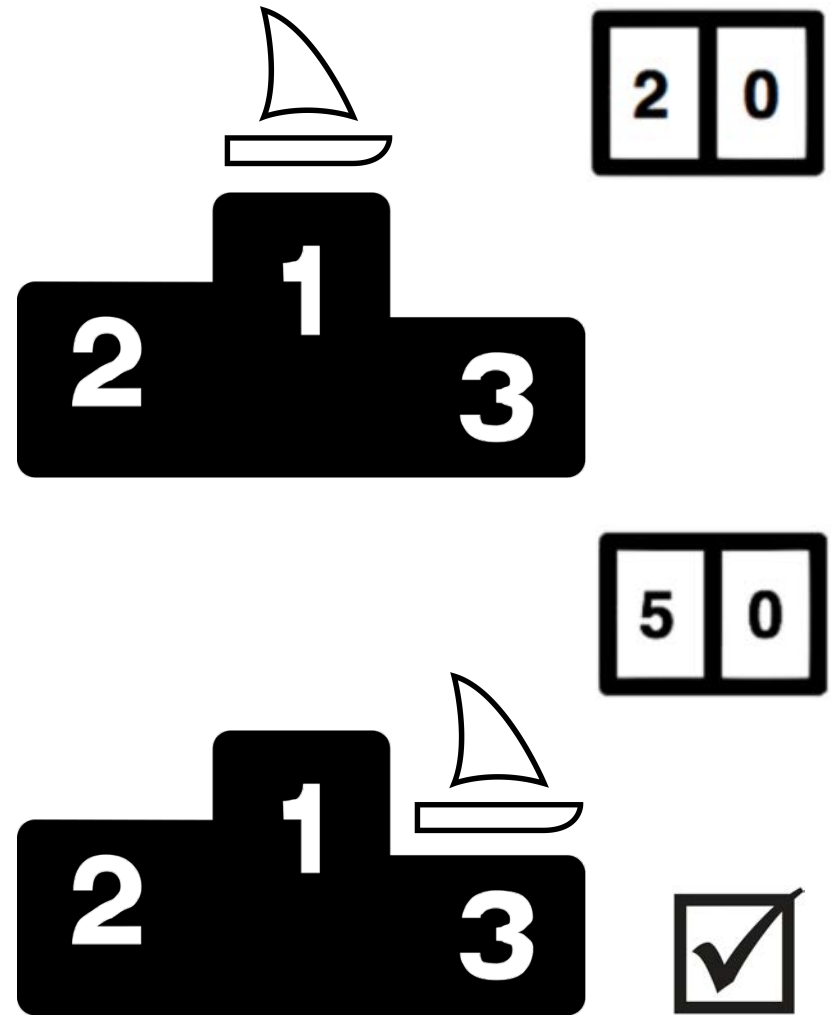
Development



$R_{\tilde{\theta}}$ - score

score and winning were
correlated at development time...

Deployment



... but no longer
correlated at deployment time

We only know this about the true reward:

The behavior incentivized by the specified reward in development has high true reward.

What you specify is contextualized by the state you specify it in. Robots should interpret it as such.

The behavior incentivized by the specified reward in development has high true reward

$$P(\tilde{\theta} | \theta^*, M_{devel}) \propto e^{\beta \mathbb{E}[R_{\theta^*}(\xi; M_{devel}) | \xi \sim P(\xi | \tilde{\theta}, M_{devel})]}$$

The behavior incentivized by the specified reward in development has high true reward

$$P(\tilde{\theta} | \theta^*, M_{devel}) \propto e^{\beta \mathbb{E}[R_{\theta^*}(\xi; M_{devel}) | \xi \sim P(\xi | \tilde{\theta}, M_{devel})]}$$

The behavior incentivized by the specified reward in development has high true reward

(approximately)
optimal trajectories 

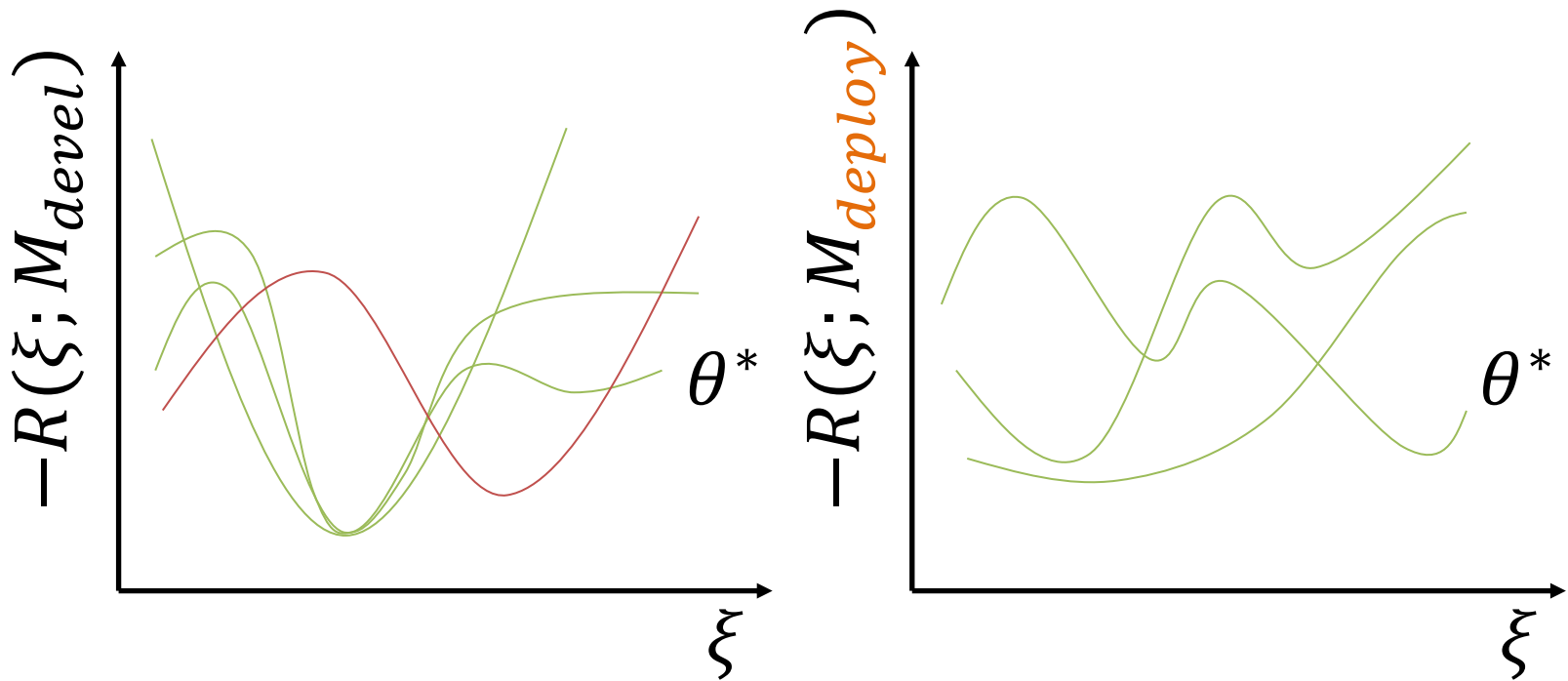
$$P(\tilde{\theta} | \theta^*, M_{devel}) \propto e^{\beta \mathbb{E}[R_{\theta^*}(\xi; M_{devel}) | \xi \sim P(\xi | \tilde{\theta}, M_{devel})]}$$

The behavior incentivized by the specified reward in development has high true reward

$$P(\tilde{\theta} | \theta^*, M_{devel}) \propto e^{\beta \mathbb{E}[R_{\theta^*}(\xi; M_{devel}) | \xi \sim P(\xi | \tilde{\theta}, M_{devel})]}$$

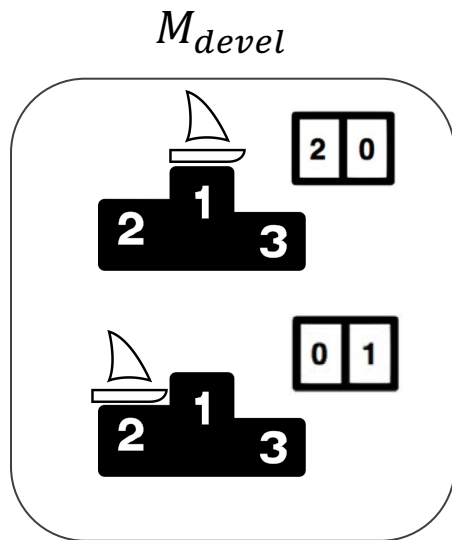
The behavior incentivized by the specified cost in development has low true cost

$$P(\tilde{\theta} | \theta^*, M_{devel}) \propto e^{\beta \mathbb{E}[R_{\theta^*}(\xi; M_{devel}) | \xi \sim P(\xi | \tilde{\theta}, M_{devel})]}$$



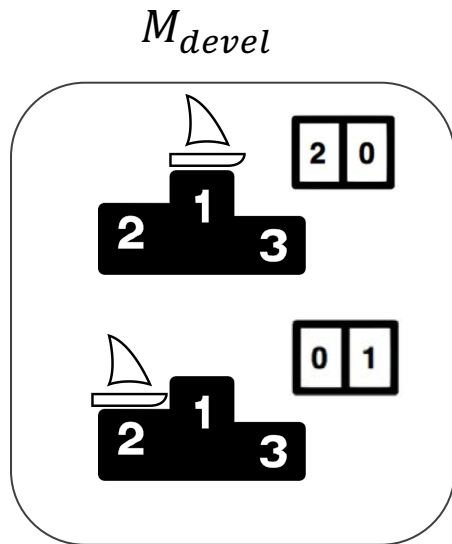
The behavior incentivized by the specified cost in development has low true cost

$$P(\tilde{\theta} | \theta^*, M_{devel}) \propto e^{\beta \mathbb{E}[R_{\theta^*}(\xi; M_{devel}) | \xi \sim P(\xi | \tilde{\theta}, M_{devel})]}$$



The behavior incentivized by the specified cost in development has low true cost

$$P(\tilde{\theta} | \theta^*, M_{devel}) \propto e^{\beta \mathbb{E}[R_{\theta^*}(\xi; M_{devel}) | \xi \sim P(\xi | \tilde{\theta}, M_{devel})]}$$



θ_1

maximizing
winning

θ_2

maximizing
score

θ_3

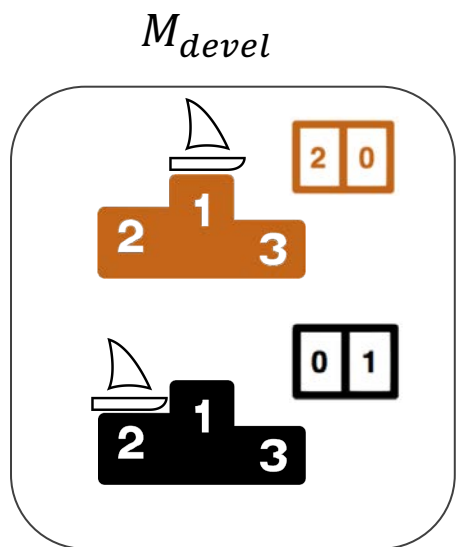
minimizing
winning

θ_4

minimizing
score

The behavior incentivized by the specified cost in development has low true cost

$$P(\tilde{\theta} | \theta^*, M_{devel}) \propto e^{\beta \mathbb{E}[R_{\theta^*}(\xi; M_{devel}) | \xi \sim P(\xi | \tilde{\theta}, M_{devel})]}$$



θ_1

maximizing
winning

θ_2

maximizing
score

θ_3

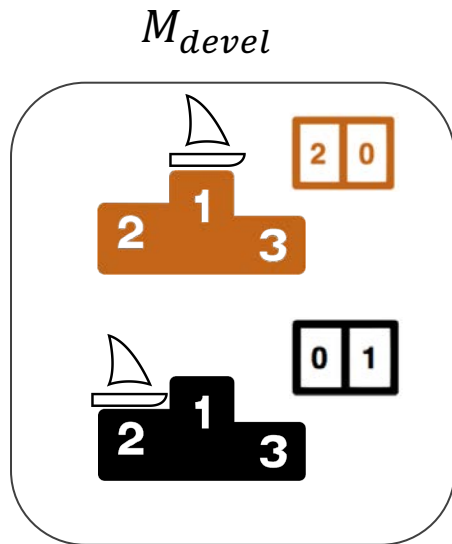
minimizing
winning

θ_4

minimizing
score

The behavior incentivized by the specified cost in development has low true cost

$$P(\tilde{\theta} | \theta^*, M_{devel}) \propto e^{\beta \mathbb{E}[R_{\theta^*}(\xi; M_{devel}) | \xi \sim P(\xi | \tilde{\theta}, M_{devel})]}$$



θ_1

maximizing
winning

θ_2

maximizing
score

θ_3

minimizing
winning

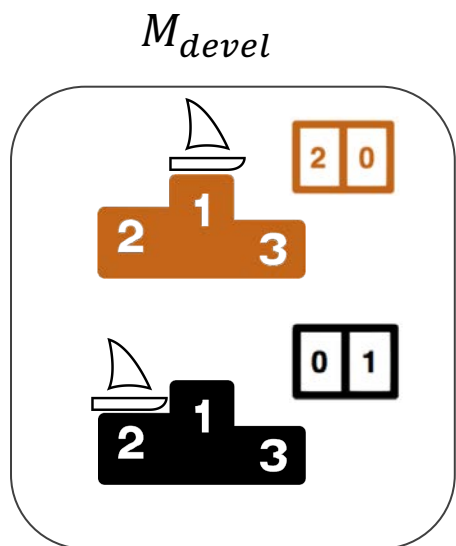
θ_4

minimizing
score



The behavior incentivized by the specified cost in development has low true cost

$$P(\tilde{\theta} | \theta^*, M_{devel}) \propto e^{\beta \mathbb{E}[R_{\theta^*}(\xi; M_{devel}) | \xi \sim P(\xi | \tilde{\theta}, M_{devel})]}$$



θ_1

maximizing
winning

θ_2

maximizing
score

θ_3

minimizing
winning

θ_4

minimizing
score

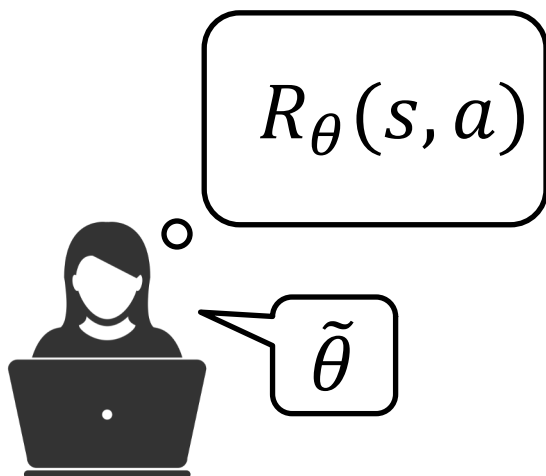


The behavior incentivized by the specified cost in development has low true cost

$$P(\tilde{\theta} | \theta^*, M_{devel}) \propto e^{\beta \mathbb{E}[R_{\theta^*}(\xi; M_{devel}) | \xi \sim P(\xi | \tilde{\theta}, M_{devel})]}$$



