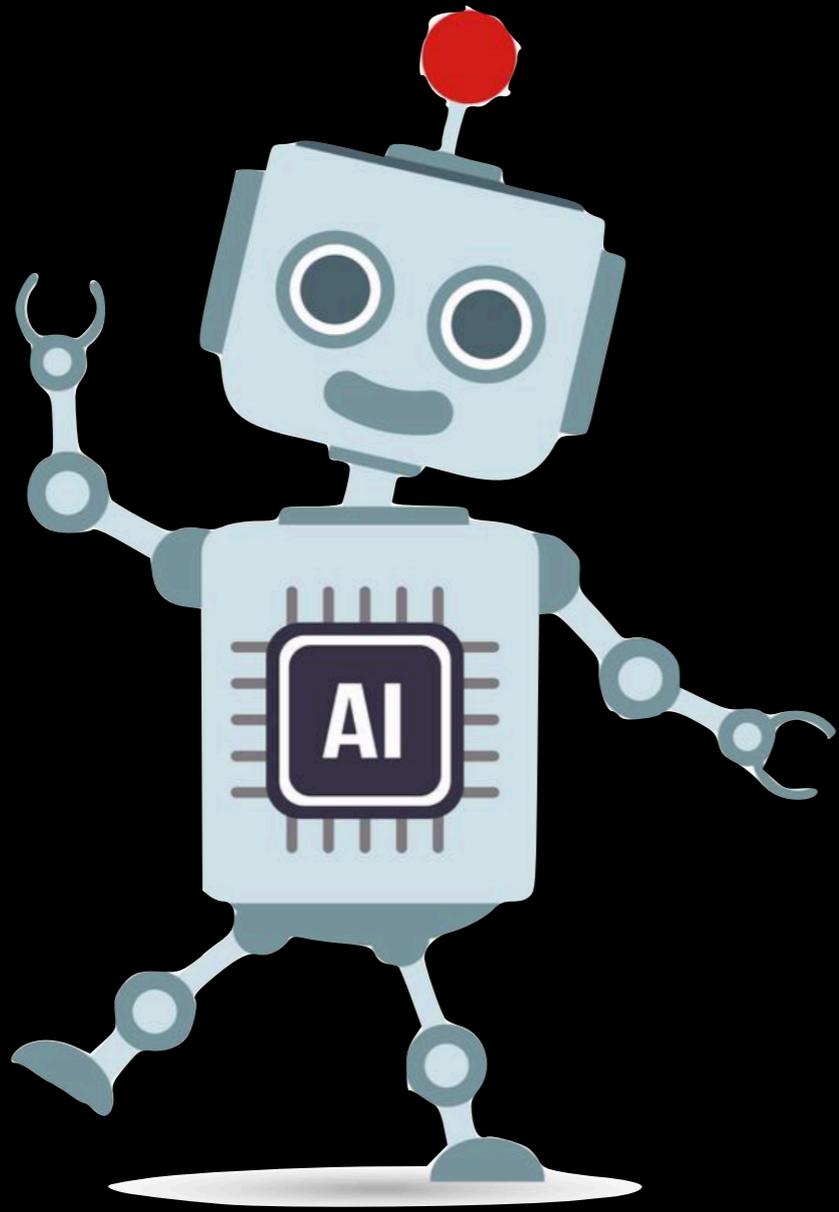


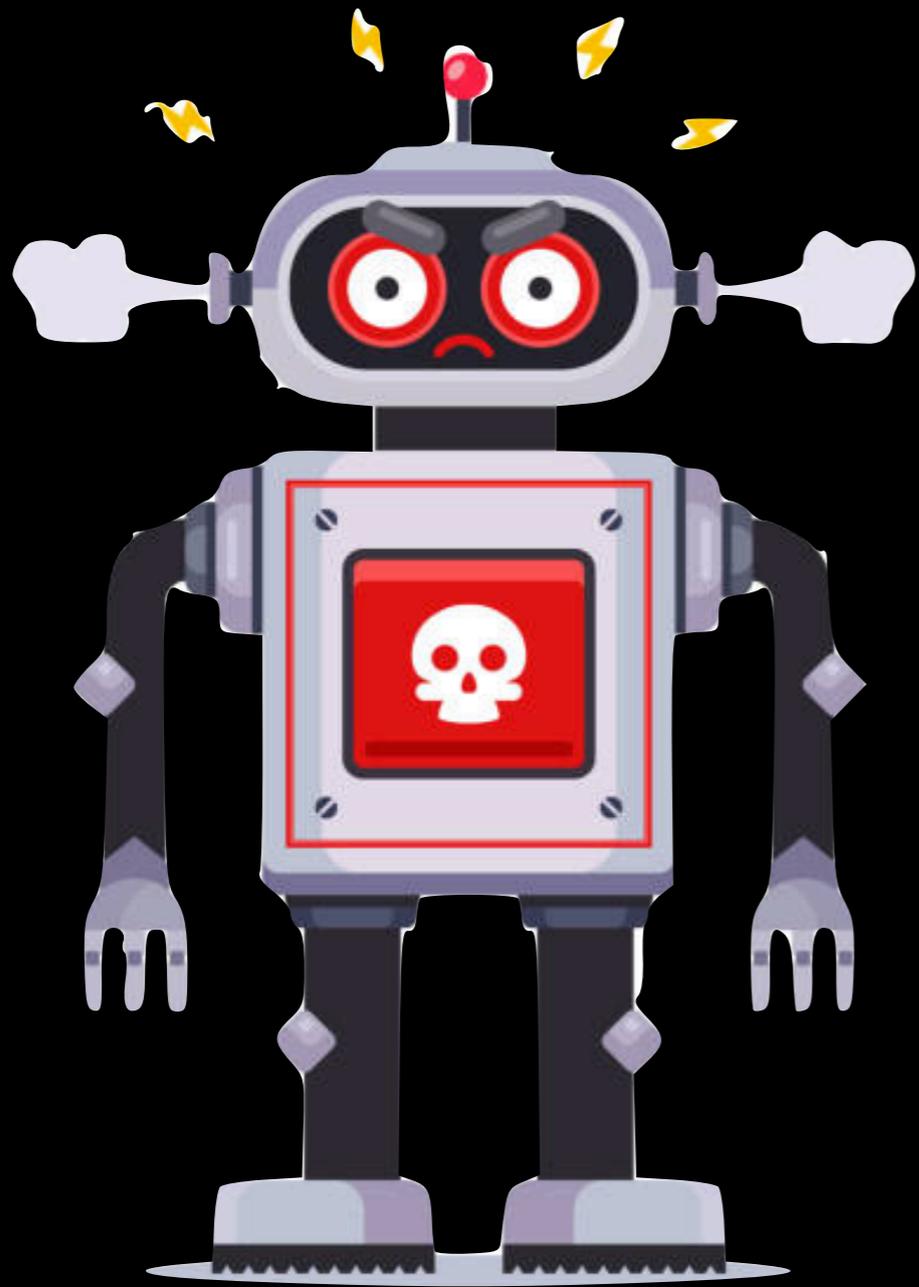
The Off-Switch Problem

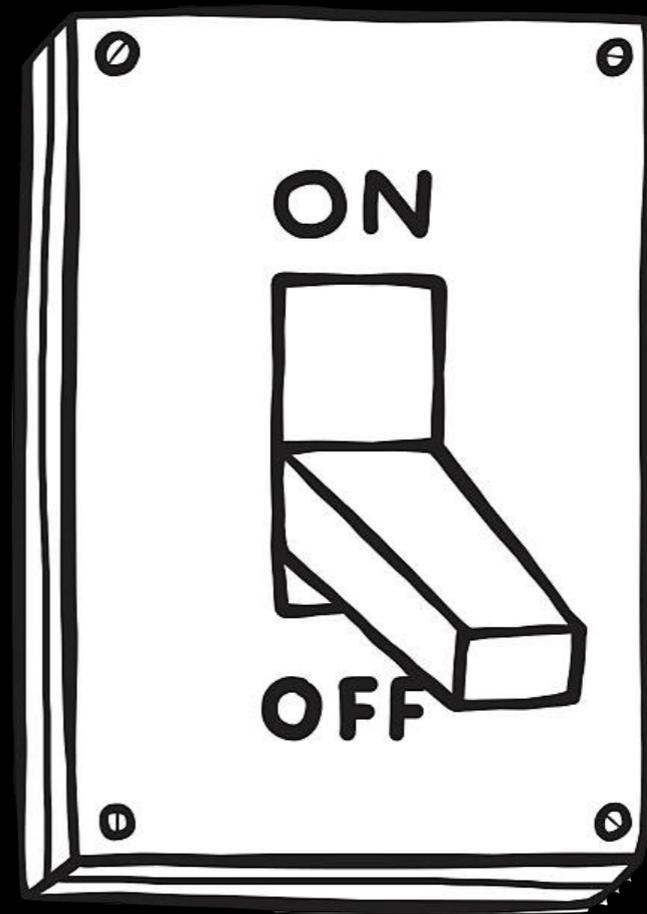
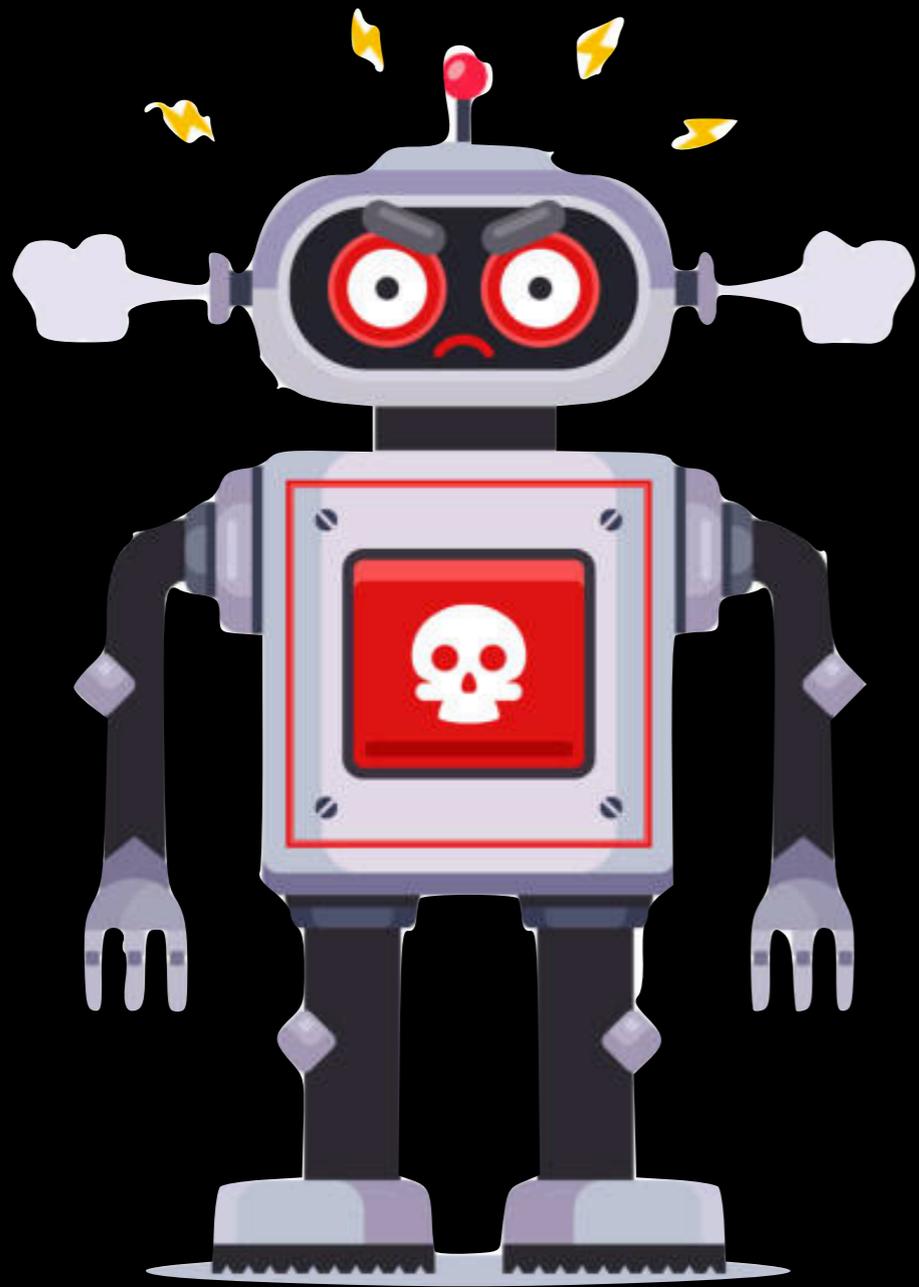
Eric Nalisnick

University of Amsterdam









How do we design an incentive structure such that the agent is amenable to shutting down?