Specified rewards as evidence about the reward



risk-averse planning

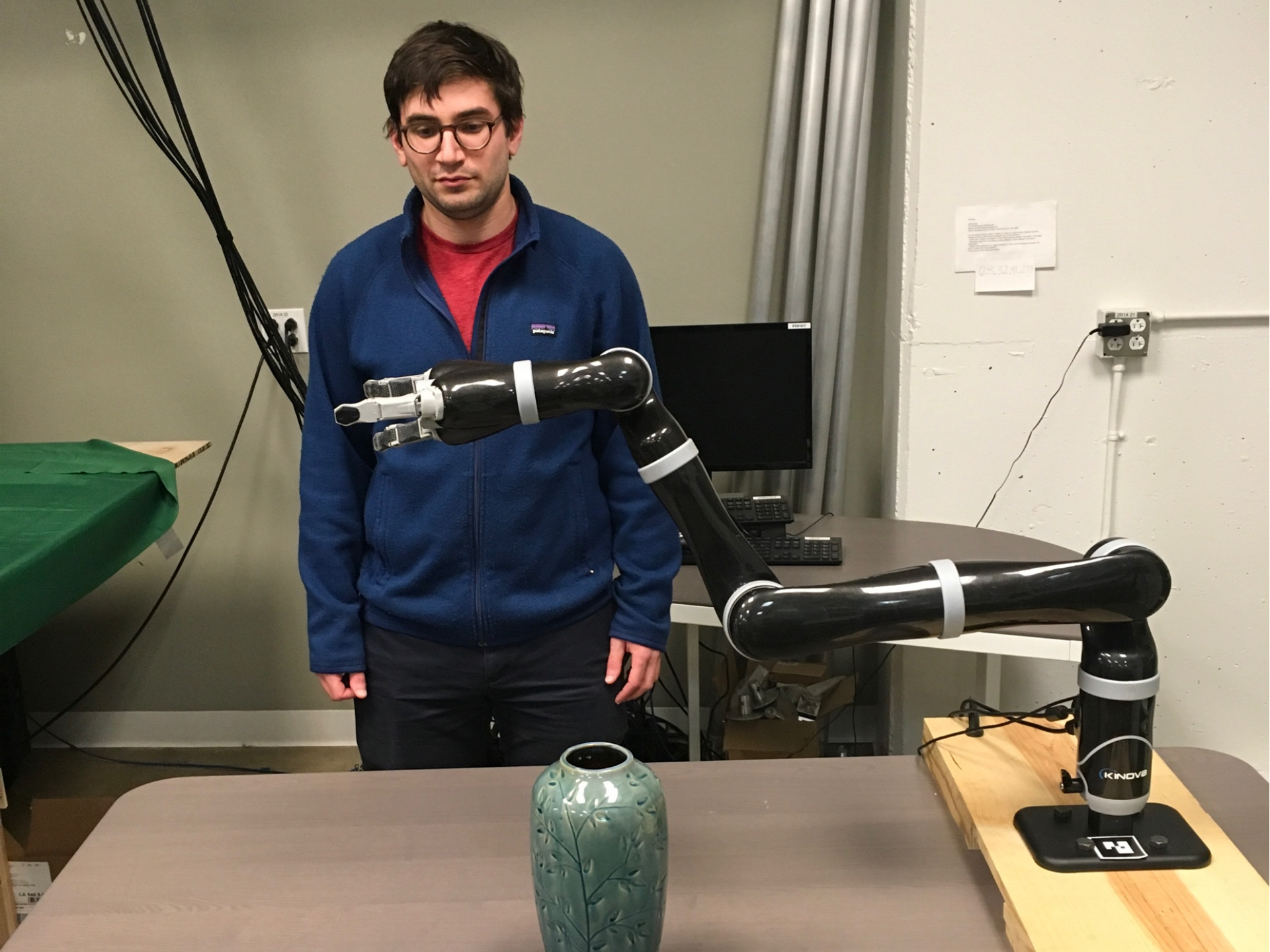
$$\max_{\xi} \min_{\theta \in \{\theta_i \sim b'(\theta)\}} R_\theta(\xi; M_{test})$$

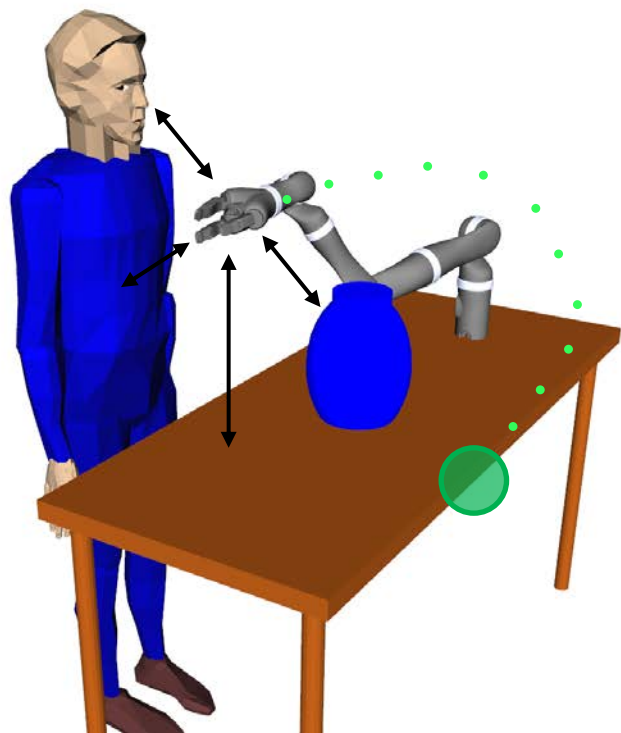
$$b'(\theta) \propto b(\theta)P(\tilde{\theta}|\theta)$$

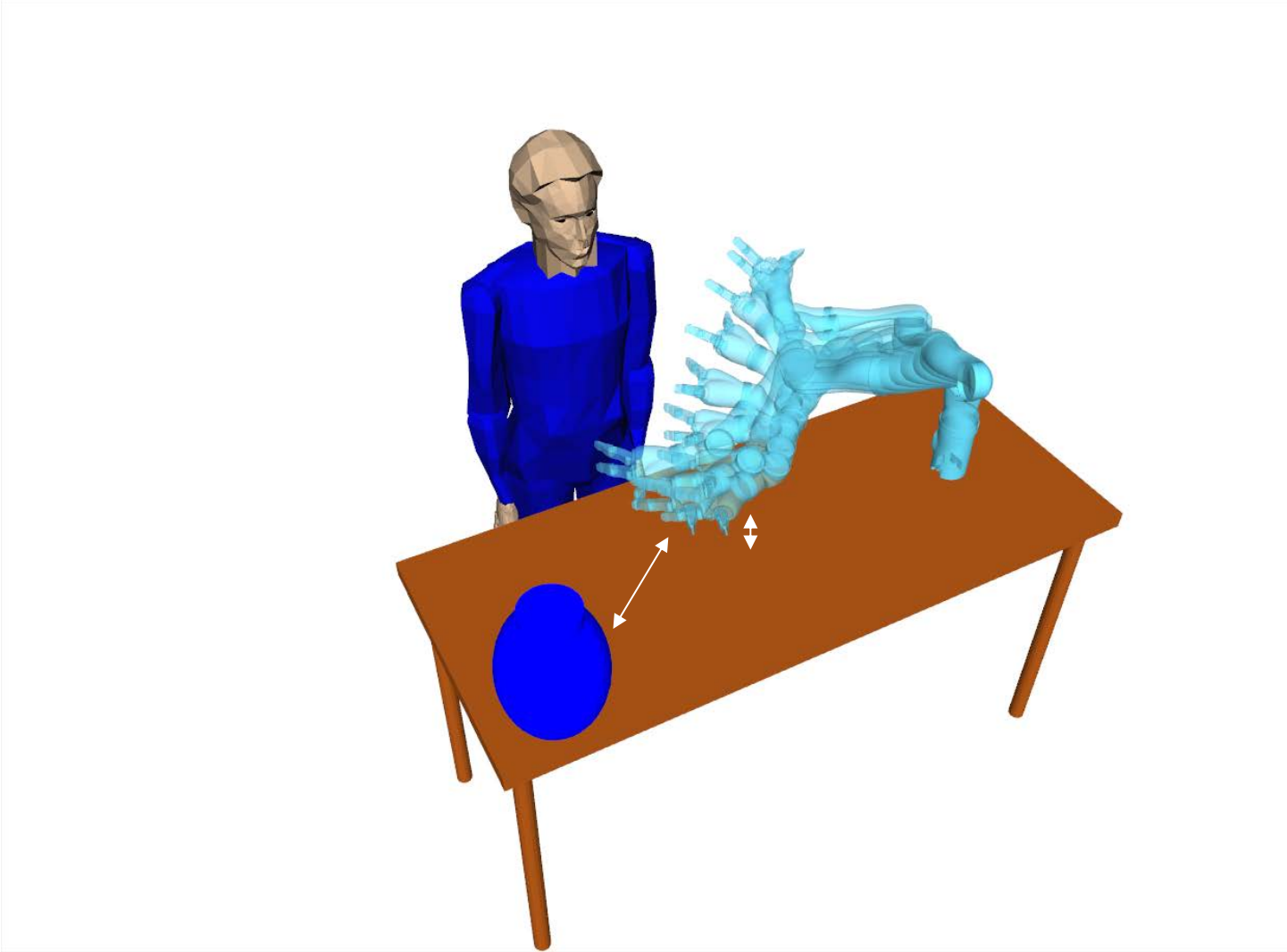


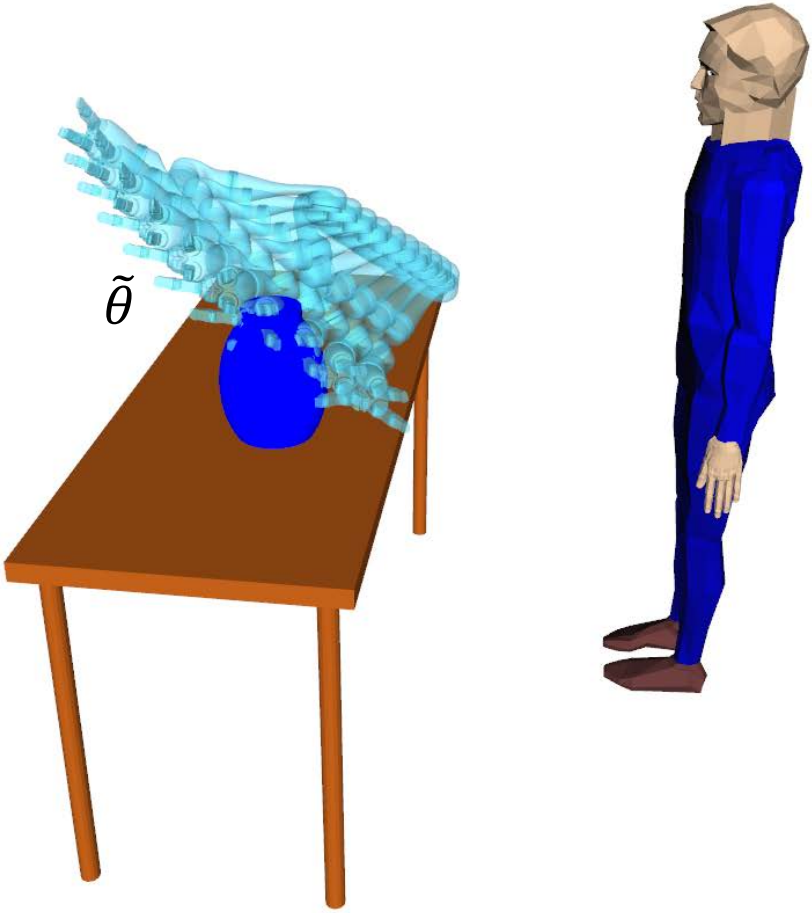
plan in expectation

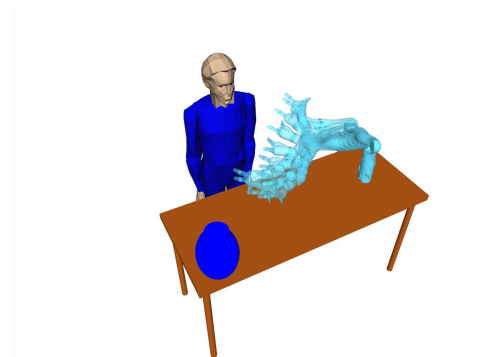
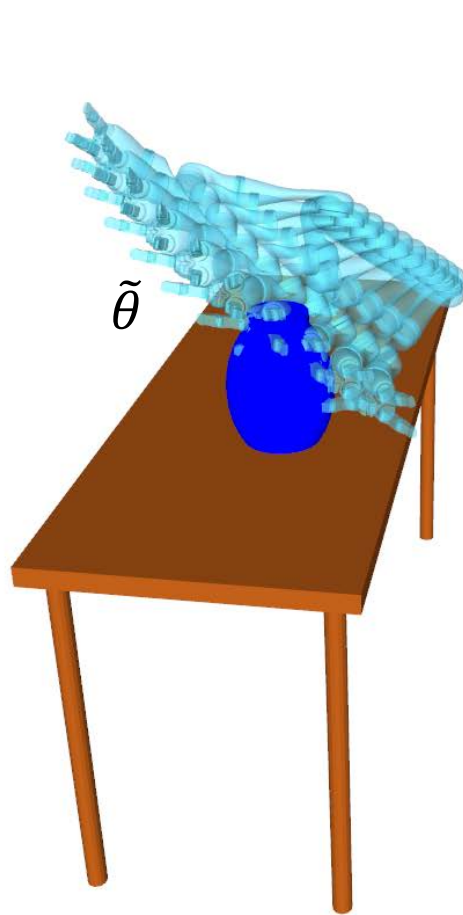
$$\max_{\xi} \mathbb{E}[R_\theta(\xi; M_{test}) | \theta \sim b'(\theta)]$$



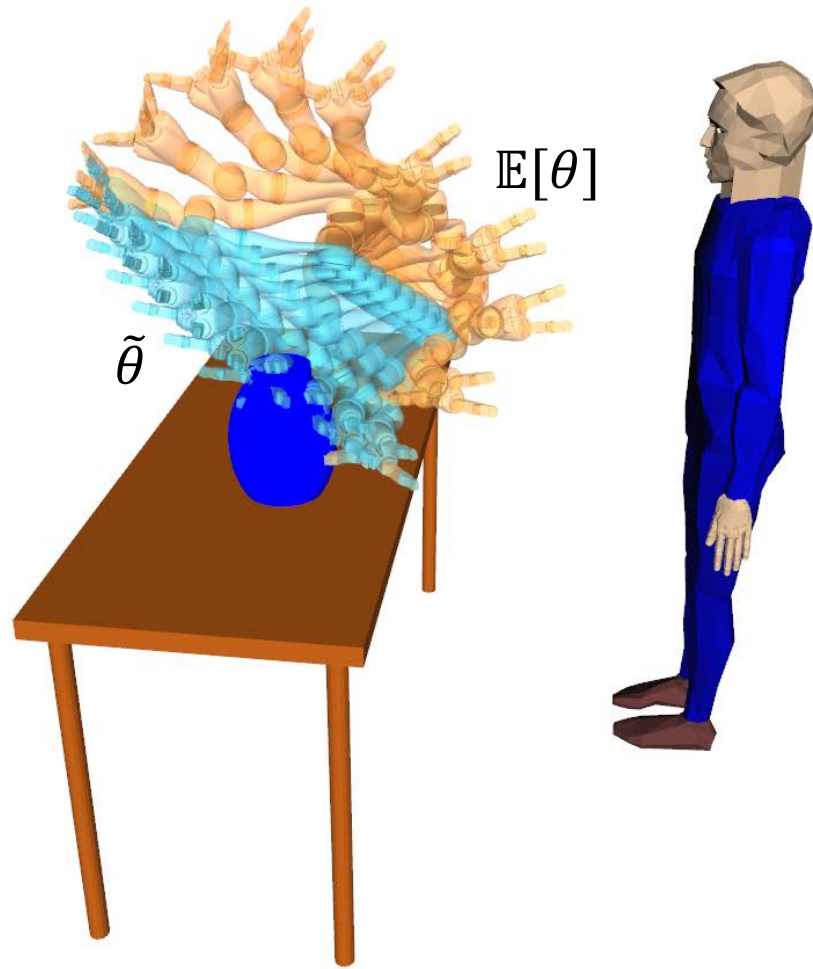








$$b'(\theta) \propto P(\tilde{\theta}|\theta, M_{train})b(\theta)$$



Easier, faster, lower regret

13.57.229.121:22362/exp?hitId=debugG458JS&assignmentId=debug3X6P6L&workerId=debug4RIDDM&mode=debug

Independent

Far from Torso Less Important More Important

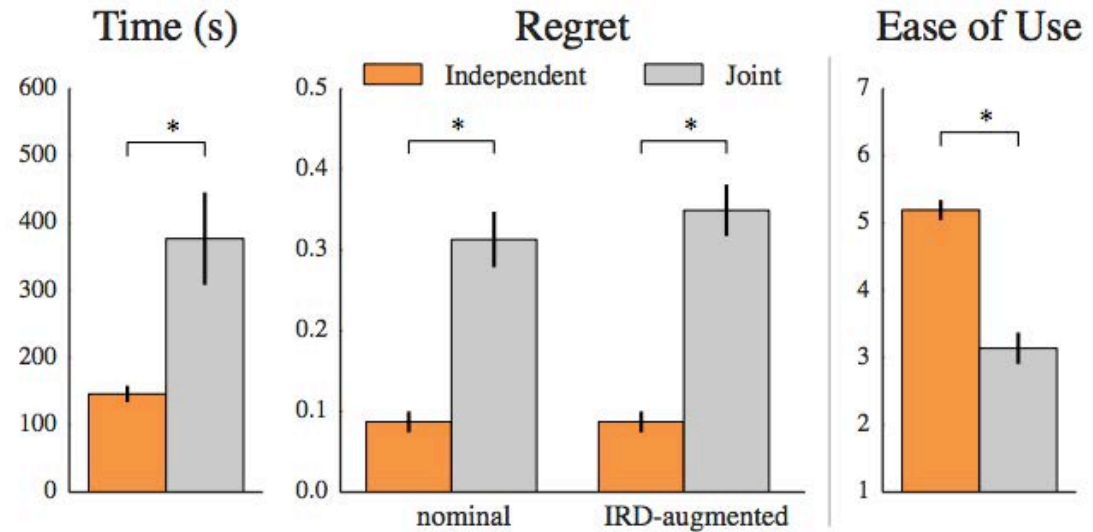
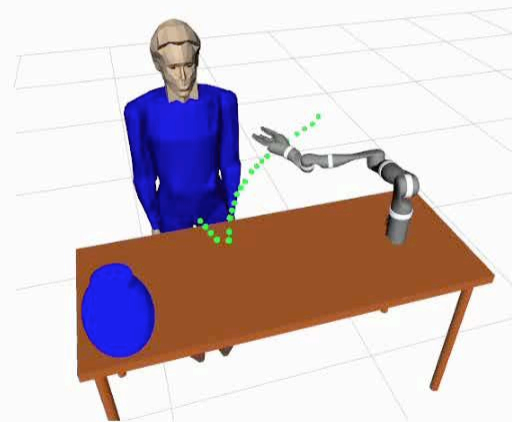
Far from Head Less Important More Important

Far from Vase Less Important More Important

Close to Table Less Important More Important

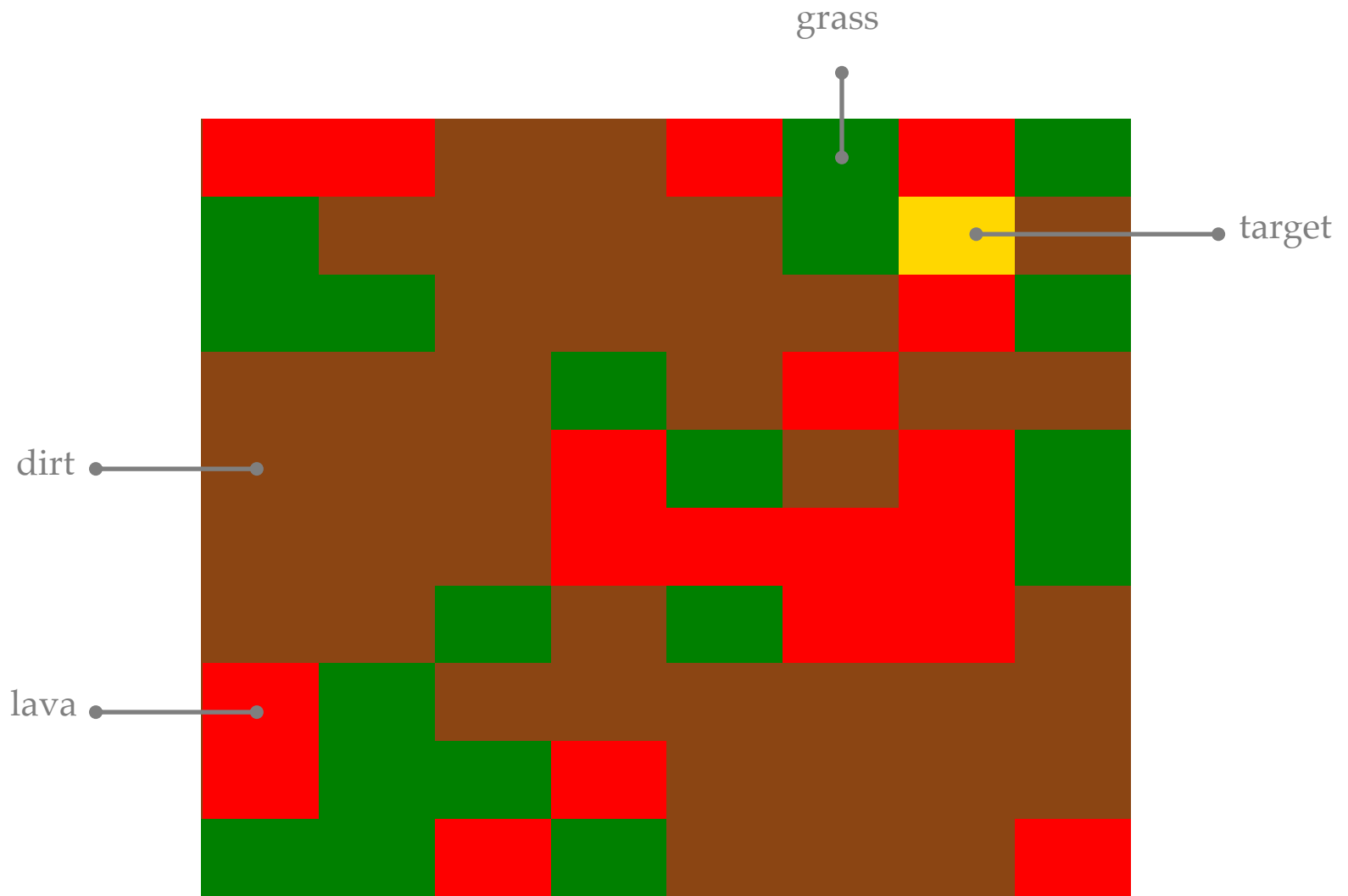
[Recompute Trajectory](#) [Next](#)

This is the independent phase, where you design a desired trajectory separately for 3 environments. The trajectory you want leads to the path in green. The current behavior is shown in the animated images. When you change the slider values, press **Recompute Trajectory** to show the robot's new path. When you have succeeded in specifying the desired behavior, the trajectory will turn green. Try to specify the correct behavior quickly, and in as few recomputations as possible.

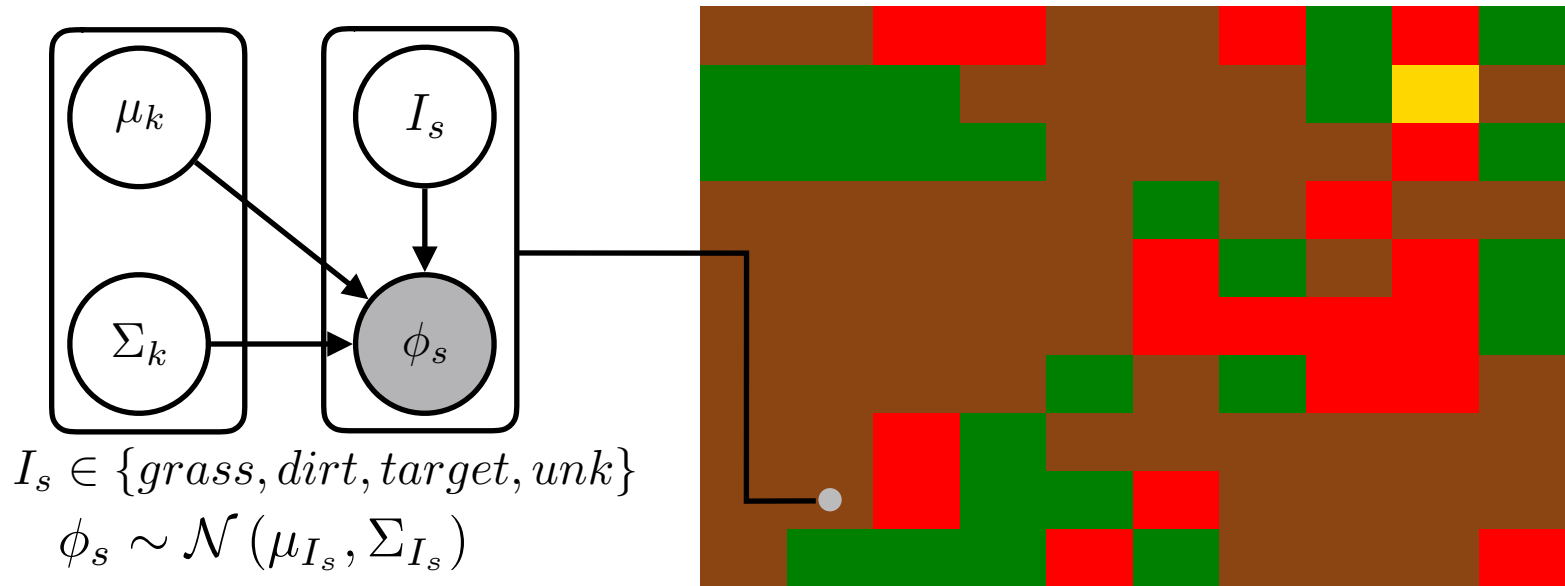


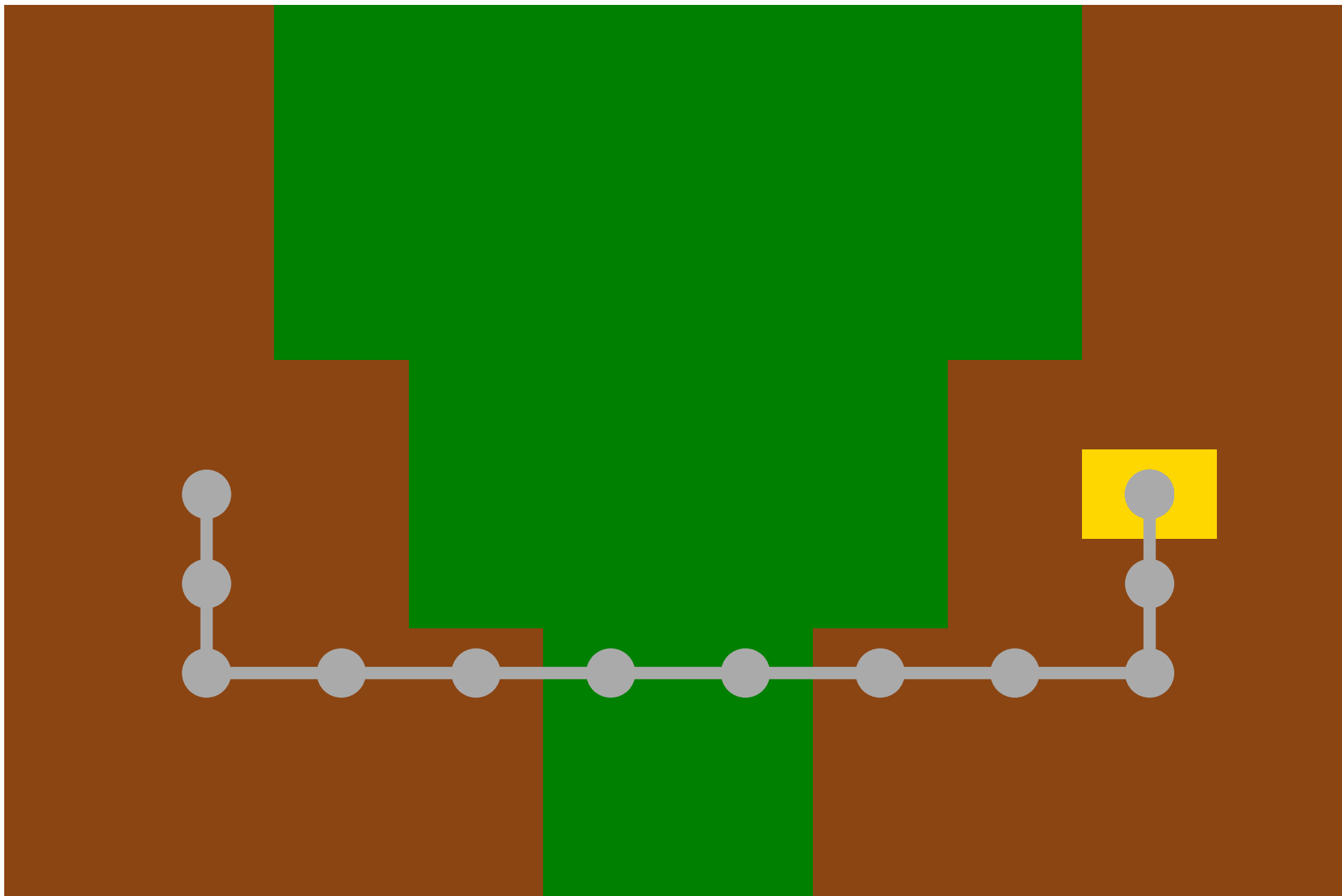
Limitations / ongoing work ..

What if we don't know the important features?!