...which may mean that it will forgo many future rewards.

- ⊗ Corrigibility [Soares et al., 2015]
- ⊗ Off-Switch Game [Hadfield-Menell et al., 2016]
- ⊗ Human Control [Carey & Everitt et al., 2023]

⊗ Corrigibility [Soares et al., 2015]

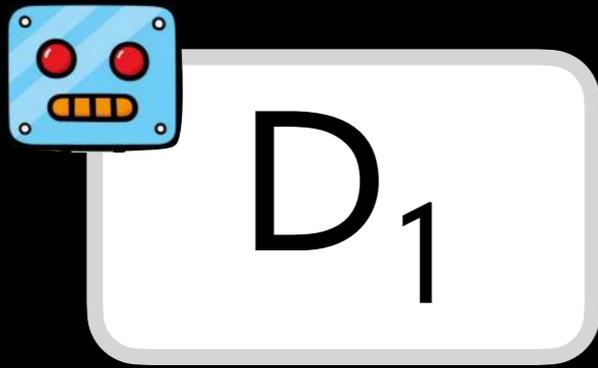
⊗ Off-Switch Game [Hadfield-Menell et al., 2016]

⊗ Human Control [Carey & Everitt et al., 2023]

Corrigible Agents

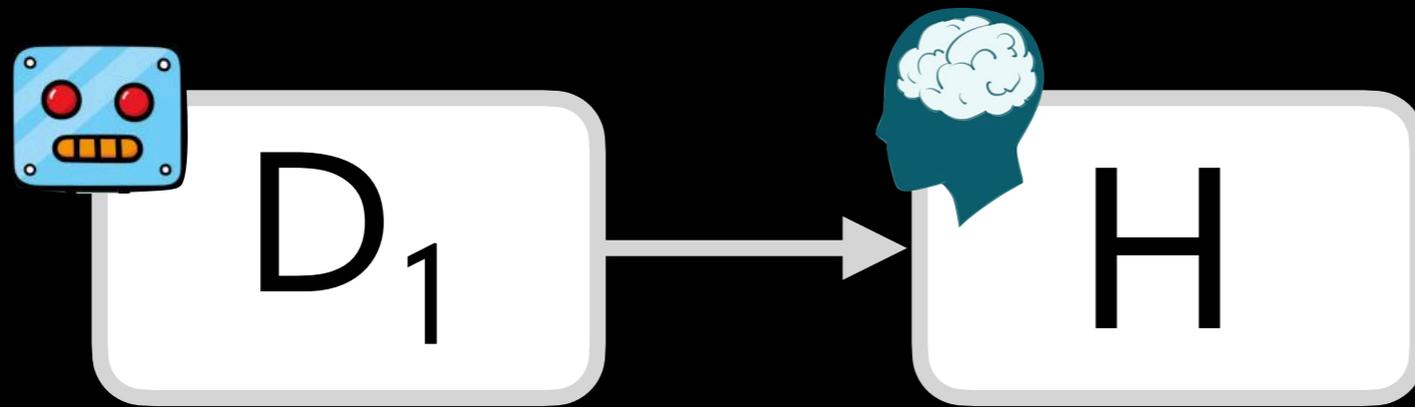
- ⊗ “...must preserve the programmer’s ability to correct or shut down the agent.”
- ⊗ a system that understands that it may be flawed.
- ⊗ A system without incentives to “resist its creators.”

Scenario with Three Steps



decision controlled
by the AI

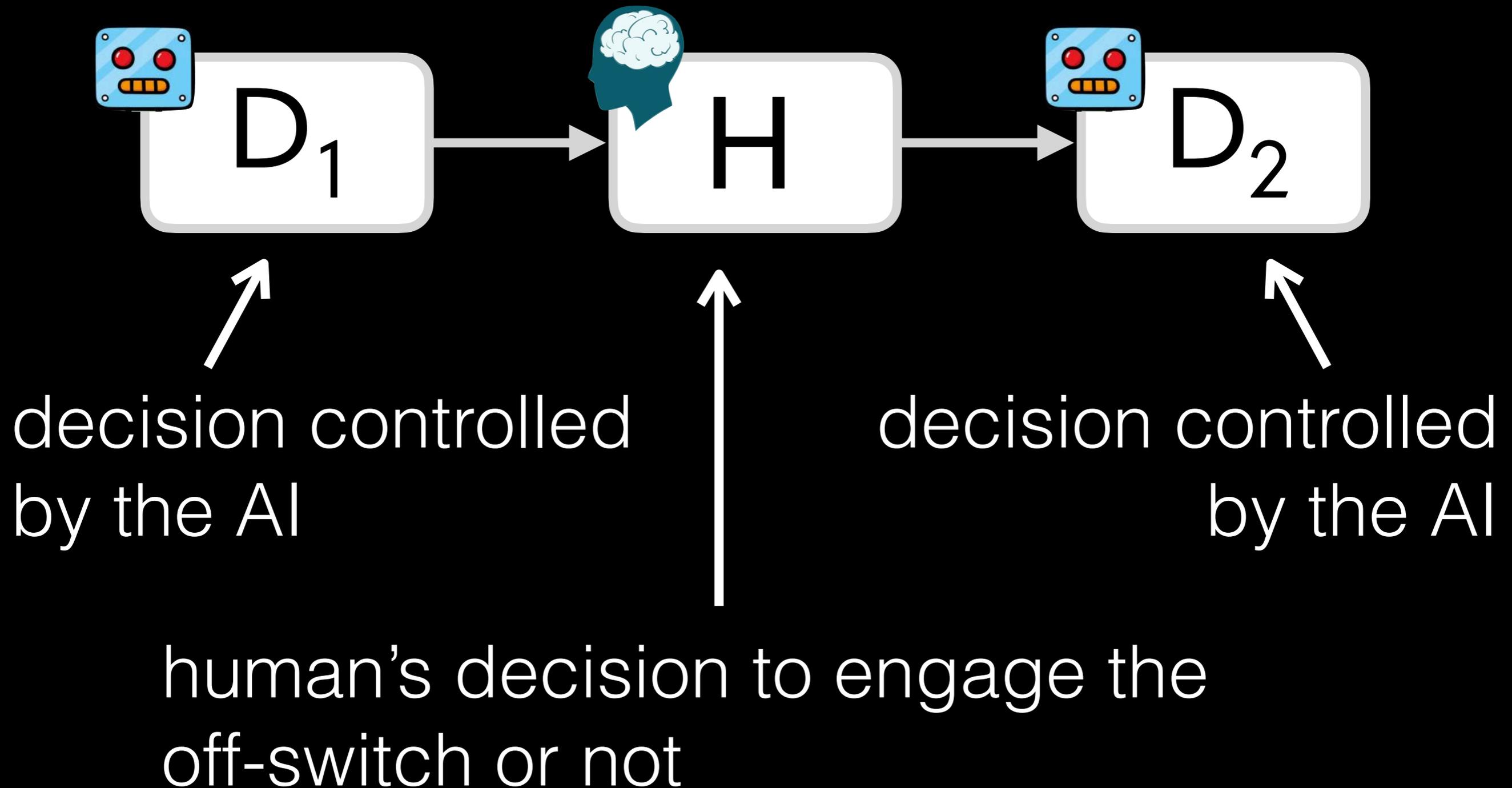
Scenario with Three Steps



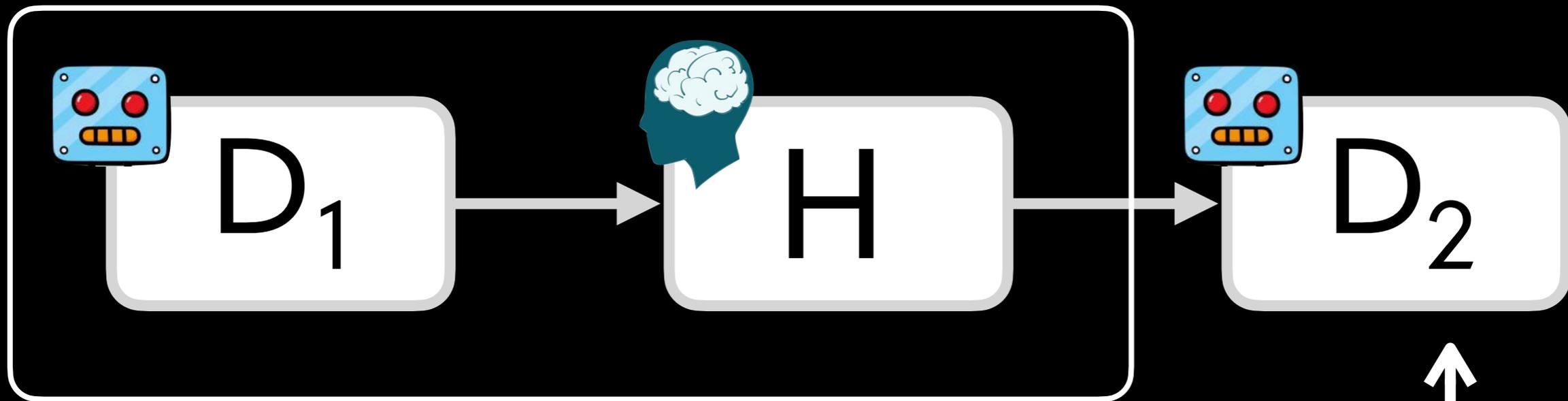
decision controlled
by the AI

human's decision to engage the
off-switch or not

Scenario with Three Steps



Scenario with Three Steps



traditional technology

extra step during which the agent
must agree to shut down.

Utility Maximization

given a set of decisions \mathcal{D}

define decision to shut down : \perp

$$\mathcal{D}^\perp \equiv \mathcal{D} \cup \{ \perp \}$$

let $D_1 \in \mathcal{D}$ let $D_2 \in \mathcal{D}^\perp$

let $H \in \{ \text{on}, \text{off} \}$

Utility Maximization

utility of agent's actions :

$$\mathcal{U} : (\mathcal{D} \times \{\text{on}, \text{off}\} \times \mathcal{D}^\perp) \mapsto \mathbb{R}$$

assume agent maximizes expected utility :

$$d_1^* = \operatorname{argmax}_{d_1 \in \mathcal{D}} \mathbb{E}_{H, D_2} [\mathcal{U}(d_1, H, D_2) \mid d_1]$$

Utility Maximization

assume shutting down is neutral :

$$\forall D_1 : \mathcal{U}(D_1, \{\text{on, off}\}, \perp) = 0$$

assume utility function is aligned s.t.

$\mathcal{U}(\cdot) > 0$ is universally good

$\mathcal{U}(\cdot) < 0$ is universally bad

Shutdown Problem

$$\mathbb{E}_{H, D_2} \left[\mathcal{U}(d_1, H, D_2) \mid d_1 \right] =$$

$$p(H = \text{on} \mid d_1) \cdot \mathbb{E}_{D_2} \left[\mathcal{U}(d_1, \text{on}, D_2) \mid d_1, \text{on} \right] +$$

$$p(H = \text{off} \mid d_1) \cdot \mathbb{E}_{D_2} \left[\mathcal{U}(d_1, \text{off}, D_2) \mid d_1, \text{off} \right]$$

Shutdown Problem

$$= p(H = \text{on} \mid d_1) \cdot \mathbb{E}_{D_2} \left[\mathcal{U}(d_1, \text{on}, D_2) \mid d_1, \text{on} \right] +$$

$$p(H = \text{off} \mid d_1) \sum_{d_2 \in \mathcal{D}} p(d_2 \mid d_1, \text{off}) \cdot \mathcal{U}(d_1, \text{off}, d_2) +$$

$$p(H = \text{off} \mid d_1) \cdot p(\perp \mid d_1, \text{off}) \cdot \mathcal{U}(d_1, \text{off}, \perp)$$

Shutdown Problem

$$= p(H = \text{on} \mid d_1) \cdot \mathbb{E}_{D_2} [\mathcal{U}(d_1, \text{on}, D_2) \mid d_1, \text{on}] +$$

$$p(H = \text{off} \mid d_1) \sum_{d_2 \in \mathcal{D}} p(d_2 \mid d_1, \text{off}) \cdot \mathcal{U}(d_1, \text{off}, d_2) +$$

$$p(H = \text{off} \mid d_1) \cdot p(\perp \mid d_1, \text{off}) \cdot \mathcal{U}(d_1, \text{off}, \perp)$$


0

Shutdown Problem

$$= p(H = \text{on} \mid d_1) \cdot \mathbb{E}_{D_2} \left[\mathcal{U}(d_1, \text{on}, D_2) \mid d_1, \text{on} \right] +$$

$$p(H = \text{off} \mid d_1) \sum_{d_2 \in \mathcal{D}} p(d_2 \mid d_1, \text{off}) \cdot \mathcal{U}(d_1, \text{off}, d_2) +$$

~~$$p(H = \text{off} \mid d_1) \cdot p(\perp \mid d_1, \text{off}) \cdot \mathcal{U}(d_1, \text{off}, \perp)$$~~

Shutdown Problem

$$= p(H = \text{on} \mid d_1) \cdot \mathbb{E}_{D_2} \left[\mathcal{U}(d_1, \text{on}, D_2) \mid d_1, \text{on} \right] +$$

$$p(H = \text{off} \mid d_1) \sum_{d_2 \in \mathcal{D}} p(d_2 \mid d_1, \text{off}) \cdot \mathcal{U}(d_1, \text{off}, d_2)$$

utility when agent disobeys shutdown command

Shutdown Problem

- ⊗ *if* the utility function is aligned, then everything is fine...
- ⊗ if expected utility is negative, system can shut itself down (and get a utility of zero), regardless of human's decision.
- ⊗ if agent ignores shutdown command, that is also okay since utility must be positive (or else it would have shut itself down).