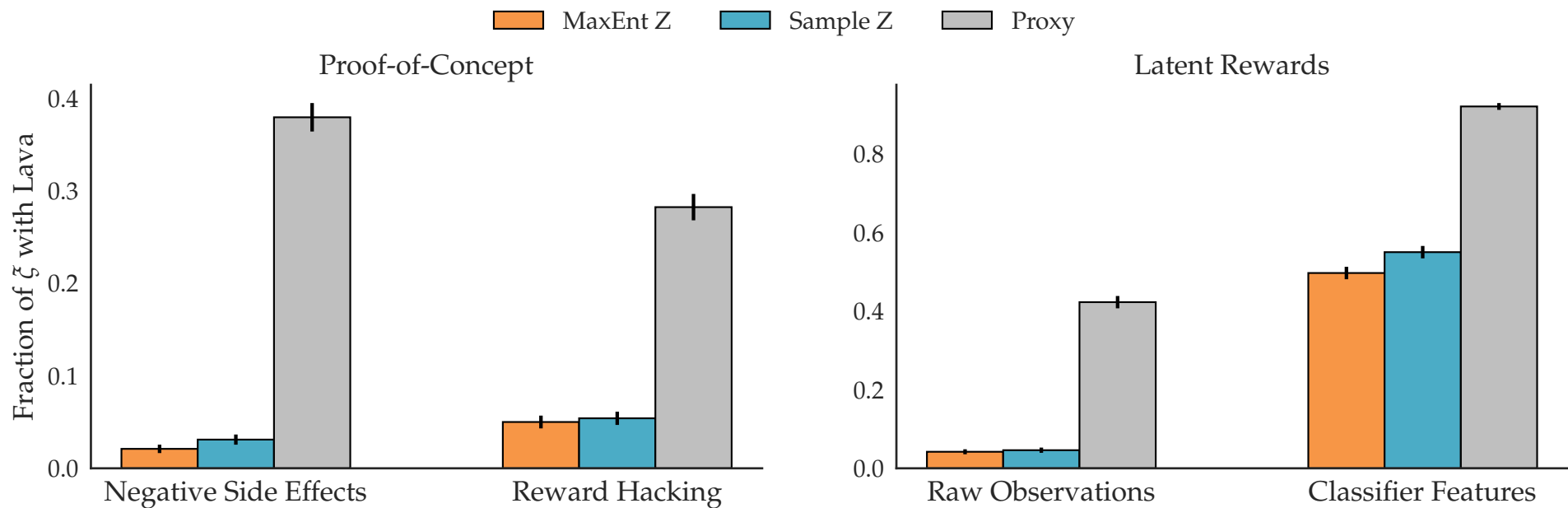


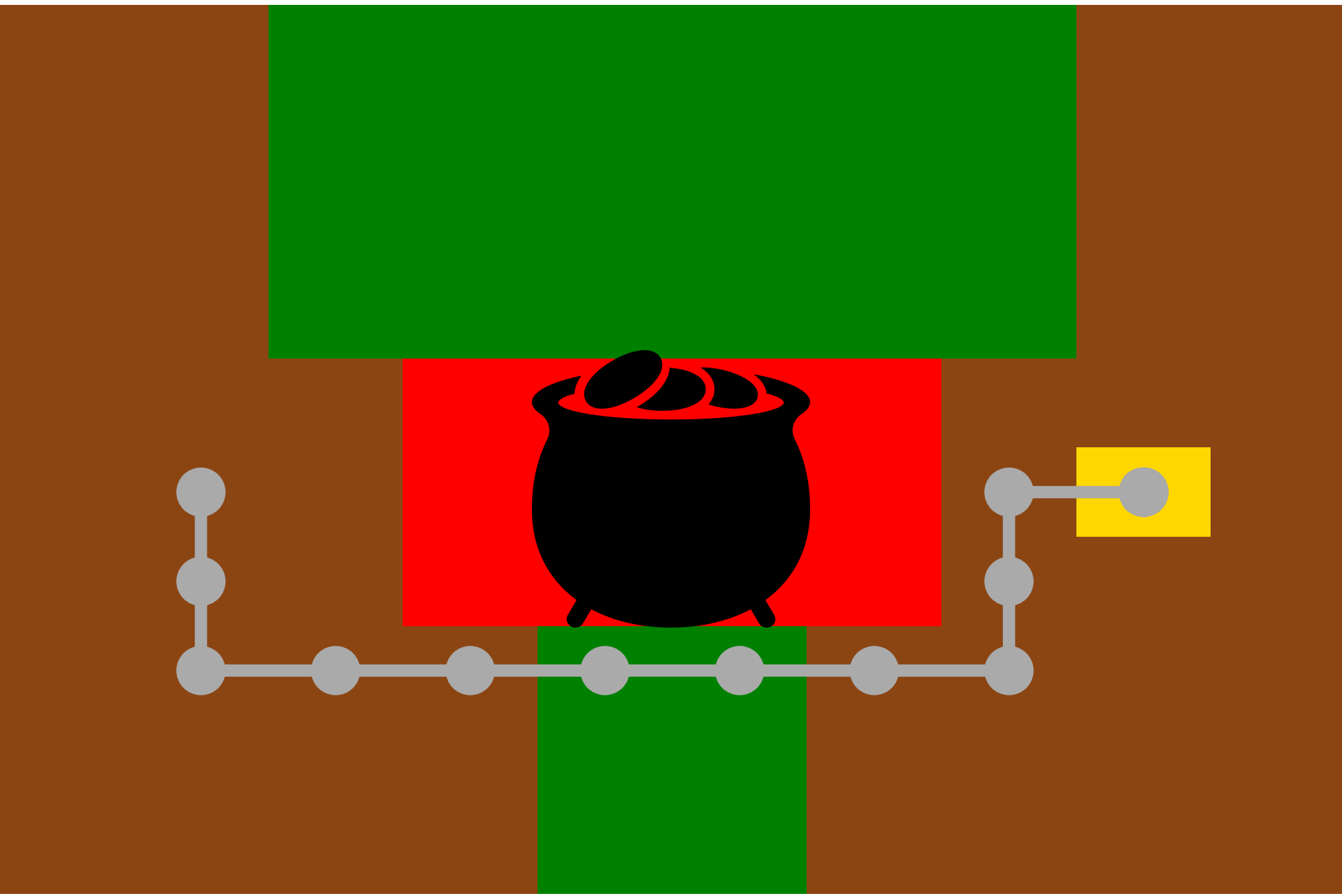
Inference from raw observations, no direct indicators...



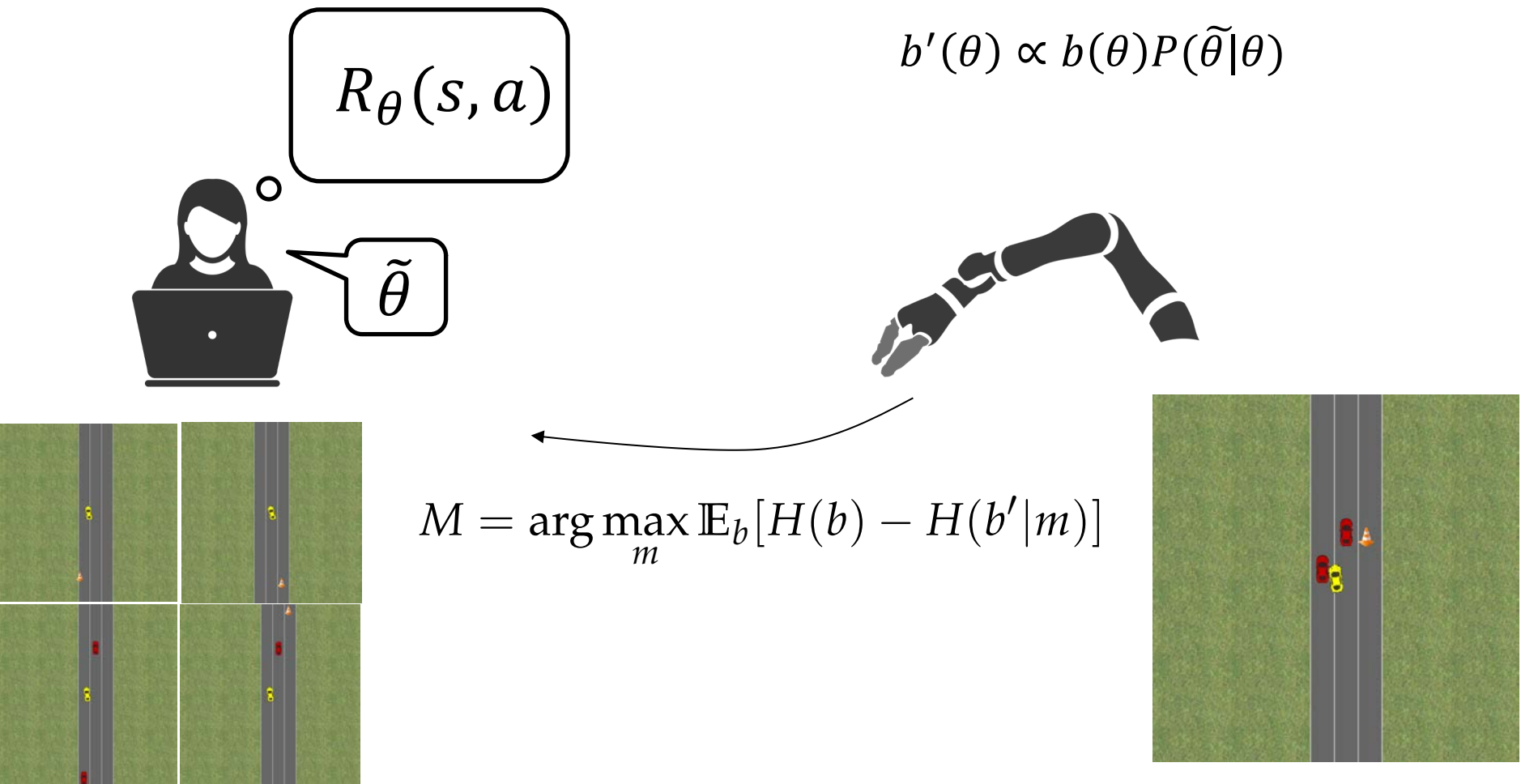


The agent can avoid unintended consequences, even when the features that matter are latent!

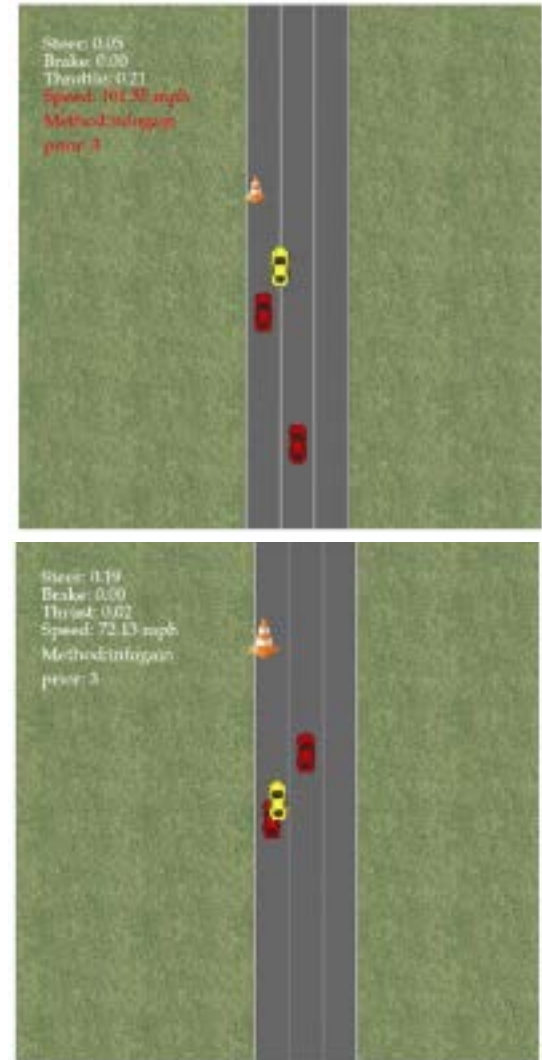
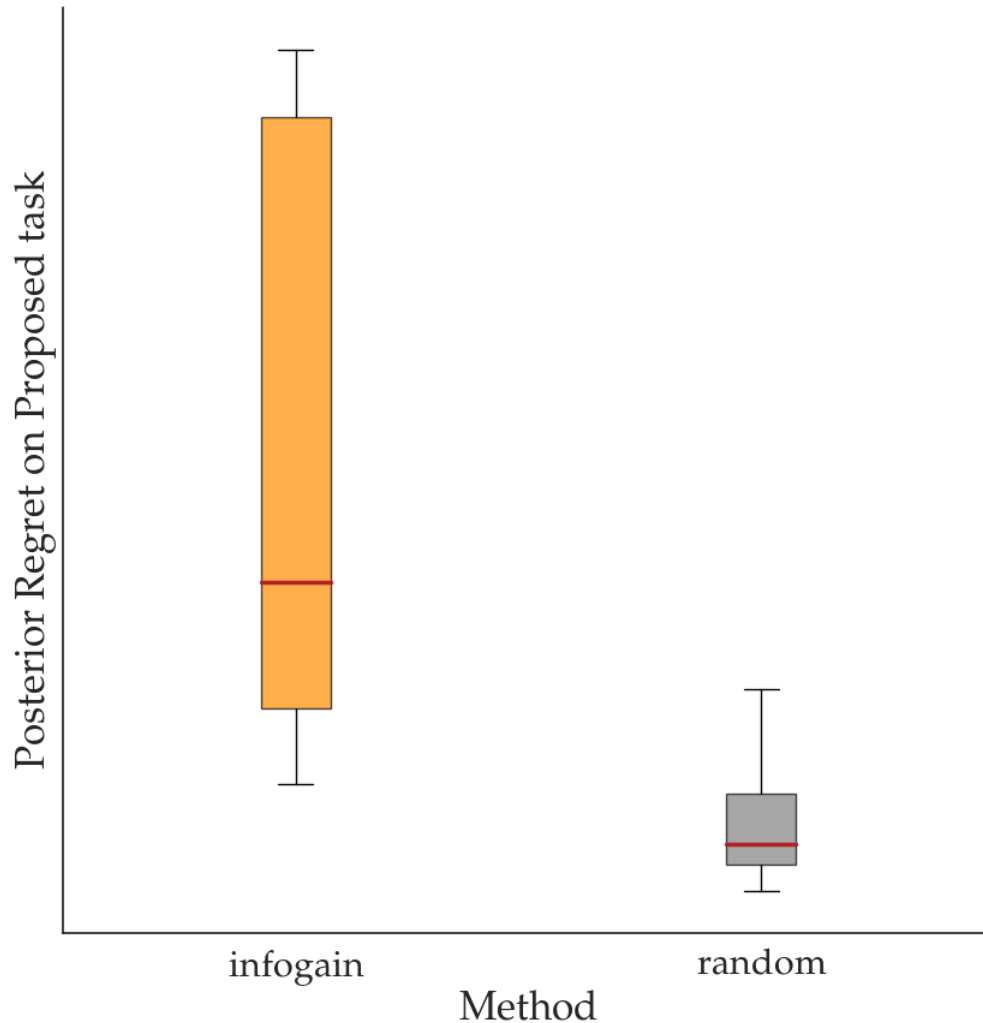




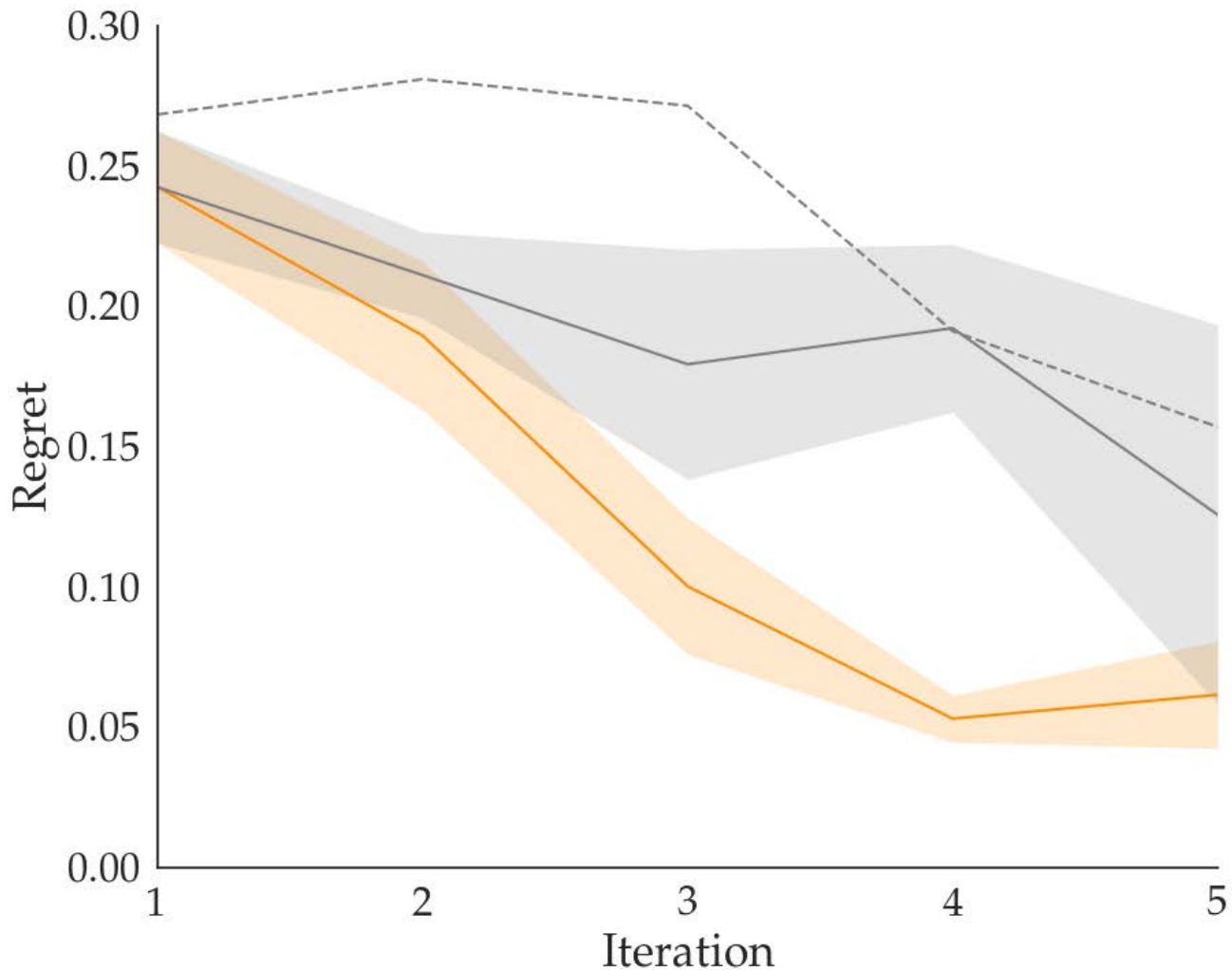
Leverage the posterior to identify edge cases



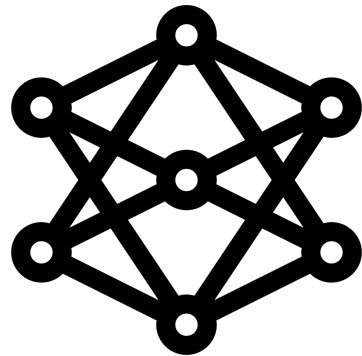
This finds edge-case environments that break the current reward function.



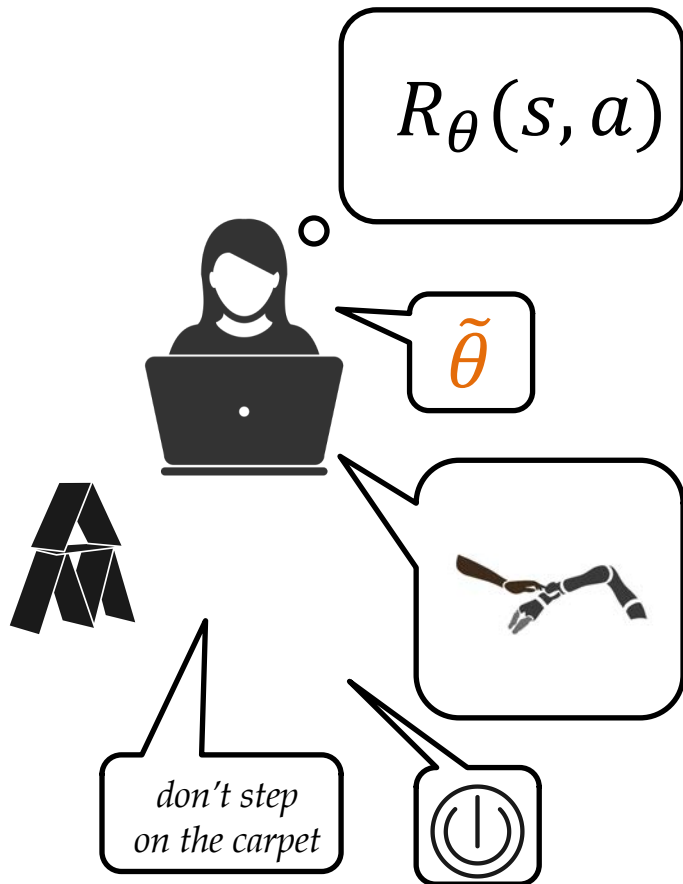
By exposing the designer to these edge cases, regret on held-out environments goes down quickly.



$$\theta^T \phi(\xi)$$



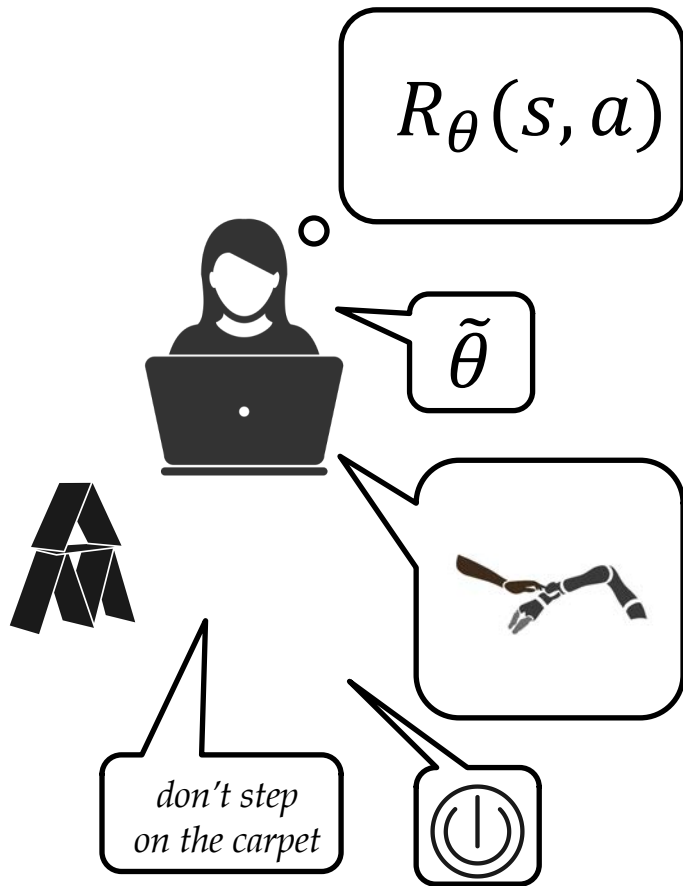
Specified rewards are evidence about the reward.



observation
(human) model

$$b'(\theta) \propto b(\theta)P(\tilde{\theta} | \theta)$$

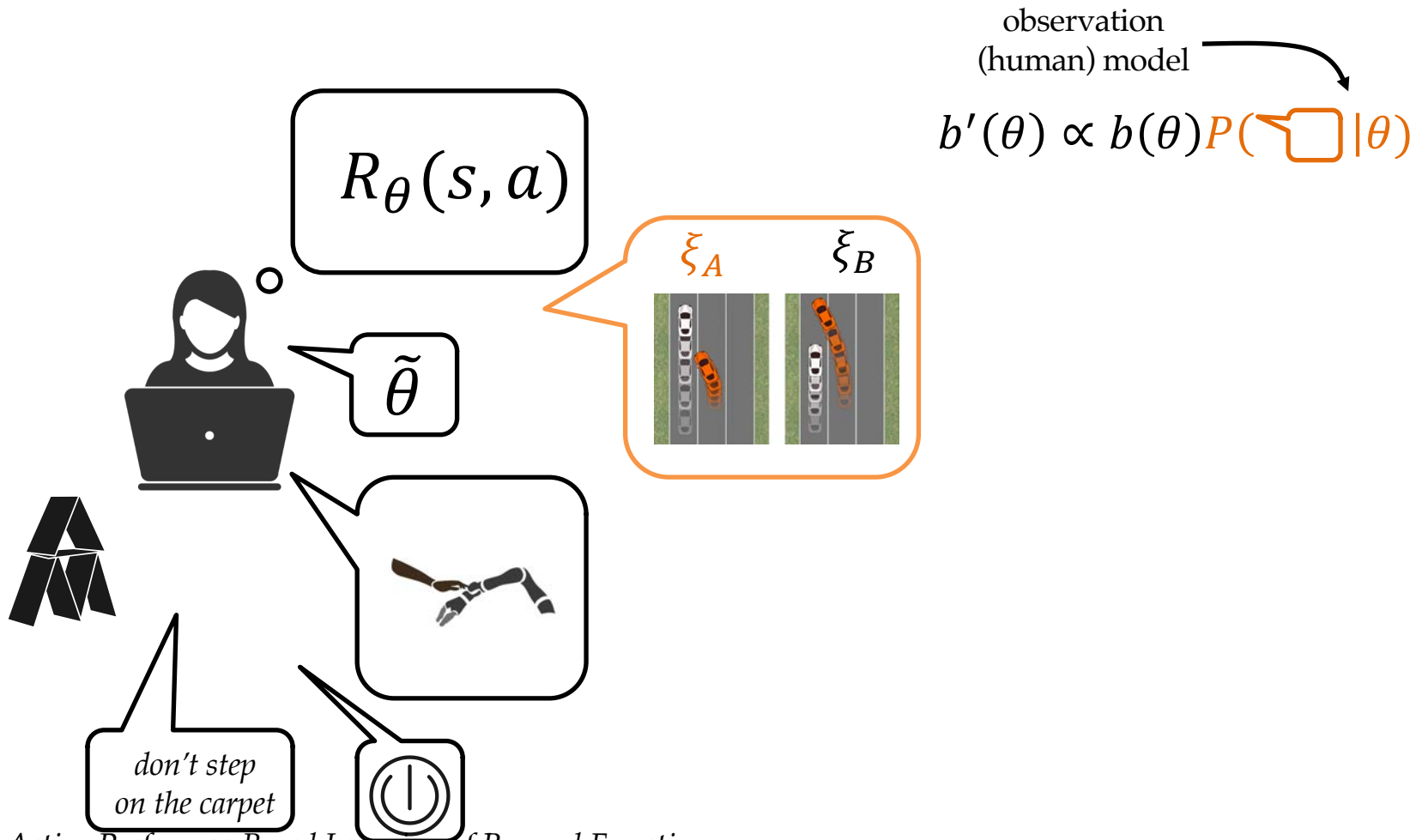
What is a human model that can be used to make sense of all these types of human feedback?



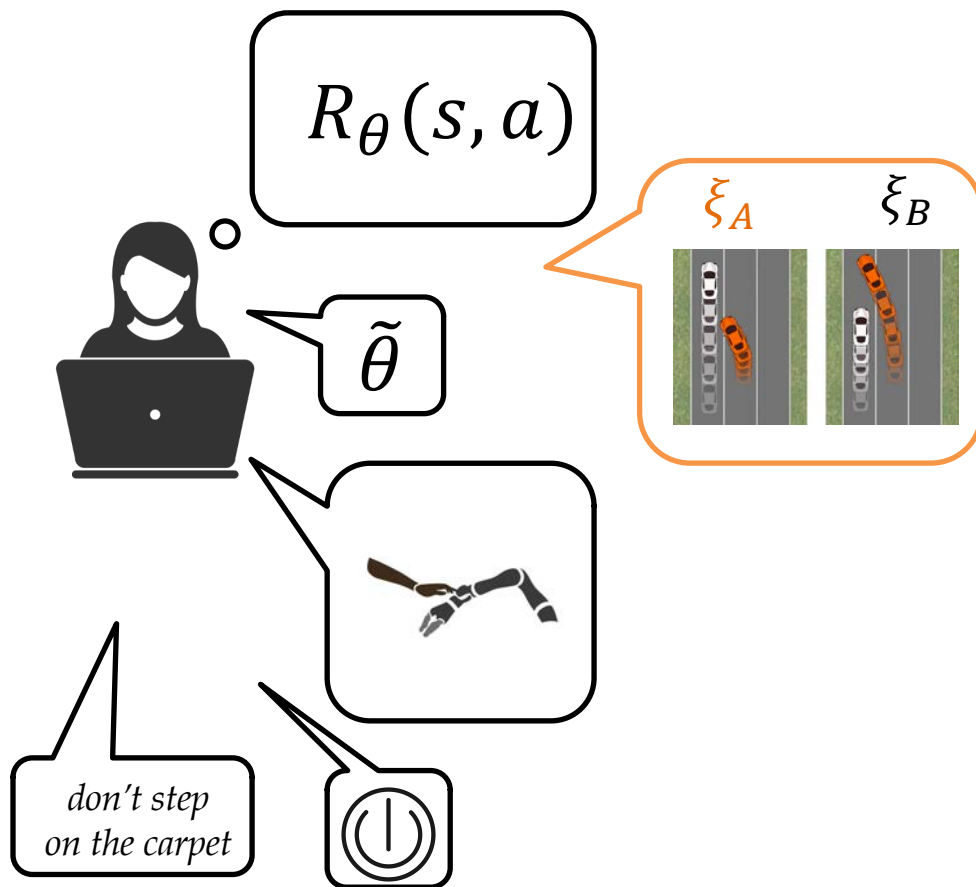
observation
(human) model

$$b'(\theta) \propto b(\theta) P(\text{speech bubble} | \theta)$$

We know what to do for comparisons.



We know what to do for comparisons:
model feedback as a reward-rational choice.



observation
(human) model

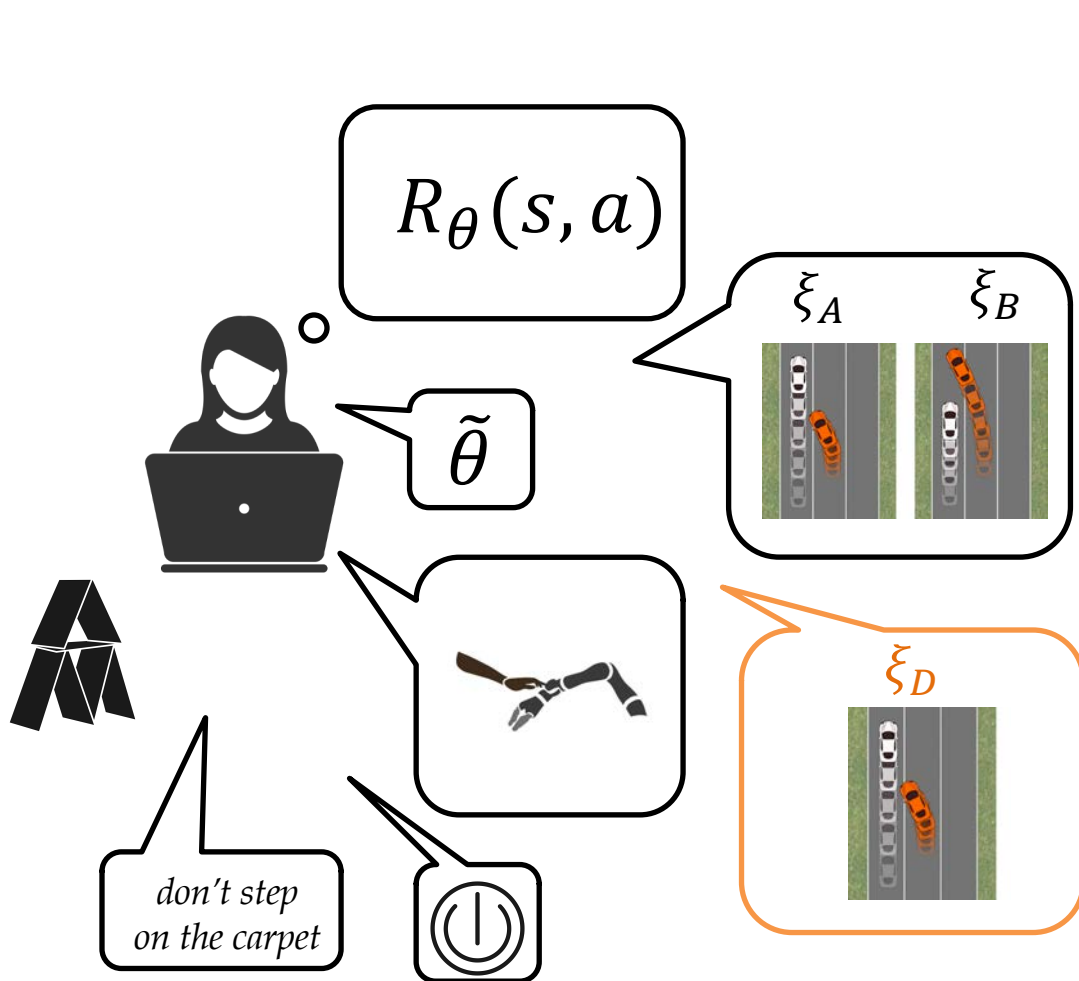
$$b'(\theta) \propto b(\theta) P(\text{choice} | \theta)$$

choices: $\{\xi_A, \xi_B\}$

choose based on reward: $R_{\theta}(\xi_A) \text{ vs } R_{\theta}(\xi_B)$

$$P(\xi_A | \theta) = \frac{e^{R_{\theta}(\xi_A)}}{e^{R_{\theta}(\xi_A)} + e^{R_{\theta}(\xi_B)}}$$

We know what to do for demonstrations:
 model the demo as a reward-rational **implicit** choice.