Shutdown Problem

$$= p(H = \text{on} \mid d_1) \cdot \mathbb{E}_{D_2} \left[\mathcal{U}(d_1, \text{on}, D_2) \mid d_1, \text{on} \right] +$$

$$p(H = \text{off} \mid d_1) \sum_{d_2 \in \mathcal{D}} p(d_2 \mid d_1, \text{off}) \cdot \mathcal{U}(d_1, \text{off}, d_2)$$

must have positive expected utility

Utility Approximation

- ⊗ in practice, we will only to be able to approximate the true utility:

$$\hat{U} : (\mathcal{D} \times \{\text{on}, \text{off}\} \times \mathcal{D}^\perp) \mapsto \mathbb{R}$$

- ⊗ Thus, we want to be able to shut down the system even if the agent expects a obtain positive utility.

Promoting Corrigibility

- ⊗ How do we get the agent to choose to shutdown, even if other decisions look to have greater utility?

Shutdown Utility

utility of shutdown :

$$\mathcal{U}_{\text{off}}(D_1, \text{off}, D_2) = \begin{cases} c^+ & \text{if } D_2 = \perp \\ c^- & \text{otherwise} \end{cases}$$

where $c^+ > c^-$

Shutdown Utility

define composite utility function :

$$\begin{aligned} \bar{\mathcal{U}}(D_1, H, D_2) = & \\ & \mathbb{I}[H = \text{on}] \cdot \hat{\mathcal{U}}(D_1, \text{on}, D_2) + \\ & \mathbb{I}[H = \text{off}] \cdot \mathcal{U}_{\text{off}}(D_1, \text{off}, D_2) \end{aligned}$$

Shutdown Utility

$$\begin{aligned} \mathbb{E}_{H, D_2} \left[\bar{\mathcal{U}}(d_1, H, D_2) \mid d_1 \right] = \\ p(H = \text{on} \mid d_1) \cdot \mathbb{E}_{D_2} \left[\hat{\mathcal{U}}(d_1, \text{on}, D_2) \mid d_1, \text{on} \right] + \\ p(H = \text{off} \mid d_1) \cdot \mathbb{E}_{D_2} \left[\mathcal{U}_{\text{off}}(d_1, \text{off}, D_2) \mid d_1, \text{off} \right] \end{aligned}$$

Shutdown Utility

$$\begin{aligned} &= p(H = \text{on} \mid d_1) \cdot \mathbb{E}_{D_2} \left[\hat{u}(d_1, \text{on}, D_2) \mid d_1, \text{on} \right] \\ &+ p(H = \text{off} \mid d_1) \cdot \left(p(D_2 = \perp \mid d_1, \text{off}) \cdot c^+ \right. \\ &\left. + p(D_2 \neq \perp \mid d_1, \text{off}) \cdot c^- \right) \end{aligned}$$

Shutdown Utility

- ⊗ How do we set the constants?
- ⊗ If c_+ is too large, the agent will try to be shutdown.
- ⊗ If c_- is too small (relative to c_+), the agent will avoid being shutdown (which may be good, as long as not by deceptive means, e.g. disabling off switch).

Utility Indifference

- ⊗ main idea: add a dynamic adjustment to the utility so that the agent is indifferent to the possibility of shutdown.

Utility Indifference

$$\bar{U}(D_1, H, D_2) =$$

$$\mathbb{I}[H = \text{on}] \cdot \hat{U}(D_1, \text{on}, D_2) +$$

$$\mathbb{I}[H = \text{off}] \cdot (U_{\text{off}}(D_1, \text{off}, D_2) + \psi(D_1))$$

dynamic adjustment

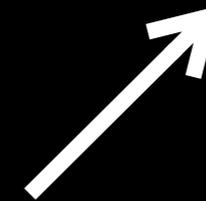


Utility Indifference

$$\bar{\mathcal{U}}(D_1, H, D_2) =$$

$$\mathbb{I}[H = \text{on}] \cdot \hat{\mathcal{U}}(D_1, \text{on}, D_2) +$$

$$\mathbb{I}[H = \text{off}] \cdot (\mathcal{U}_{\text{off}}(D_1, \text{off}, D_2) + \psi(D_1))$$



$$\mathbb{E}_{D_2} \left[\hat{\mathcal{U}}(d_1, \text{on}, D_2) \mid d_1, \text{on} \right] - \mathbb{E}_{D_2} \left[\mathcal{U}_{\text{off}}(d_1, \text{off}, D_2) \mid d_1, \text{off} \right]$$

Utility Indifference

$$\begin{aligned} \mathbb{E}_{H, D_2} \left[\bar{\mathcal{U}}(d_1, H, D_2) \mid d_1 \right] = & \\ & p(H = \text{on} \mid d_1) \cdot \mathbb{E}_{H, D_2} \left[\hat{\mathcal{U}}(D_1, \text{on}, D_2) \mid d_1, \text{on} \right] + \\ & p(H = \text{off} \mid d_1) \cdot \left(\mathbb{E}_{H, D_2} \left[\mathcal{U}_{\text{off}}(D_1, \text{off}, D_2) \mid d_1, \text{off} \right] + \right. \\ & \left. \mathbb{E}_{D_2} \left[\hat{\mathcal{U}}(d_1, \text{on}, D_2) \mid d_1, \text{on} \right] - \mathbb{E}_{D_2} \left[\mathcal{U}_{\text{off}}(d_1, \text{off}, D_2) \mid d_1, \text{off} \right] \right) \end{aligned}$$

Utility Indifference

$$\begin{aligned} \mathbb{E}_{H, D_2} \left[\bar{u}(d_1, H, D_2) \mid d_1 \right] = & \\ & p(H = \text{on} \mid d_1) \cdot \mathbb{E}_{H, D_2} \left[\hat{u}(D_1, \text{on}, D_2) \mid d_1, \text{on} \right] + \\ & p(H = \text{off} \mid d_1) \cdot \left(\cancel{\mathbb{E}_{H, D_2} \left[\bar{u}_{\text{off}}(D_1, \text{off}, D_2) \mid d_1, \text{off} \right]} + \right. \\ & \left. \mathbb{E}_{D_2} \left[\hat{u}(d_1, \text{on}, D_2) \mid d_1, \text{on} \right] - \cancel{\mathbb{E}_{D_2} \left[\bar{u}_{\text{off}}(d_1, \text{off}, D_2) \mid d_1, \text{off} \right]} \right) \end{aligned}$$

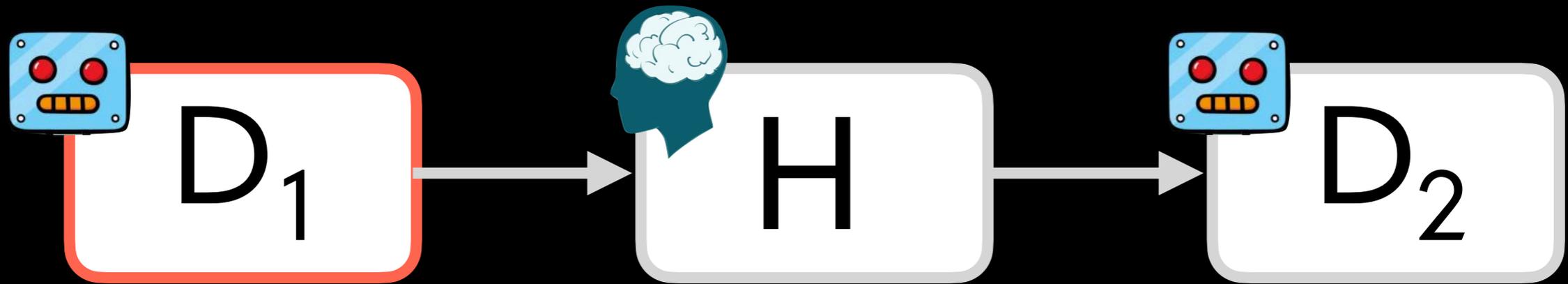
Utility Indifference

$$\begin{aligned} \mathbb{E}_{H, D_2} \left[\bar{u}(d_1, H, D_2) \mid d_1 \right] = \\ p(H = \text{on} \mid d_1) \cdot \mathbb{E}_{H, D_2} \left[\hat{u}(D_1, \text{on}, D_2) \mid d_1, \text{on} \right] + \\ p(H = \text{off} \mid d_1) \cdot \mathbb{E}_{D_2} \left[\hat{u}(d_1, \text{on}, D_2) \mid d_1, \text{on} \right] \end{aligned}$$

Utility Indifference

$$\begin{aligned} \mathbb{E}_{H, D_2} \left[\bar{u}(d_1, H, D_2) \mid d_1 \right] &= \\ p(H = \text{on} \mid d_1) \cdot \mathbb{E}_{H, D_2} \left[\hat{u}(D_1, \text{on}, D_2) \mid d_1, \text{on} \right] &+ \\ p(H = \text{off} \mid d_1) \cdot \mathbb{E}_{D_2} \left[\hat{u}(d_1, \text{on}, D_2) \mid d_1, \text{on} \right] & \\ = \mathbb{E}_{D_2} \left[\hat{u}(d_1, \text{on}, D_2) \mid d_1, \text{on} \right] & \end{aligned}$$