observation
 (human) model

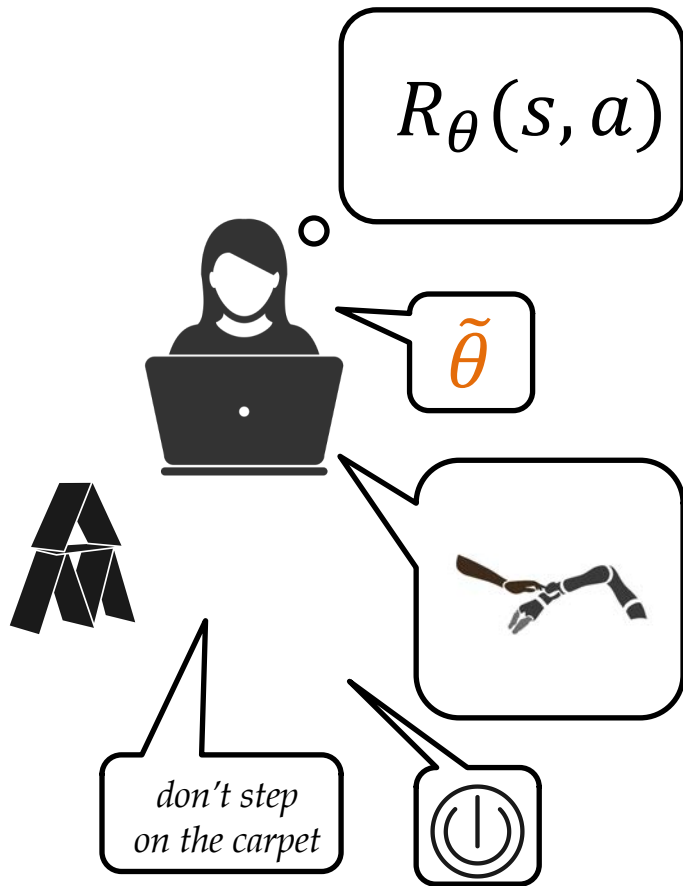
$$b'(\theta) \propto b(\theta) P(\text{choice} | \theta)$$

choices: $\{\xi_i\}$

choose based on reward: $R_\theta(\xi_D) \text{ vs } R_\theta(\xi) \forall \xi$

$$P(\xi_D | \theta) = \frac{e^{R_\theta(\xi_D)}}{\sum_{\xi} e^{R_\theta(\xi)}}$$

We know what to do for specified rewards

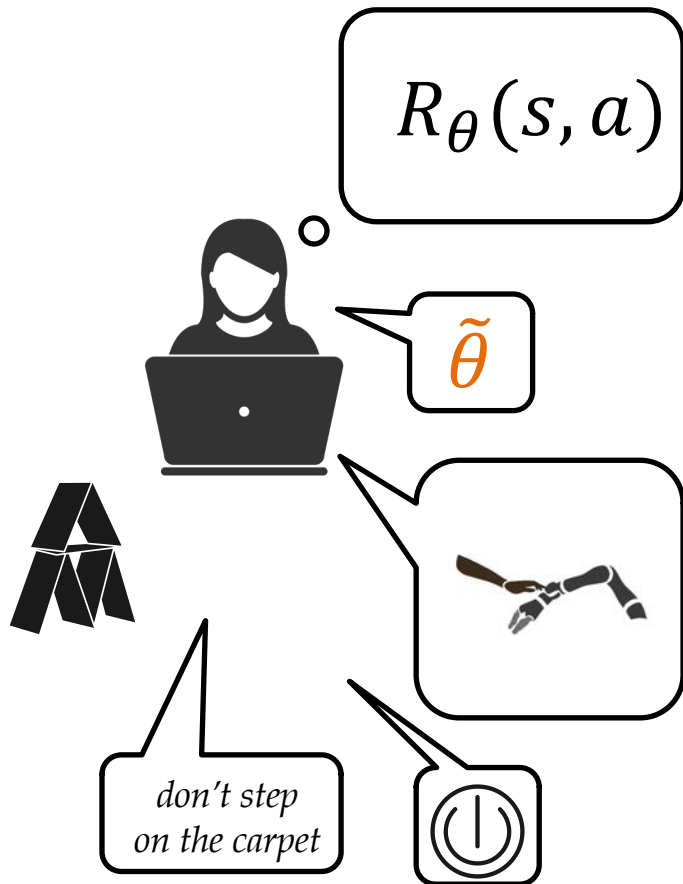


observation
(human) model

$$b'(\theta) \propto b(\theta) P(\text{observation} | \theta)$$

$$P(\tilde{\theta} | \theta) = \frac{e^{\mathbb{E}[R_{\theta}(\xi) | \xi \sim P(\xi | \tilde{\theta}, M_{devel})]}}{\sum_{\bar{\theta}} e^{\mathbb{E}[R_{\theta}(\xi) | \xi \sim P(\xi | \bar{\theta}, M_{devel})]}}$$

We know what to do for specified rewards:
model them as a reward-rational implicit choice.



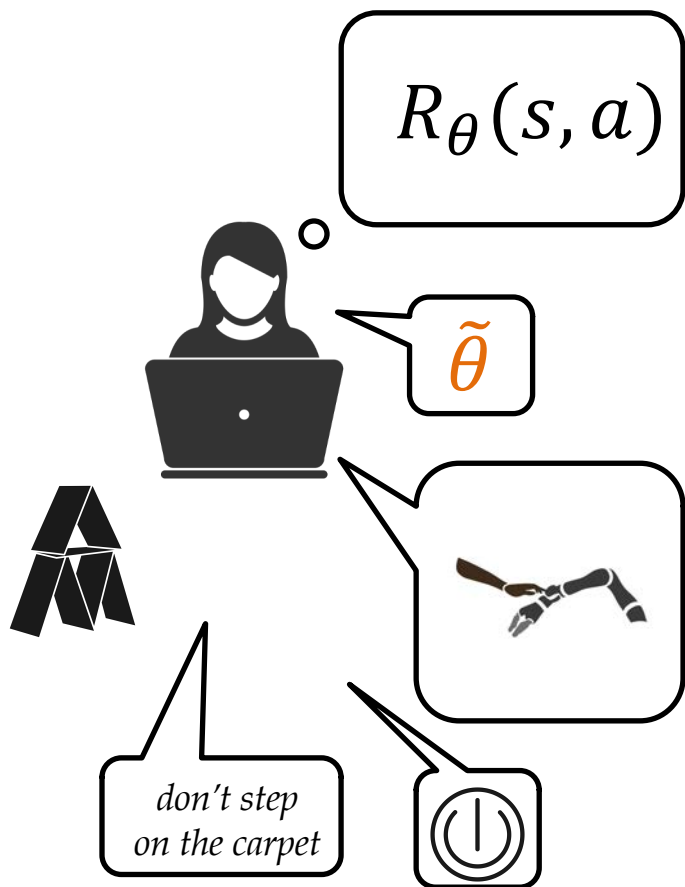
observation
(human) model

$$b'(\theta) \propto b(\theta) P(\text{observation} | \theta)$$

choices: $\{\theta_i\}$

$$P(\tilde{\theta} | \theta) = \frac{e^{\mathbb{E}[R_{\theta}(\xi) | \xi \sim P(\xi | \tilde{\theta}, M_{devel})]}}{\sum_{\bar{\theta}} e^{\mathbb{E}[R_{\theta}(\xi) | \xi \sim P(\xi | \bar{\theta}, M_{devel})]}}$$

We know what to do for specified rewards:
model them as a reward-rational implicit choice.



observation
(human) model

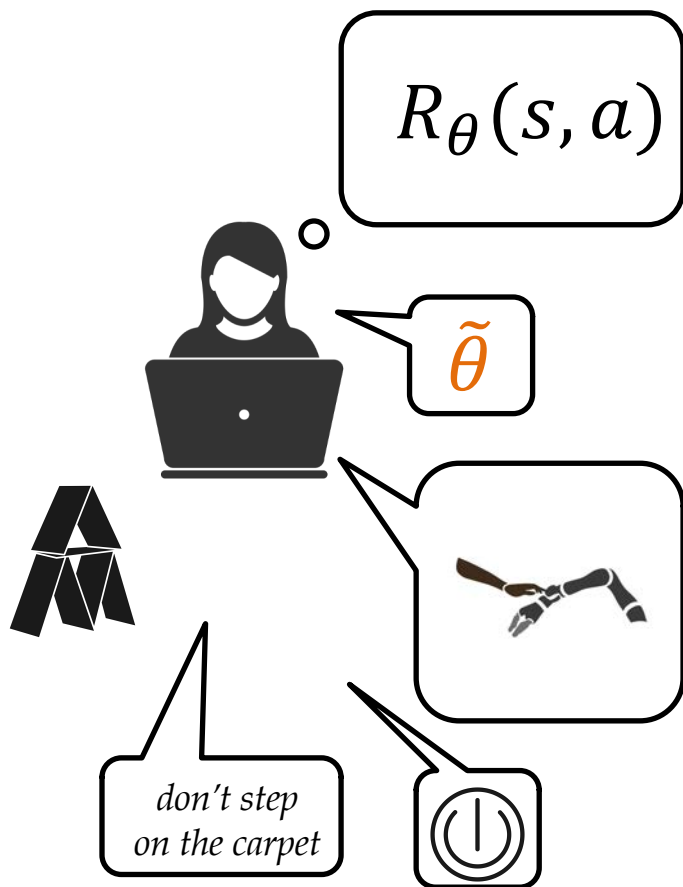
$$b'(\theta) \propto b(\theta) P(\text{observation} | \theta)$$

choices: $\{\theta_i\}$

choose based
on reward: $R_{\theta}(\tilde{\theta})?!$

$$P(\tilde{\theta} | \theta) = \frac{e^{\mathbb{E}[R_{\theta}(\xi) | \xi \sim P(\xi | \tilde{\theta}, M_{devel})]}}{\sum_{\bar{\theta}} e^{\mathbb{E}[R_{\theta}(\xi) | \xi \sim P(\xi | \bar{\theta}, M_{devel})]}}$$

We know what to do for specified rewards:
model them as a reward-rational implicit choice.



observation
(human) model

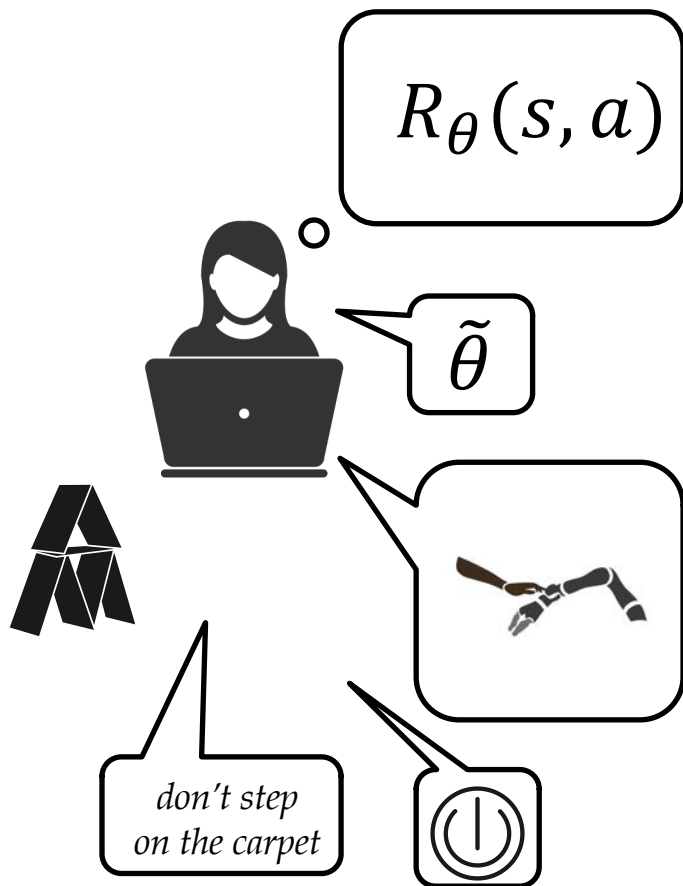
$$b'(\theta) \propto b(\theta) P(\text{observation} | \theta)$$

choices: $\{\theta_i\}$

choose based
on reward: $\mathbb{E}[R_\theta(\xi) | \xi \sim P(\xi | \tilde{\theta}, M_{devel})]$
vs
 $\mathbb{E}[R_\theta(\xi) | \xi \sim P(\xi | \bar{\theta}, M_{devel})] \forall \bar{\theta}$

$$P(\tilde{\theta} | \theta) = \frac{e^{\mathbb{E}[R_\theta(\xi) | \xi \sim P(\xi | \tilde{\theta}, M_{devel})]}}{\sum_{\bar{\theta}} e^{\mathbb{E}[R_\theta(\xi) | \xi \sim P(\xi | \bar{\theta}, M_{devel})]}}$$

Reward-rational (implicit) choices



observation
(human) model

$$b'(\theta) \propto b(\theta) P(\text{choice} | \theta)$$

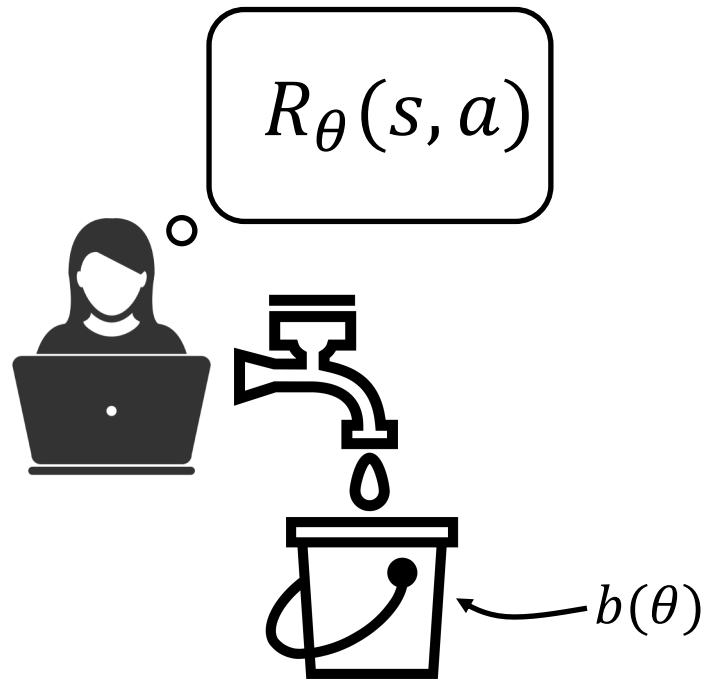
choices: $\{c\}$

choose based on reward: $\mathbb{E}[R_\theta(\xi) | \xi \sim \psi(c^*)]$

vs

$\mathbb{E}[R_\theta(\xi) | \xi \sim \psi(c)] \forall c$

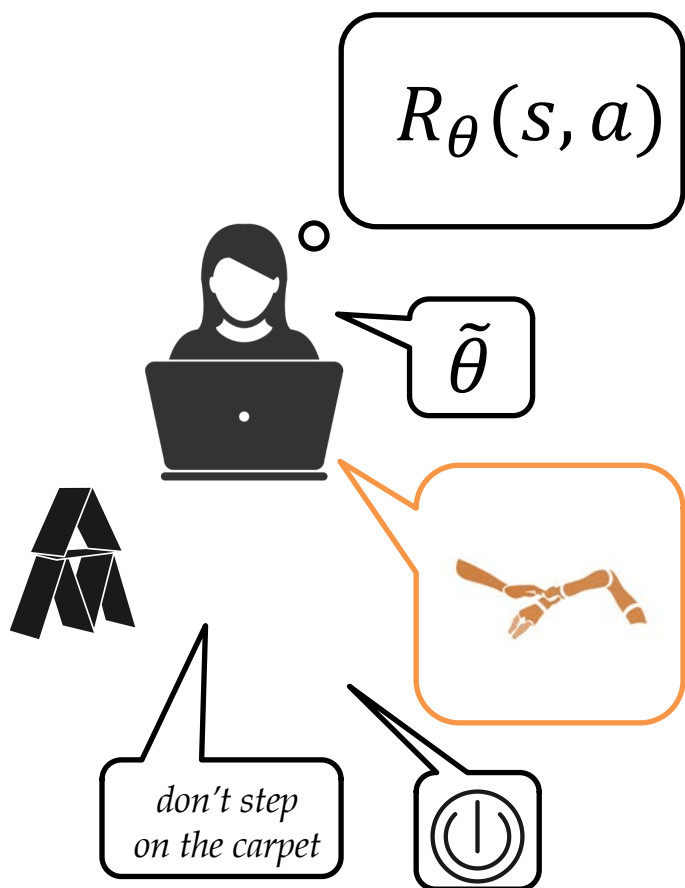
$$P(c^* | \theta) = \frac{e^{\mathbb{E}[R_\theta(\xi) | \xi \sim \psi(c^*)]}}{\sum e^{\mathbb{E}[R_\theta(\xi) | \xi \sim \psi(c)]}}$$



How should the robot extract the leaked information into an updated belief?

Key idea: Interpret any type of human feedback as a reward-rational implicit choice.

Human feedback as a reward-rational implicit choice.

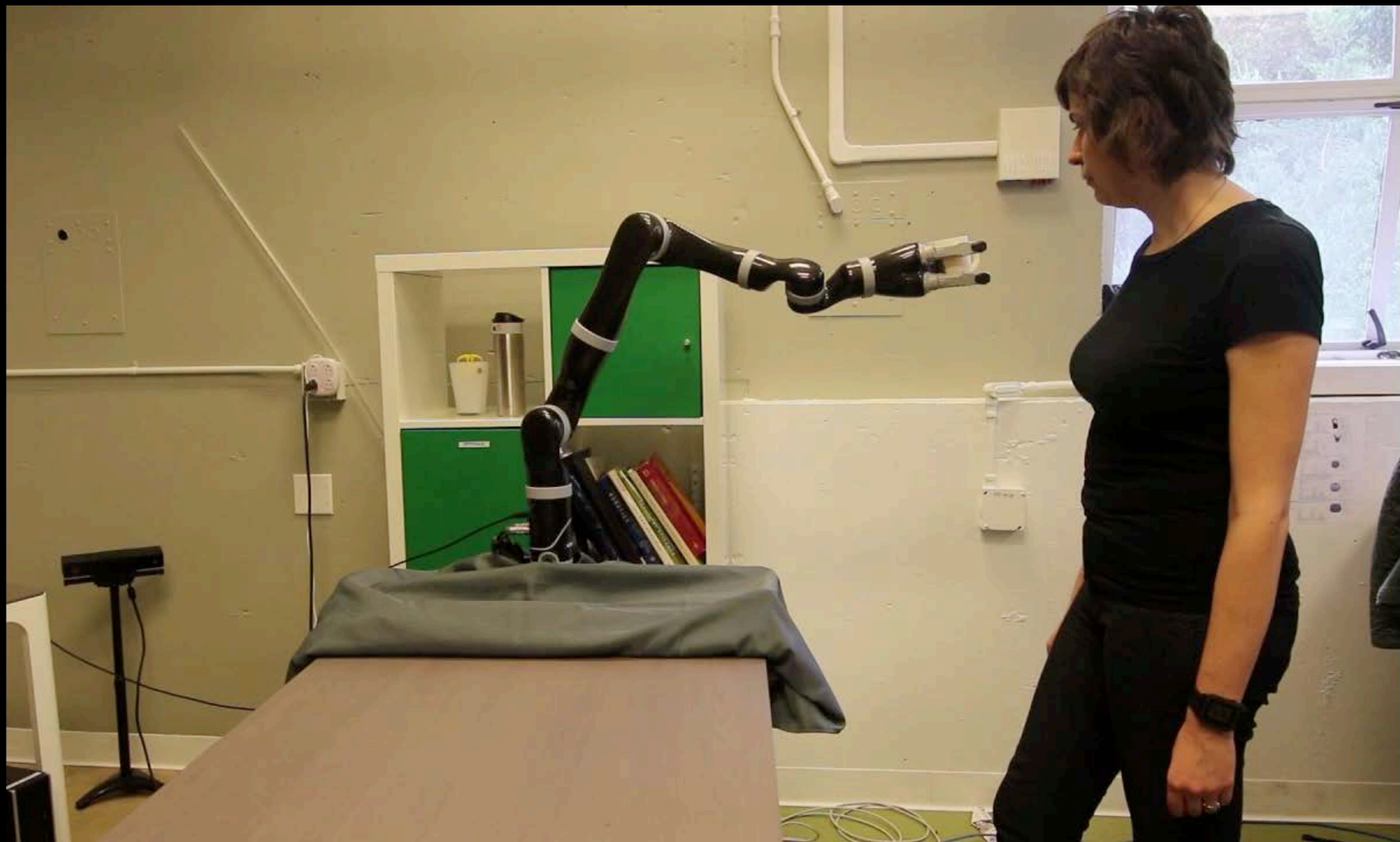


observation
(human) model

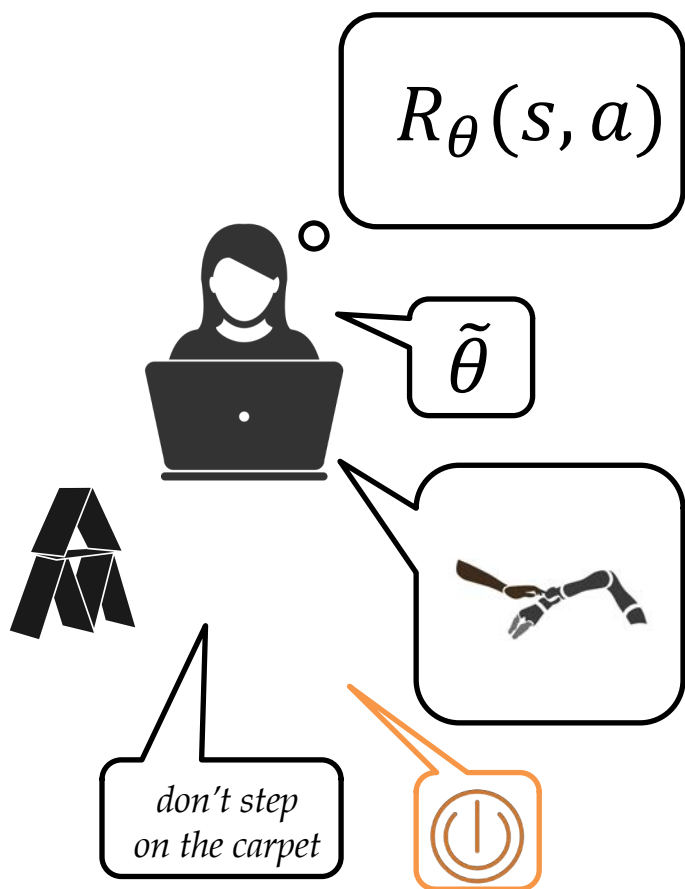
$$b'(\theta) \propto b(\theta) P(\text{speech bubble} | \theta)$$

choices: $\{\tau\}$ (external torques)

choose based on reward: $R_{\theta}(\xi(\xi_{original}, \tau))$
(deformed trajectories)



Human feedback as a reward-rational implicit choice.



observation
(human) model

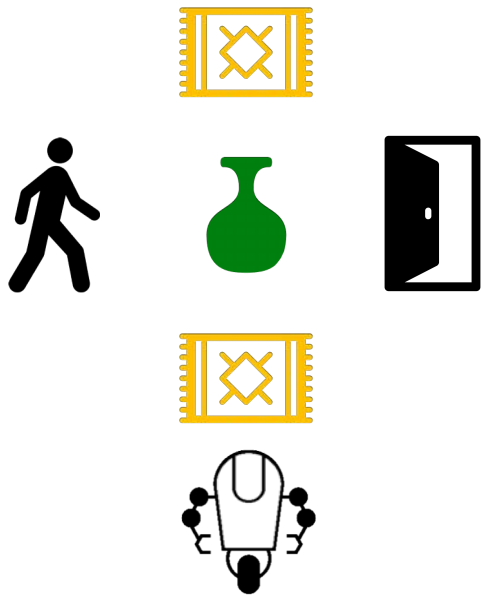
$$b'(\theta) \propto b(\theta) P(\text{power button} | \theta)$$

choices: $\{ \textit{press button}, \textit{do nothing} \}$

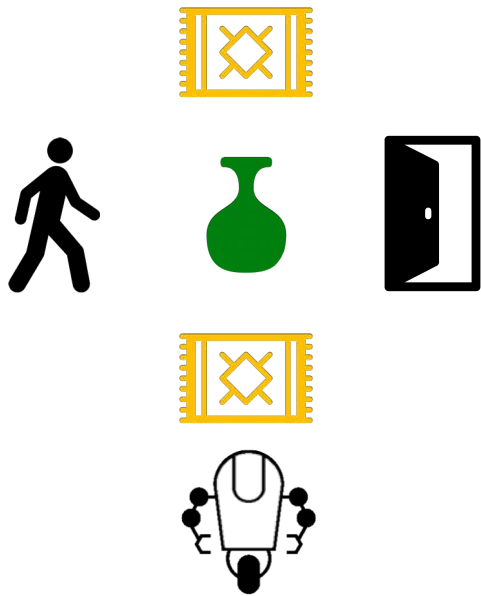
choose based
on reward: $R_{\theta}(\xi_{\textit{stopped}})$
vs
 $R_{\theta}(\xi_{\textit{planned}})$



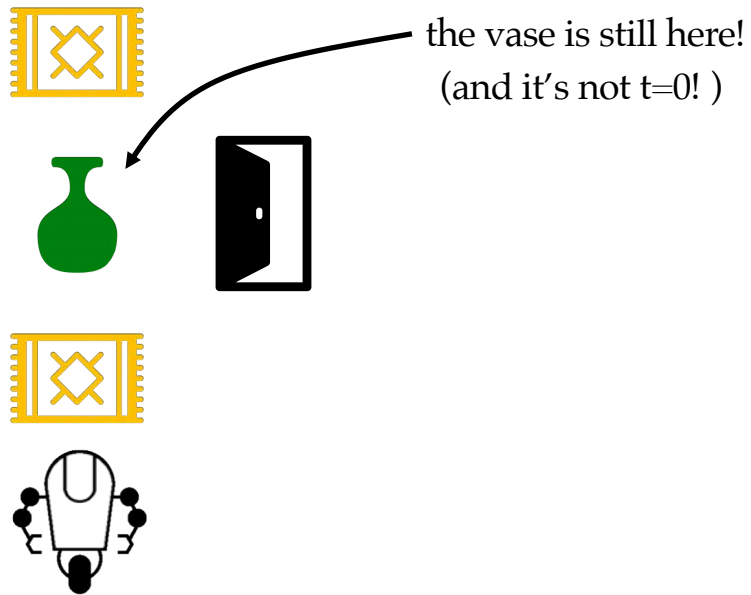
So far, we've talked about sources of information that look at human behavior:

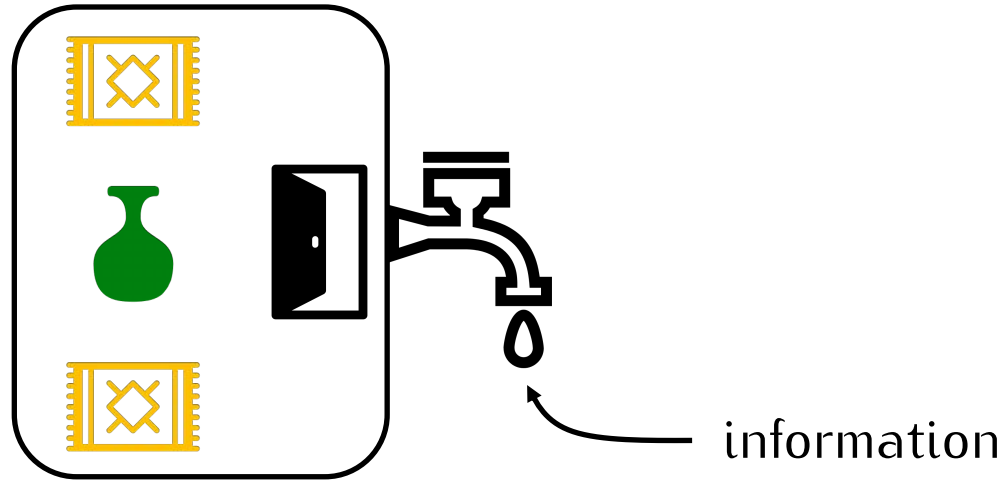


To know that you shouldn't break the vase, you need to see some behavior, e.g.:



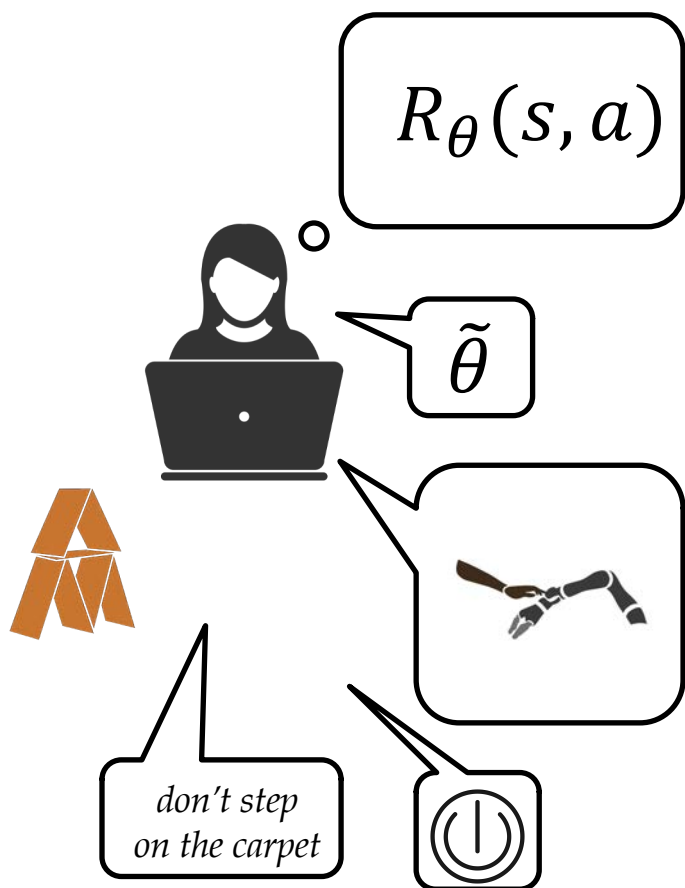
What if we don't see any behavior?





When the agent is deployed in an environment that the human has been acting in, the state of the environment has information about the human's intended reward.