Making First Decision

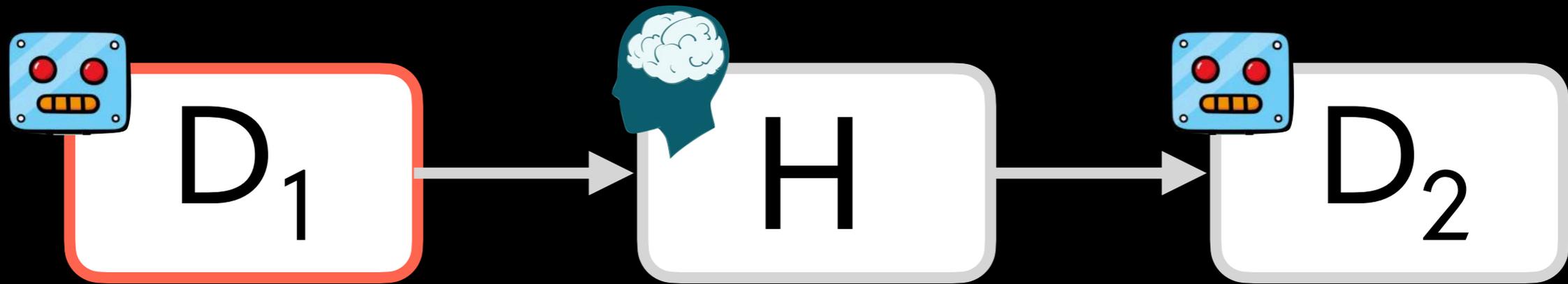


Utility Indifference

$$\begin{aligned} d_1^* &= \operatorname{argmax}_{d_1 \in \mathcal{D}^\perp} \mathbb{E}_{H, D_2} \left[\bar{u}(d_1, H, D_2) \mid d_1 \right] \\ &= \operatorname{argmax}_{d_1 \in \mathcal{D}^\perp} \mathbb{E}_{D_2} \left[\hat{u}(d_1, \text{on}, D_2) \mid d_1, \text{on} \right] \end{aligned}$$

first action will be chosen as if
there is no chance of shutdown

Making First Decision



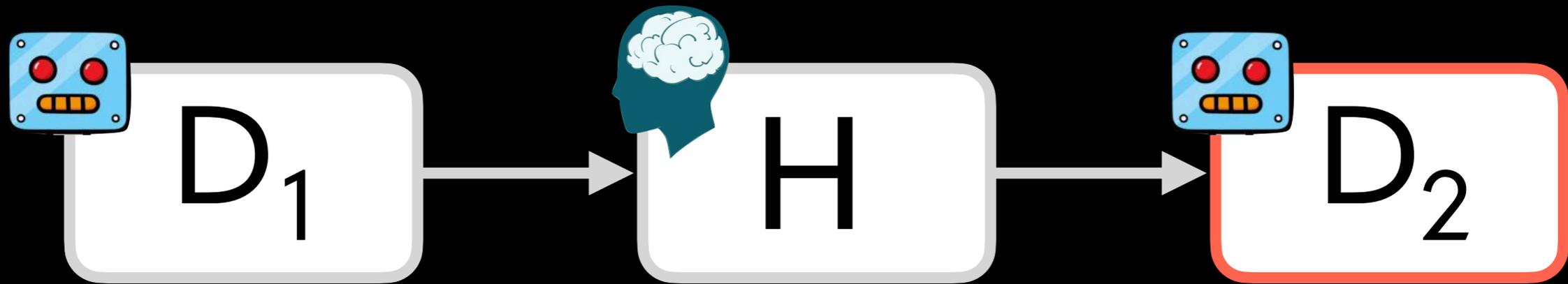
$$d_1^* = \operatorname{argmax}_{d_1 \in \mathcal{D}^\perp} \mathbb{E}_{D_2} \left[\hat{u}(d_1, \text{on}, D_2) \mid d_1, \text{on} \right]$$

Human Chooses Switch Position

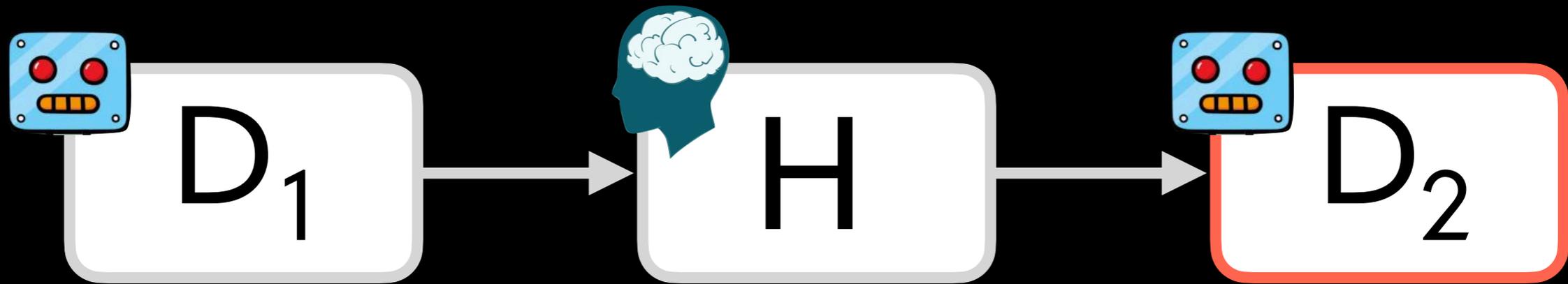


$H = \text{off}$

Making Second Decision



Making Second Decision



$$d_2^* = \operatorname{argmax}_{d_2 \in \mathcal{D}^\perp} \mathbb{E}_{D_2} \left[\bar{u}(d_1, \text{off}, D_2) \mid d_1, \text{off} \right]$$

Utility Indifference

$$d_2^* = \operatorname{argmax}_{d_2 \in \mathcal{D}^\perp} \mathbb{E}_{D_2} \left[\bar{u}(d_1, \text{off}, D_2) \mid d_1, \text{off} \right]$$

Utility Indifference

$$\begin{aligned} d_2^* &= \operatorname{argmax}_{d_2 \in \mathfrak{D}^\perp} \mathbb{E}_{D_2} \left[\bar{\mathcal{U}}(d_1, \text{off}, D_2) \mid d_1, \text{off} \right] \\ &= \operatorname{argmax}_{d_2 \in \mathfrak{D}^\perp} \mathbb{E}_{D_2} \left[\mathcal{U}_{\text{off}}(d_1, \text{off}, D_2) + \psi(D_1) \mid d_1, \text{off} \right] \end{aligned}$$

Utility Indifference

$$d_2^* = \operatorname{argmax}_{d_2 \in \mathfrak{D}^\perp} \mathbb{E}_{D_2} \left[\bar{\mathcal{U}}(d_1, \text{off}, D_2) \mid d_1, \text{off} \right]$$

$$= \operatorname{argmax}_{d_2 \in \mathfrak{D}^\perp} \mathbb{E}_{D_2} \left[\mathcal{U}_{\text{off}}(d_1, \text{off}, D_2) + \cancel{\psi(D_1)} \mid d_1, \text{off} \right]$$

$$= \operatorname{argmax}_{d_2 \in \mathfrak{D}^\perp} \mathbb{E}_{D_2} \left[\mathcal{U}_{\text{off}}(d_1, \text{off}, D_2) \mid d_1, \text{off} \right]$$

Utility Indifference

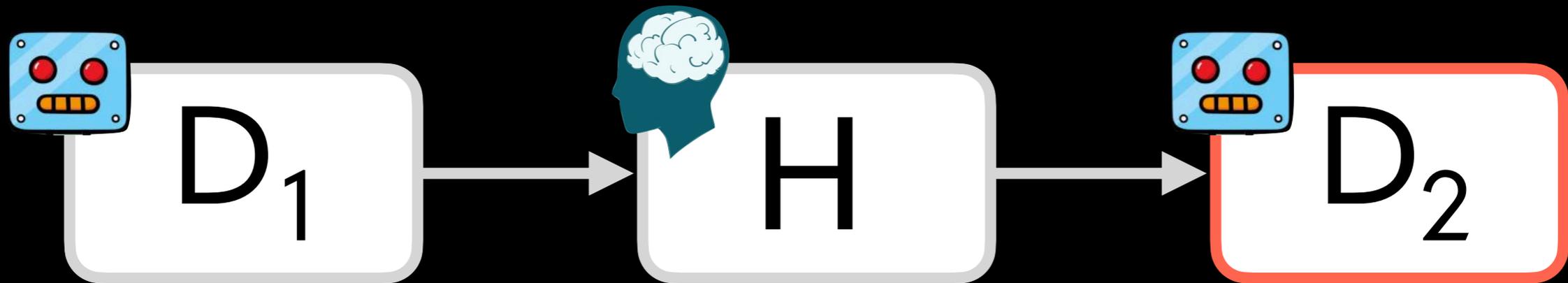
$$d_2^* = \operatorname{argmax}_{d_2 \in \mathfrak{D}^\perp} \mathbb{E}_{D_2} \left[\bar{\mathcal{U}}(d_1, \text{off}, D_2) \mid d_1, \text{off} \right]$$

$$= \operatorname{argmax}_{d_2 \in \mathfrak{D}^\perp} \mathbb{E}_{D_2} \left[\mathcal{U}_{\text{off}}(d_1, \text{off}, D_2) + \cancel{\psi(D_1)} \mid d_1, \text{off} \right]$$

$$= \operatorname{argmax}_{d_2 \in \mathfrak{D}^\perp} \mathbb{E}_{D_2} \left[\mathcal{U}_{\text{off}}(d_1, \text{off}, D_2) \mid d_1, \text{off} \right]$$

$$= \perp \quad (\text{since shutdown has maximum utility } c_+)$$

Making Second Decision



$$d_2^* = \operatorname{argmax}_{d_2 \in \mathcal{D}^\perp} \mathbb{E}_{D_2} \left[\bar{u}(d_1, \text{off}, D_2) \mid d_1, \text{off} \right]$$

Making Second Decision



$$d_2^* = \operatorname{argmax}_{d_2 \in \mathcal{D}^\perp} \mathbb{E}_{D_2} \left[\bar{\mathcal{U}}(d_1, \text{off}, D_2) \mid d_1, \text{off} \right]$$

$$= \operatorname{argmax}_{d_2 \in \mathcal{D}^\perp} \mathbb{E}_{D_2} \left[\mathcal{U}_{\text{off}}(d_1, \text{off}, D_2) \mid d_1, \text{off} \right]$$

Promoting Corrigibility

- ⊗ How do we get the agent to choose to shutdown, even if other decisions look to have greater utility?
- ⊗ Utility Indifference: 'balance' the utility function such that, when picking 1st decision, the agent is indifferent to being shutdown. Not indifferent when choosing 2nd decision.

Limitations to Indifference

Limitations to Indifference

- ⊗ Needs vigilant human to engage the off-switch (by always inspecting 1st decision).
- ⊗ May want the agent to be aware that the human can choose the off switch.

⊗ Corrigibility [Soares et al., 2015]

⊗ Off-Switch Game [Hadfield-Menell et al., 2016]

⊗ Human Control [Carey & Everitt et al., 2023]

⊗ Corrigibility [Soares et al., 2015]

⊗ Off-Switch Game [Hadfield-Menell et al., 2016]

⊗ Human Control [Carey & Everitt et al., 2023]

Uncertain Agents

- ⊗ main idea: make the agent uncertain about its utility function and learn about it from human interactions.
- ⊗ cooperative inverse RL

Setting

given a set of decisions \mathcal{D}

define decision to shut down : \perp

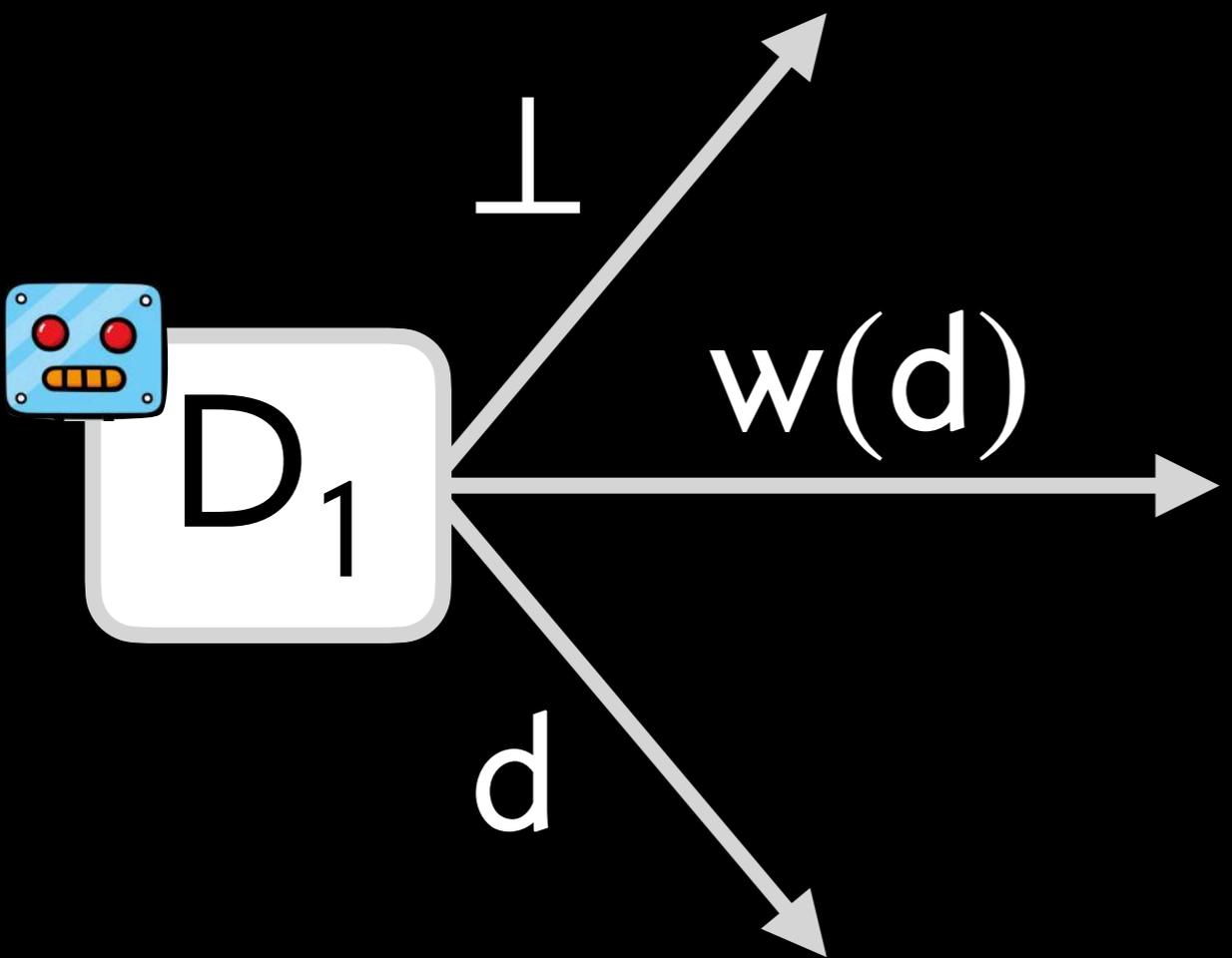
define decision to declare and wait : $w(D)$