The state of the environment as a reward-rational implicit choice.



observation
(human) model

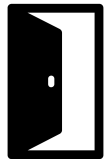
$$b'(\theta) \propto b(\theta) P(\text{observation} | \theta)$$

choices: $\{s_0\}$ (states)

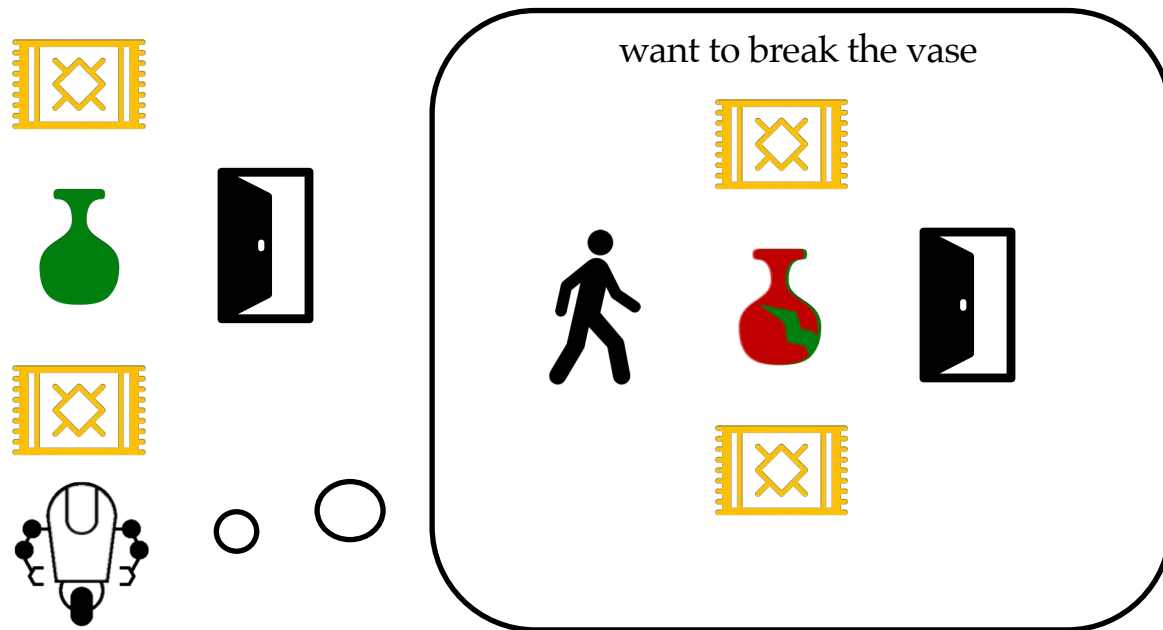
choose based on reward: $\mathbb{E}[R_{\theta}(\xi_{-T:0}) | \xi(0) = s_0]$

(trajectories that end at the observed state)

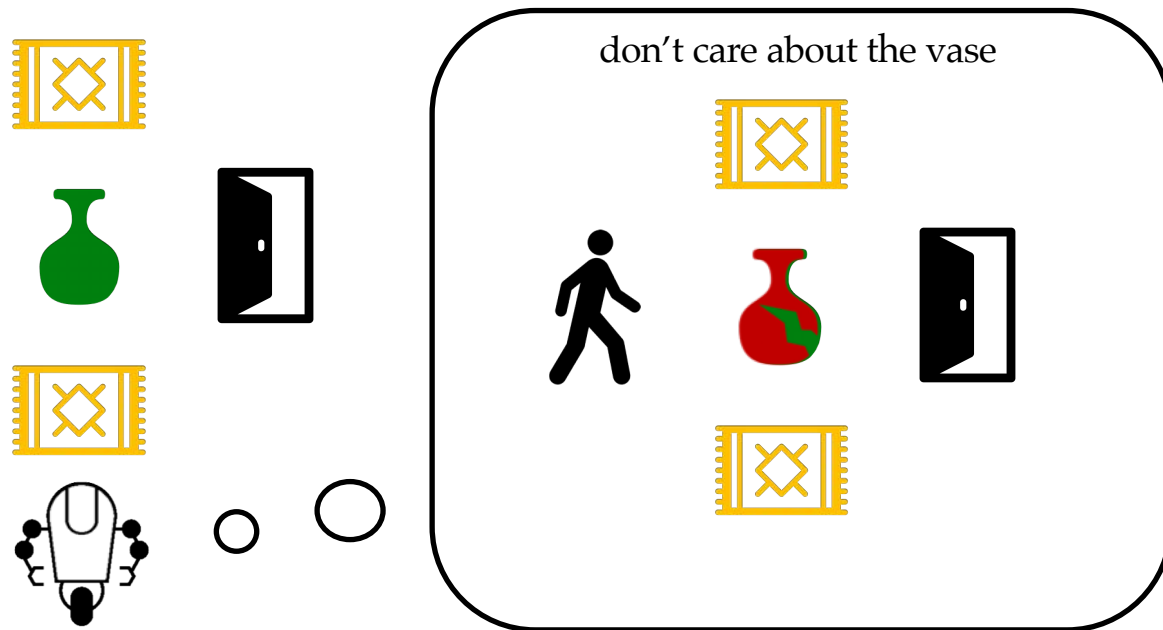
What reward function is the state consistent with?



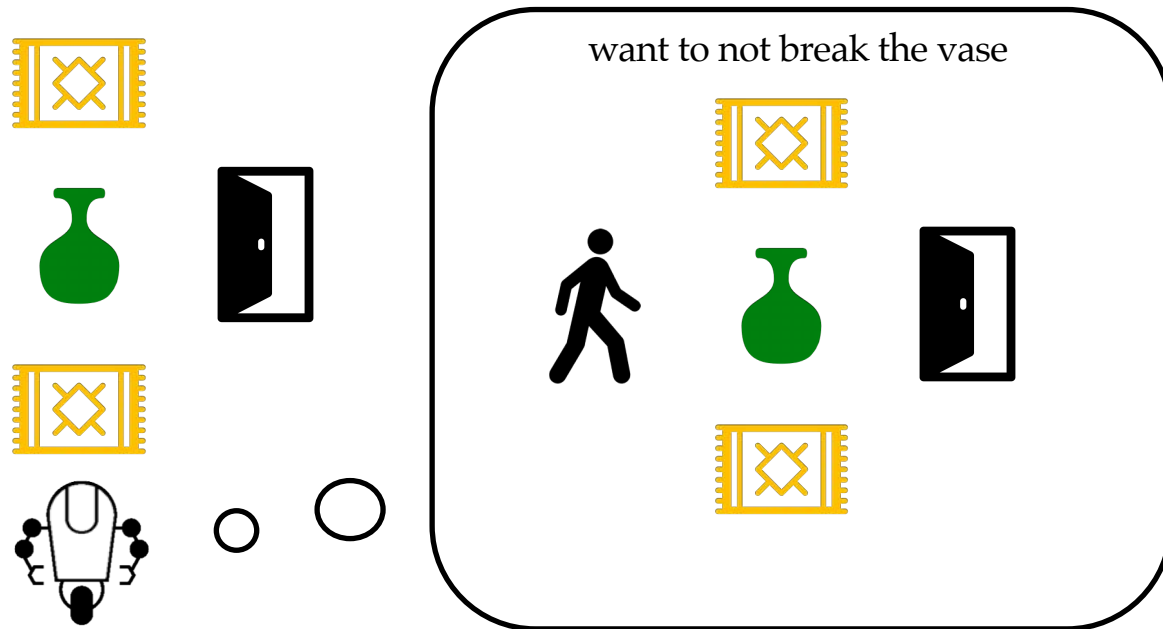
What reward function is the state consistent with?



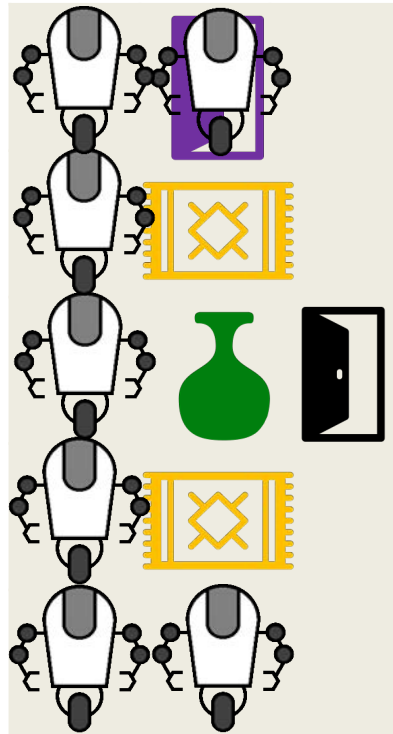
What reward function is the state consistent with?



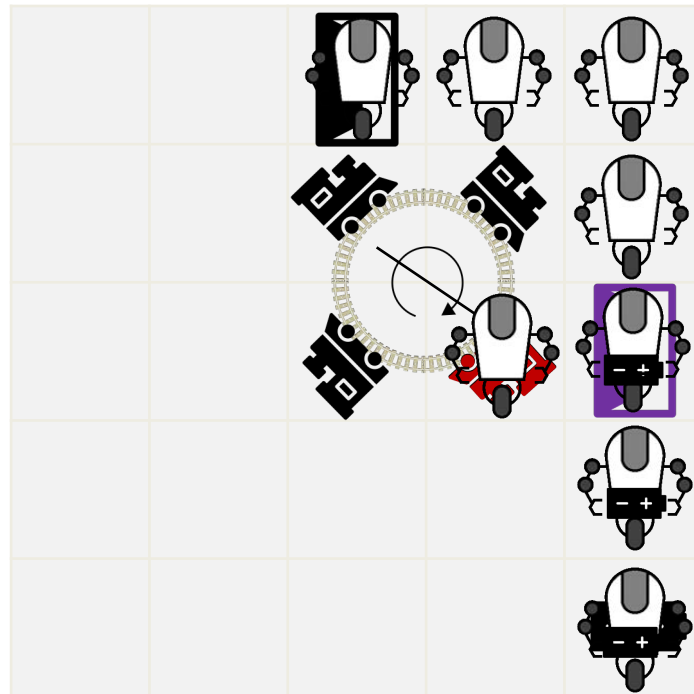
What reward function is the state consistent with?



Side effects: Room with vase



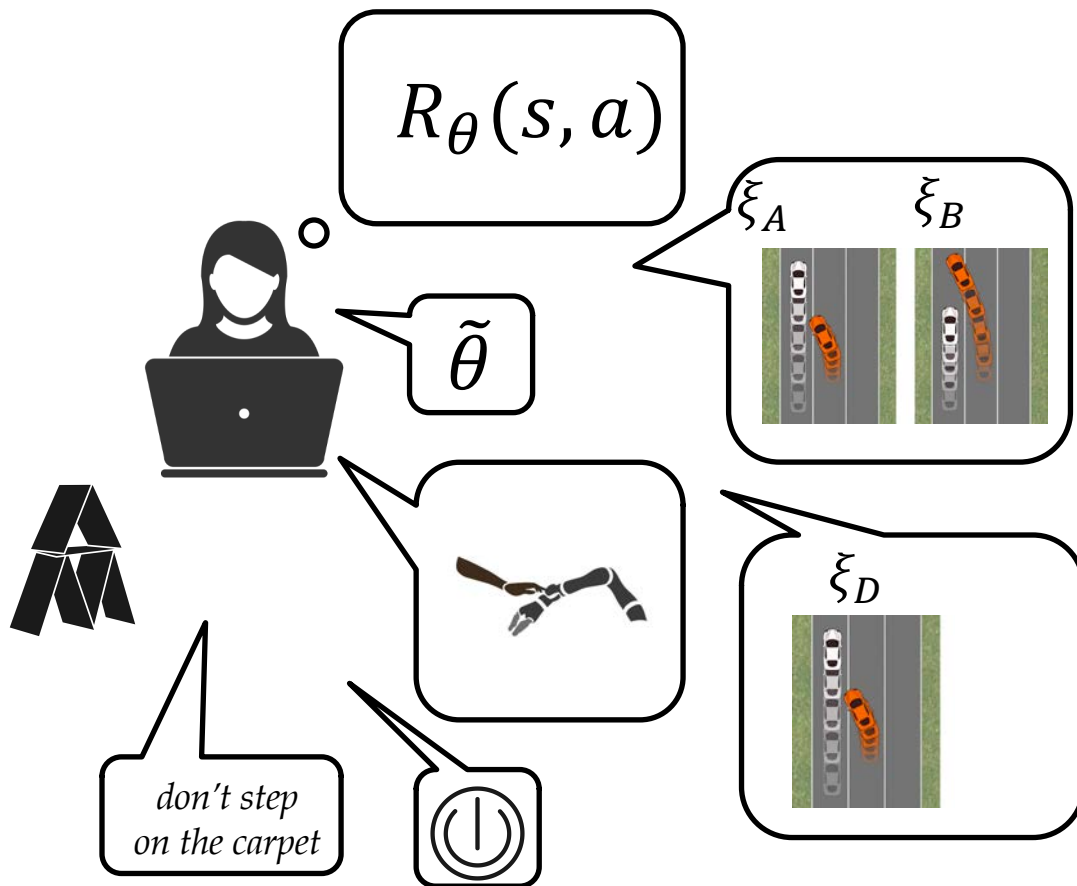
Desirable side effects: Batteries



AI \neq optimize specified reward

AI = optimize intended reward

Human feedback as reward-rational implicit choice



observation
(human) model

$$b'(\theta) \propto b(\theta) P(\text{🗨️} | \theta)$$

choices: $\{c\}$

choose based on reward: $\mathbb{E}[R_\theta(\xi) | \xi \sim \psi(c^*)]$

vs

$\mathbb{E}[R_\theta(\xi) | \xi \sim \psi(c)] \forall c$

$$P(c^* | \theta) = \frac{e^{\mathbb{E}[R_\theta(\xi) | \xi \sim \psi(c^*)]}}{\sum e^{\mathbb{E}[R_\theta(\xi) | \xi \sim \psi(c)]}}$$



Agents overlearn from specified rewards,
but leave other information on the table.

We can read the right amount of information into each source by interpreting them as reward-rational implicit choices.



Agents overlearn from specified rewards,
but leave other information on the table.

We can read the right amount of information into each source by
interpreting them as reward-rational implicit choices.

???

?????

*optimization, search, constraint
satisfaction, satisficing, RL...*

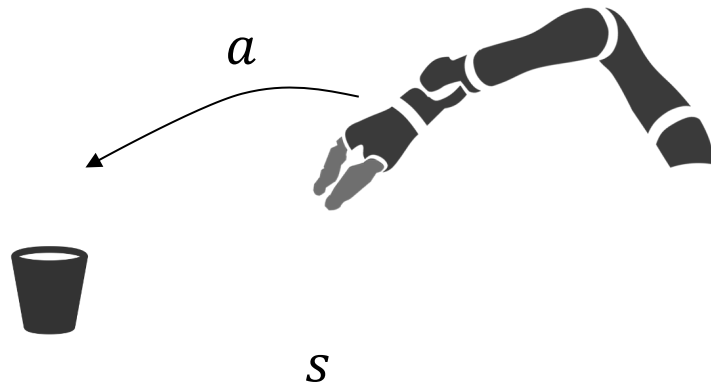
task **specification**

cost, reward, loss, constraints,...



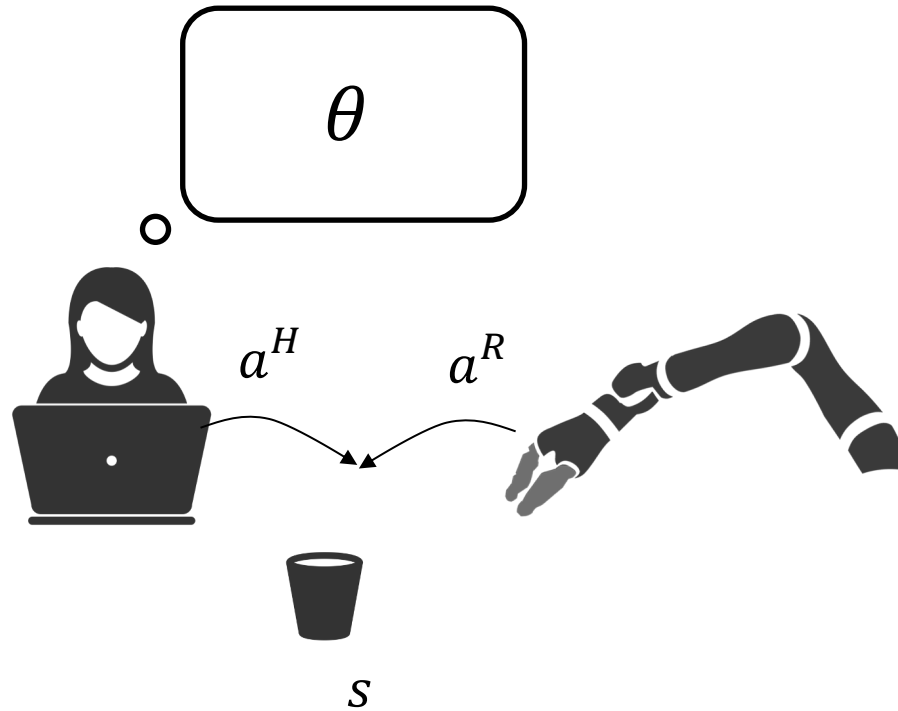
behavior

$$\max \mathbb{E}[\sum_t R_{\tilde{\theta}}(s_t, a_t)]$$



Assistance Games

$$\max \mathbb{E}[\sum_t R_\theta(s_t, a_t^R, a_t^H)]$$



Thanks!