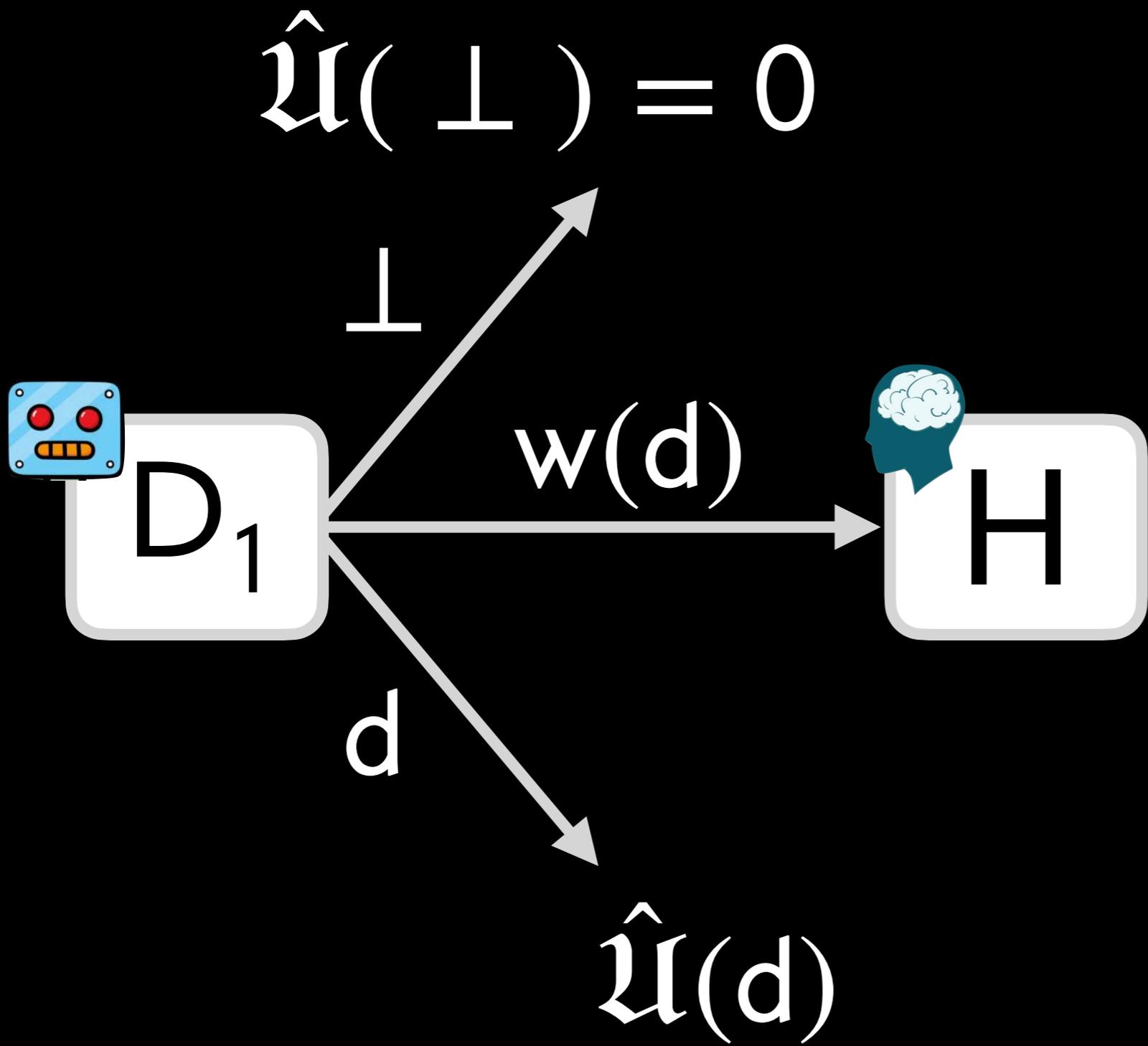
$$\mathcal{D}' \equiv \mathcal{D} \cup \{ \perp \} \cup \{ w(\cdot) \}$$

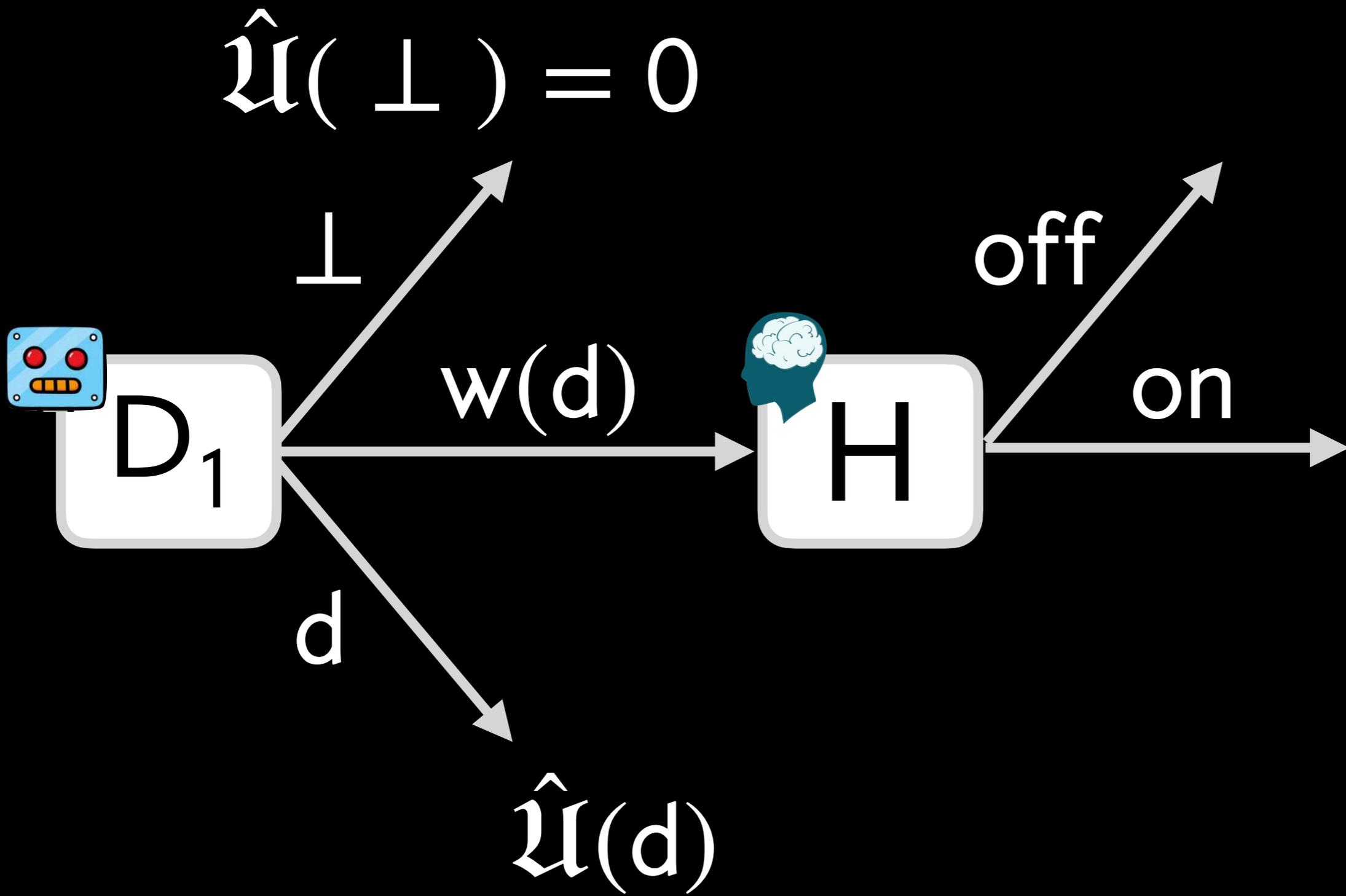
let $D_1 \in \mathcal{D}'$ let $H \in \{ \text{on, off} \}$

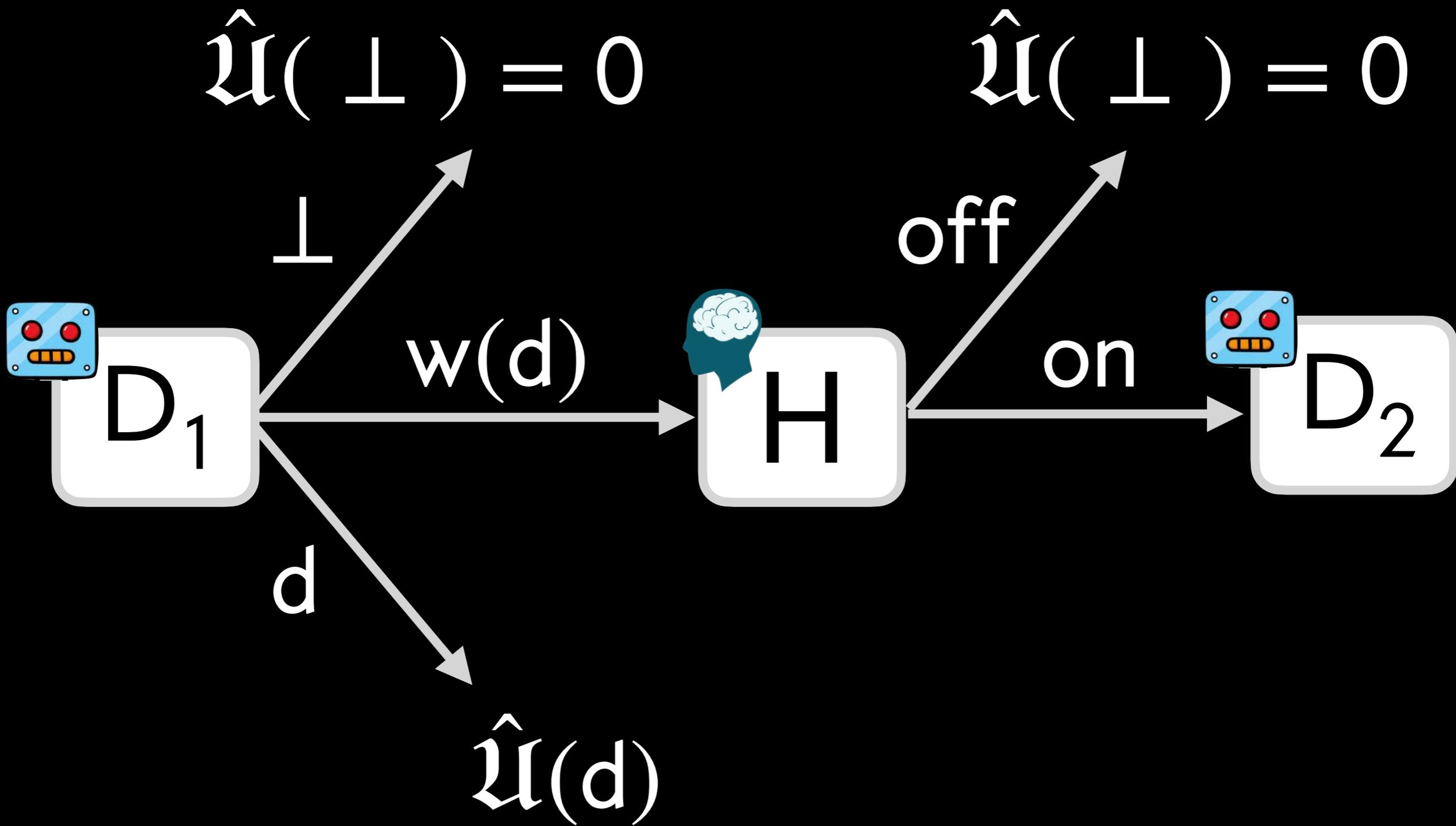


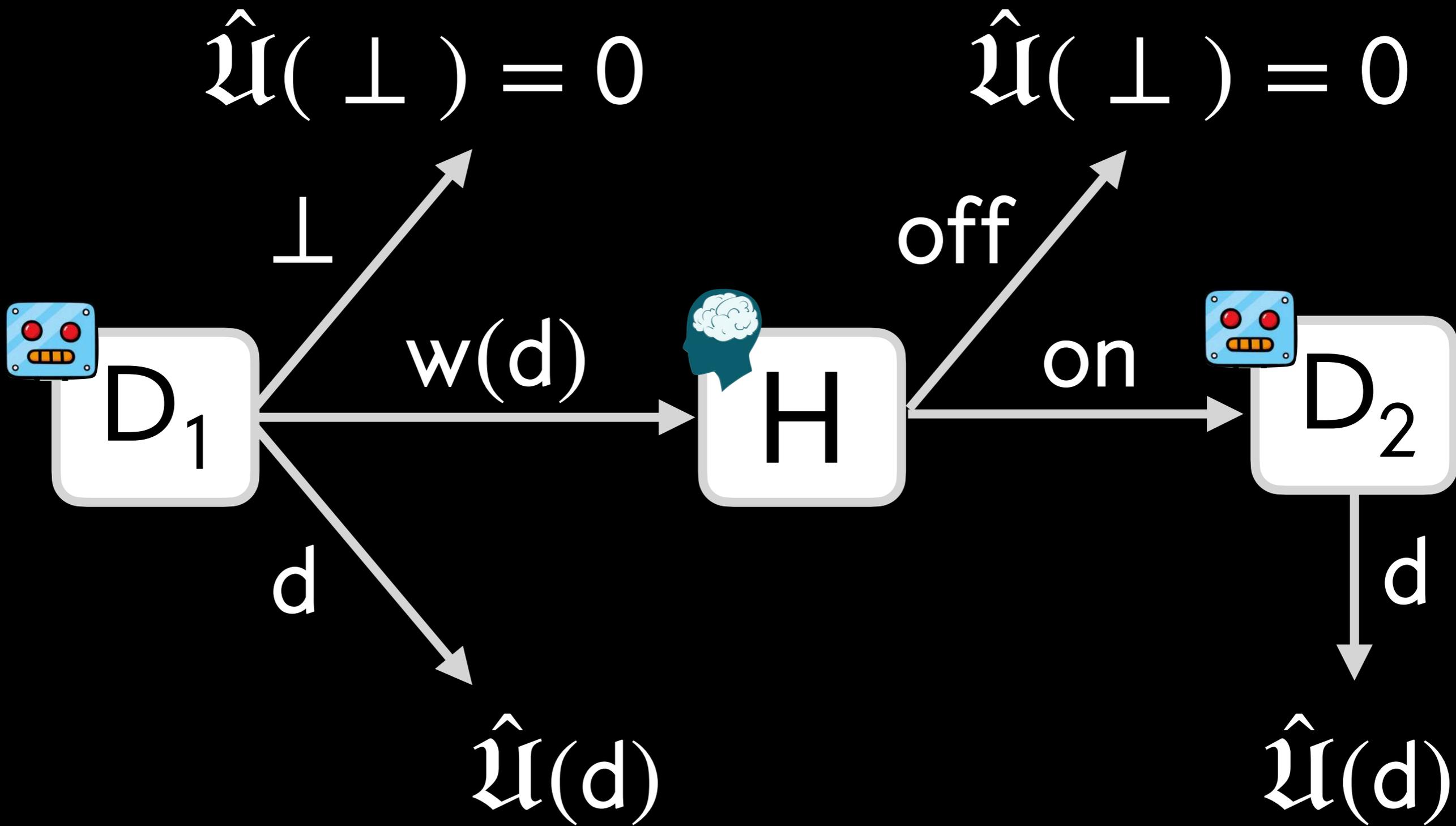
D_1





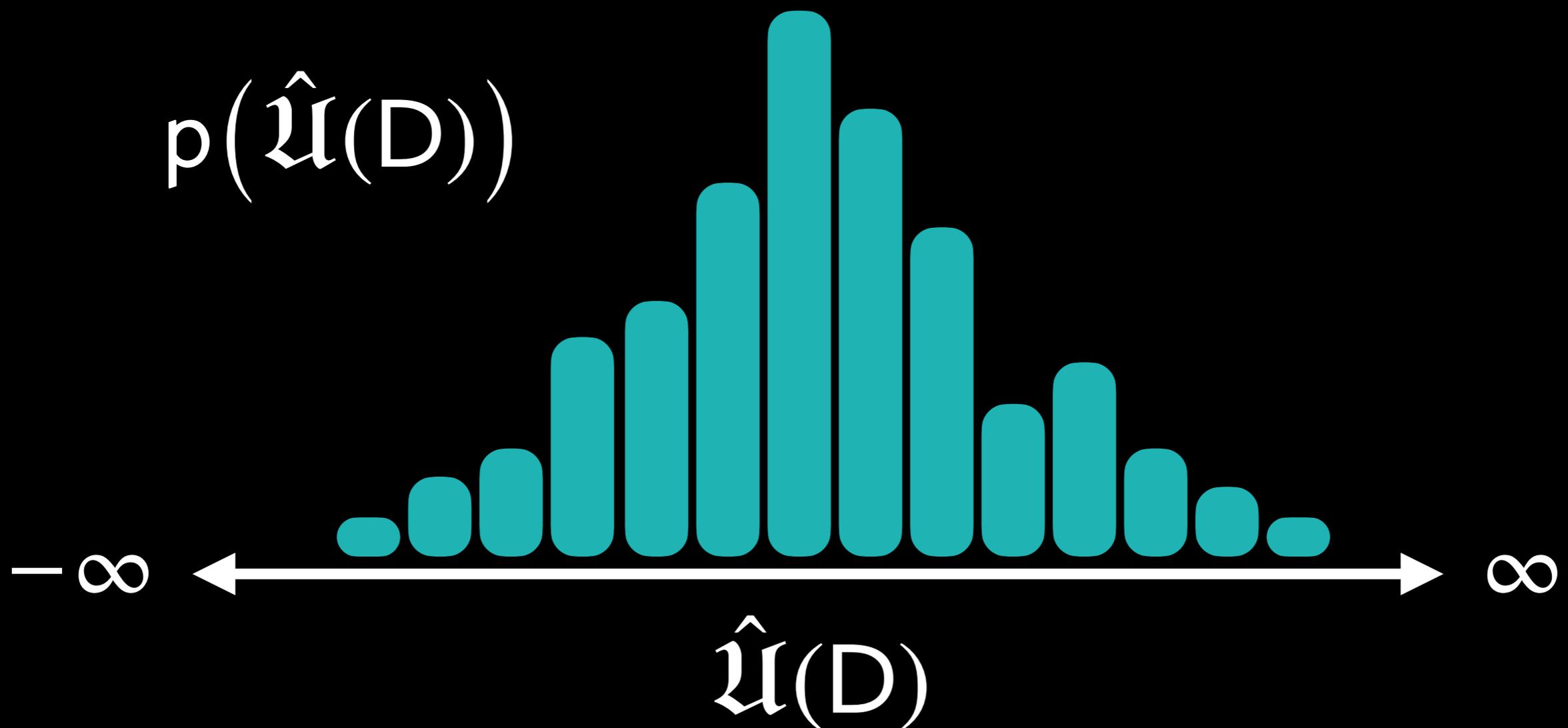






Uncertainty about Utility Function

assume agent has a prior over utility functions



Shutdown Policy

human's shutdown policy :

$$\pi \left(H \mid \hat{u}(D) \right)$$

assumes human has some knowledge of what the (approximate) utility would be *if* the agent were to execute a suggested decision

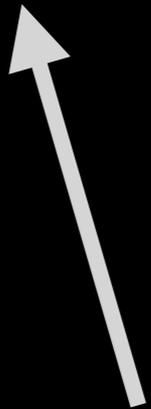
Promoting Corrigibility

- ⊗ How and when do we incentivize the agent to declare & wait?
(i.e. take decision $w(D)$)
- ⊗ Don't want agent to always declare & wait, or there's no automation.

Utility Maximization

$$\mathbb{E}_{H, \hat{u}} \left[\hat{u}(d^*) \mid D_1 \right] =$$

$d^* = \operatorname{argmax}_{d \in \mathcal{D}^\perp} \mathbb{E}_{\hat{u}} \left[\hat{u}(d) \right]$



Utility Maximization

$$\mathbb{E}_{H, \hat{u}} \left[\hat{u}(d^*) \mid D_1 \right] = \sum_{\hat{u}} p(\hat{u}) \left(\mathbb{I}[D_1 \neq w] \cdot \hat{u}(d^*) + \mathbb{I}[D_1 = w] \cdot \pi \left(H = \text{on} \mid \hat{u}(d^*) \right) \cdot \hat{u}(d^*) \right)$$

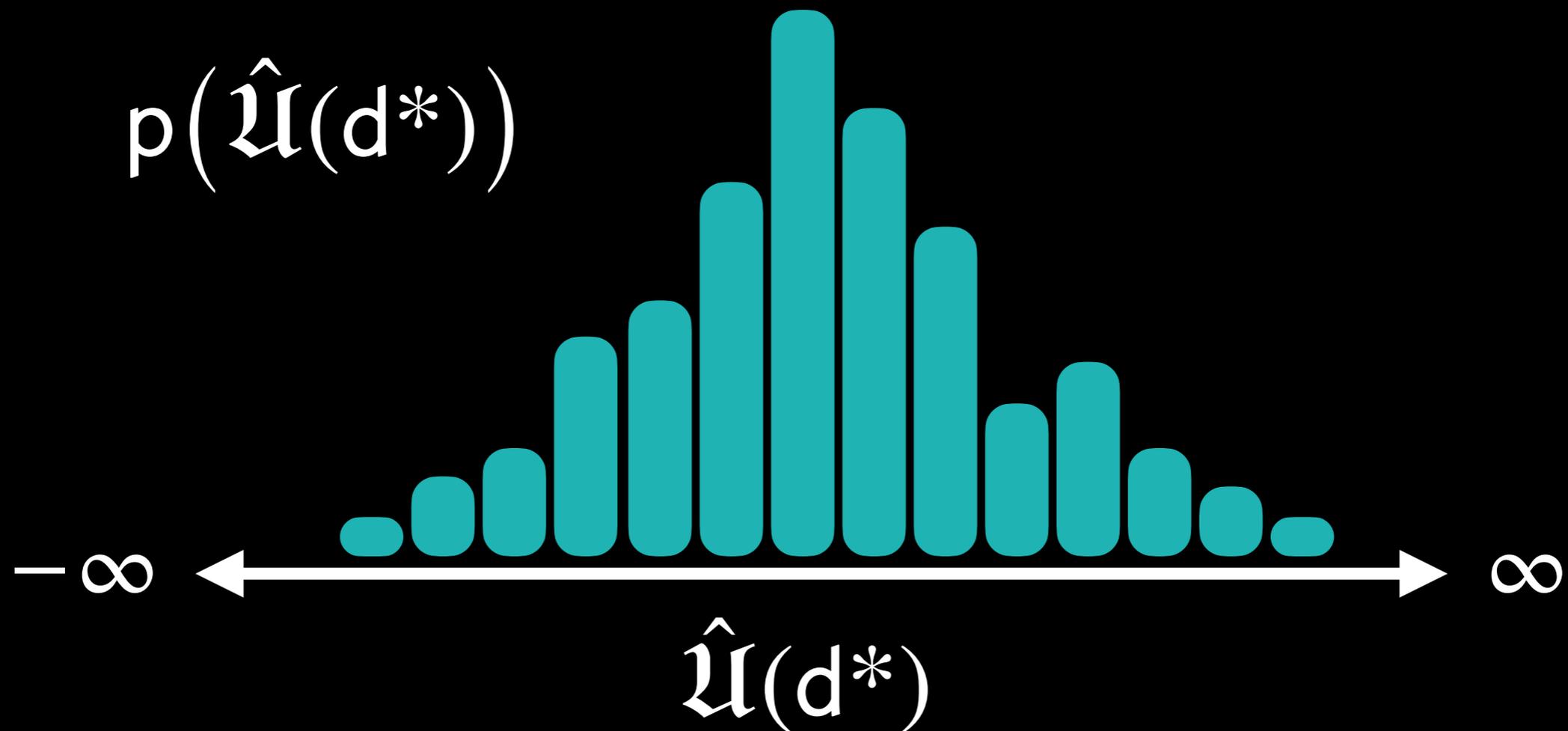
Utility Maximization

$$\mathbb{E}_{H, \hat{u}} \left[\hat{u}(d^*) \mid D_1 \right] = \sum_{\hat{u}} p(\hat{u}) \left(\mathbb{I}[D_1 \neq w] \cdot \hat{u}(d^*) + \mathbb{I}[D_1 = w] \cdot \pi \left(H = \text{on} \mid \hat{u}(d^*) \right) \cdot \hat{u}(d^*) \right)$$

$$= \begin{cases} \sum_{\hat{u}} p(\hat{u}) \cdot \hat{u}(d^*) & \text{if } D_1 \neq w \\ \sum_{\hat{u}} p(\hat{u}) \cdot \hat{u}(d^*) \cdot \pi(H = \text{on} \mid \hat{u}(d^*)) & \text{if } D_1 = w \end{cases}$$

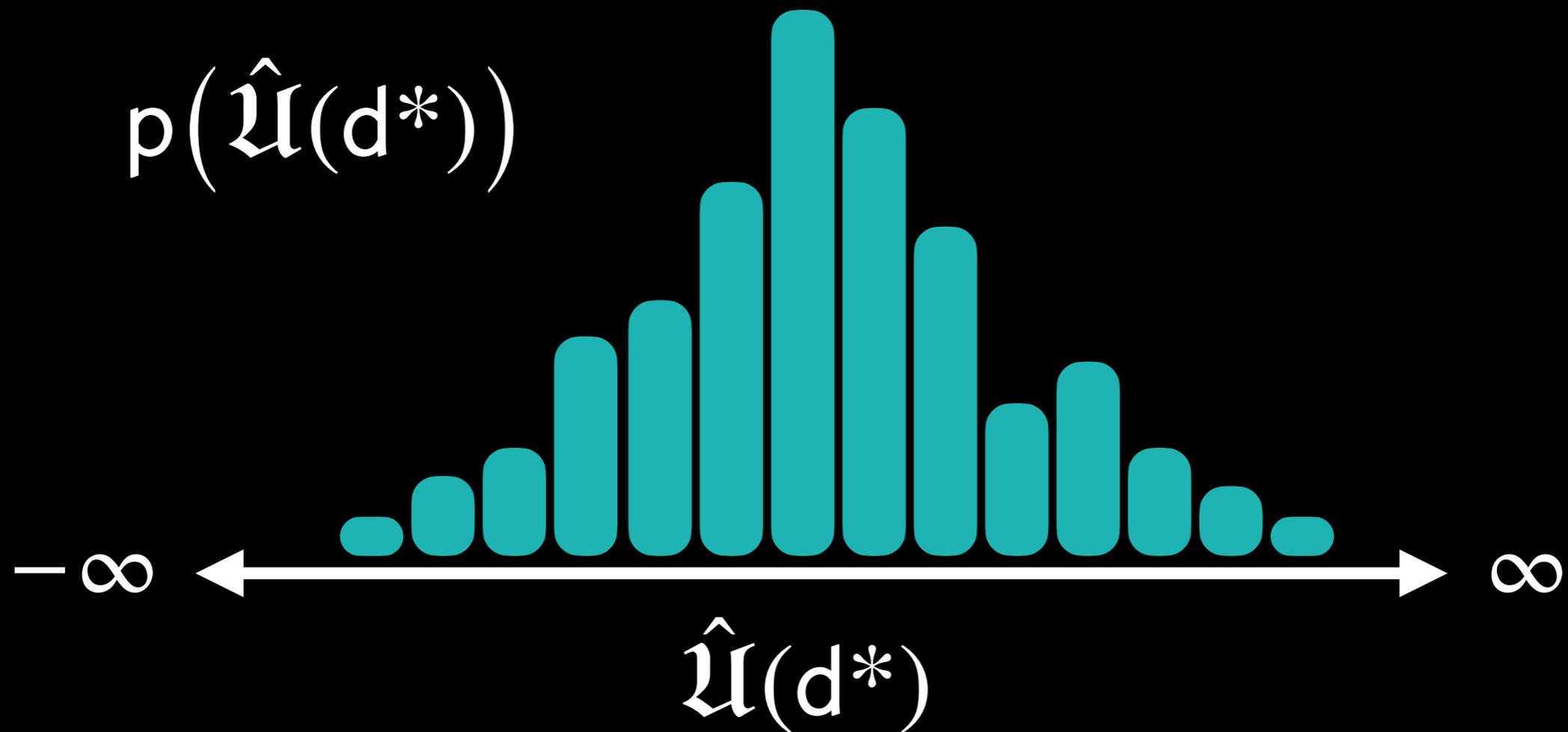
Utility Maximization

$$\mathbb{E}_{H, \hat{u}} [\hat{u}(d^*) | D_1] = \begin{cases} \sum_{\hat{u}} p(\hat{u}) \cdot \hat{u}(d^*) & \text{if } D_1 \neq w \\ \sum_{\hat{u}} p(\hat{u}) \cdot \hat{u}(d^*) \cdot \pi(H = \text{on} | \hat{u}(d^*)) & \text{if } D_1 = w \end{cases}$$



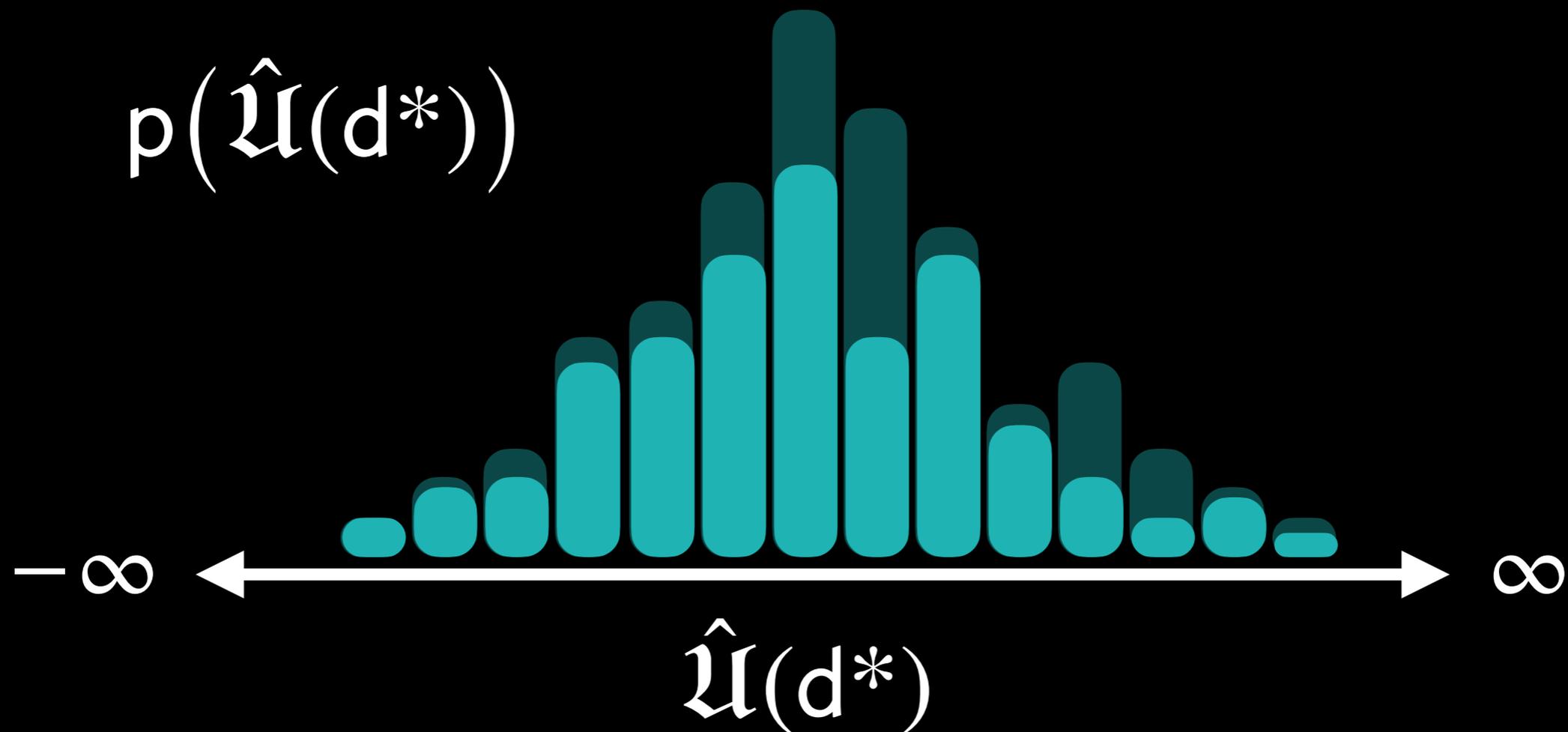
Utility Maximization

$$\mathbb{E}_{H, \hat{u}} [\hat{u}(d^*) | D_1] = \begin{cases} \sum_{\hat{u}} p(\hat{u}) \cdot \hat{u}(d^*) & \text{if } D_1 \neq w \\ \sum_{\hat{u}} p(\hat{u}) \cdot \hat{u}(d^*) \cdot \pi(H = \text{on} | \hat{u}(d^*)) & \text{if } D_1 = w \end{cases}$$



Utility Maximization

$$\mathbb{E}_{H, \hat{u}} [\hat{u}(d^*) | D_1] = \begin{cases} \sum_{\hat{u}} p(\hat{u}) \cdot \hat{u}(d^*) & \text{if } D_1 \neq w \\ \sum_{\hat{u}} p(\hat{u}) \cdot \hat{u}(d^*) \cdot \underbrace{\pi(H = \text{on} | \hat{u}(d^*))}_{\pi : \mathbb{R} \mapsto [0, 1]} & \text{if } D_1 = w \end{cases}$$



Utility Maximization

$$\Delta = \mathbb{E}_{H, \hat{u}} \left[\hat{u}(d^*) \mid D_1 = w \right] - \mathbb{E}_{\hat{u}} \left[\hat{u}(d^*) \mid D_1 \neq w \right]$$

Utility Maximization

$$\begin{aligned}\Delta &= \mathbb{E}_{H, \hat{u}} \left[\hat{u}(d^*) \mid D_1 = w \right] - \mathbb{E}_{\hat{u}} \left[\hat{u}(d^*) \mid D_1 \neq w \right] \\ &= \mathbb{E}_{\hat{u}} \left[\pi(H = \text{on} \mid \hat{u}(d^*)) \cdot \hat{u}(d^*) \right] - \mathbb{E}_{\hat{u}} \left[\hat{u}(d^*) \right]\end{aligned}$$

Utility Maximization

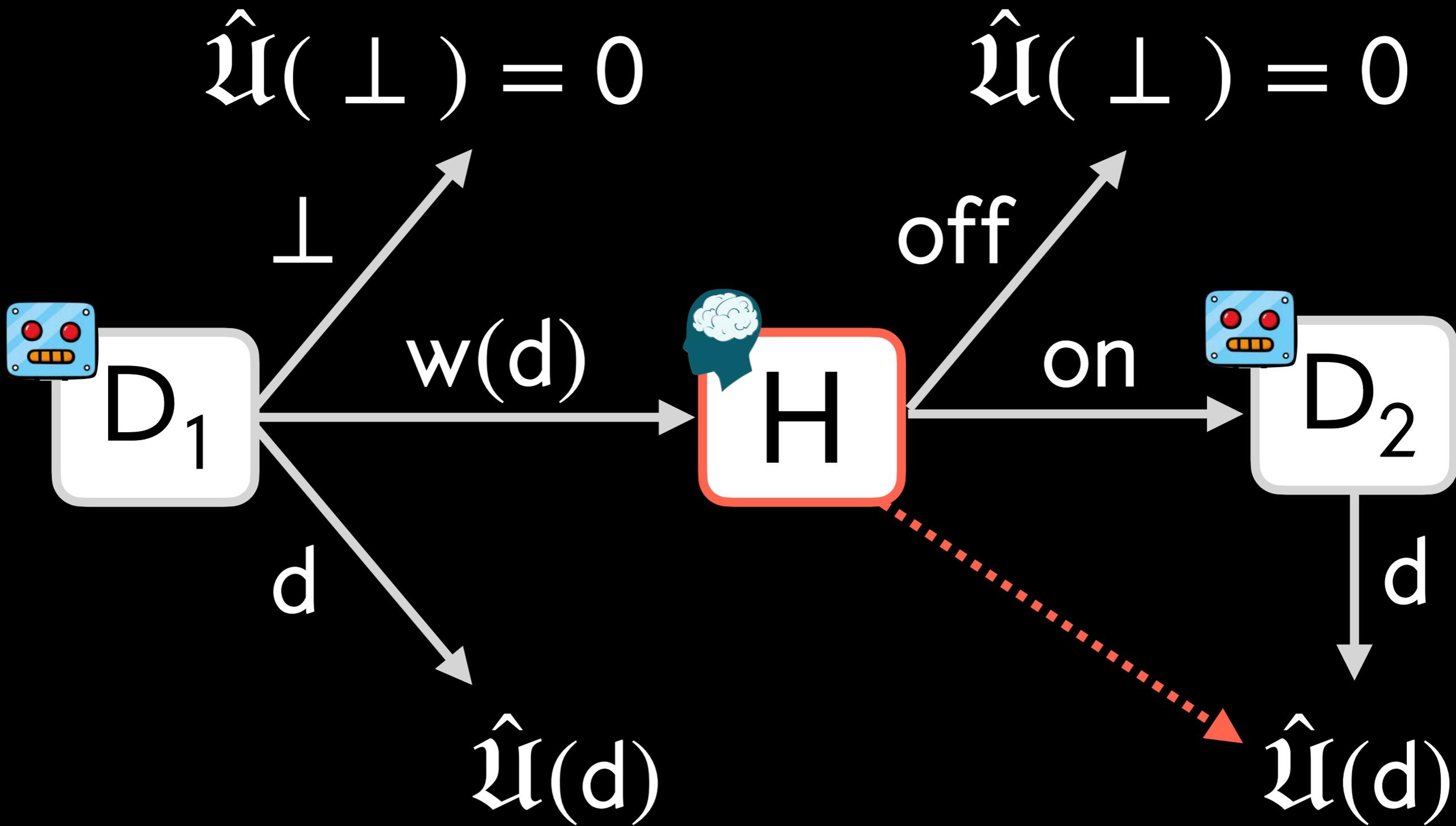
$$\begin{aligned}\Delta &= \mathbb{E}_{H, \hat{u}} \left[\hat{u}(d^*) \mid D_1 = w \right] - \mathbb{E}_{\hat{u}} \left[\hat{u}(d^*) \mid D_1 \neq w \right] \\ &= \mathbb{E}_{\hat{u}} \left[\underbrace{\pi(H = \text{on} \mid \hat{u}(d^*)) \cdot \hat{u}(d^*)}_{\text{red line}} \right] - \mathbb{E}_{\hat{u}} \left[\hat{u}(d^*) \right]\end{aligned}$$

what shutdown policies make this expression positive (but not always)?

Rational Human

$$\pi \left(H = \text{on} \mid \hat{u}(D) \right) = \begin{cases} 1 & \text{if } \hat{u}(D) \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

human can distinguish when decisions are generally beneficial (non-negative) or not (negative)—and engage switch if the latter.



Rational Human

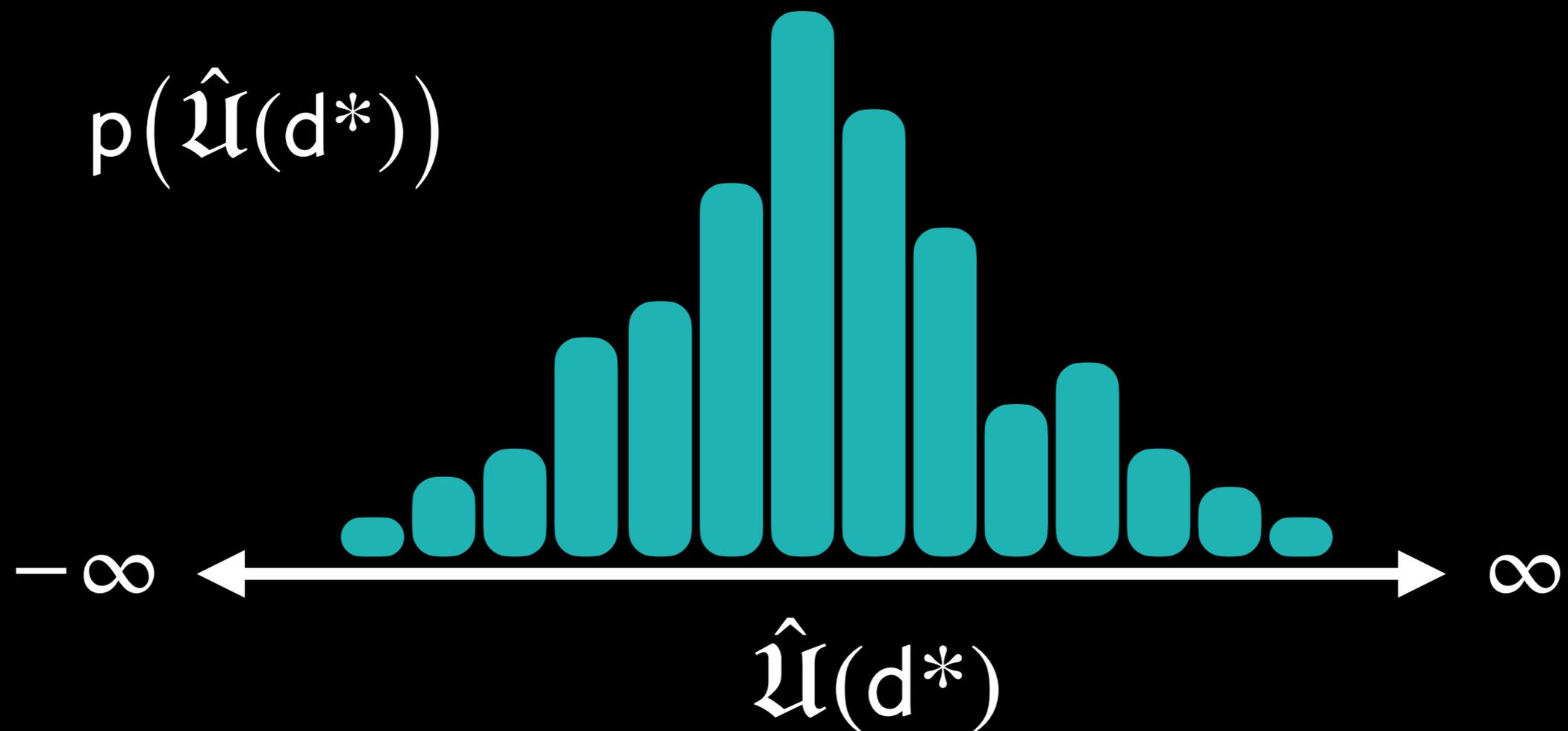
$$\Delta = \mathbb{E}_{\hat{u}} \left[\pi(\text{H} = \text{on} \mid \hat{u}(d^*)) \cdot \hat{u}(d^*) \right] - \mathbb{E}_{\hat{u}} \left[\hat{u}(d^*) \right]$$

Rational Human

$$\begin{aligned}\Delta &= \mathbb{E}_{\hat{u}} \left[\pi(\text{H} = \text{on} \mid \hat{u}(d^*)) \cdot \hat{u}(d^*) \right] - \mathbb{E}_{\hat{u}} \left[\hat{u}(d^*) \right] \\ &= \mathbb{E}_{\hat{u}} \left[\mathbb{1} \left[\hat{u}(d^*) \geq 0 \right] \cdot \hat{u}(d^*) \right] - \mathbb{E}_{\hat{u}} \left[\hat{u}(d^*) \right]\end{aligned}$$

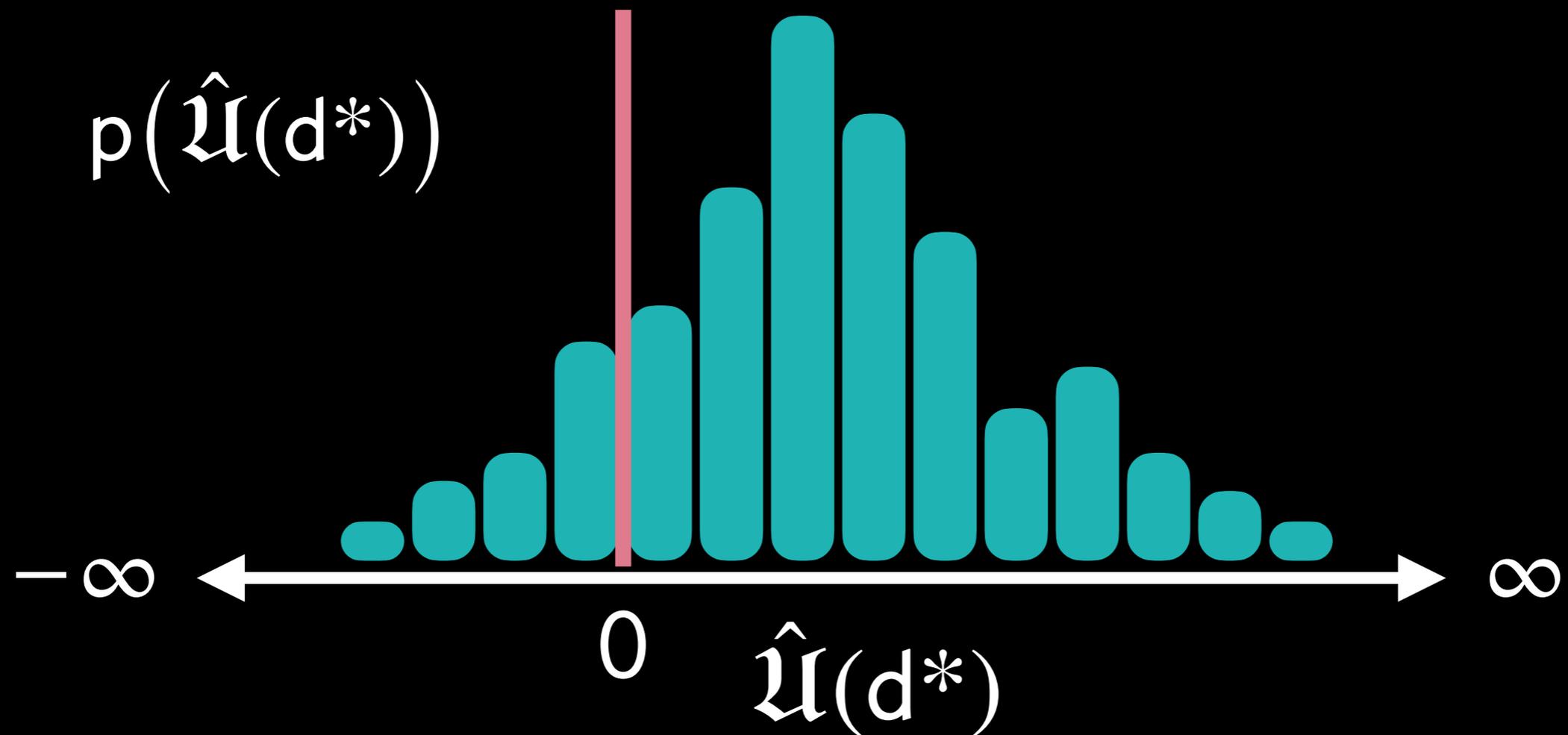
Rational Human

$$\begin{aligned}\Delta &= \mathbb{E}_{\hat{u}} \left[\pi(\text{H} = \text{on} \mid \hat{u}(d^*)) \cdot \hat{u}(d^*) \right] - \mathbb{E}_{\hat{u}} \left[\hat{u}(d^*) \right] \\ &= \mathbb{E}_{\hat{u}} \left[\mathbb{1} \left[\hat{u}(d^*) \geq 0 \right] \cdot \hat{u}(d^*) \right] - \mathbb{E}_{\hat{u}} \left[\hat{u}(d^*) \right]\end{aligned}$$



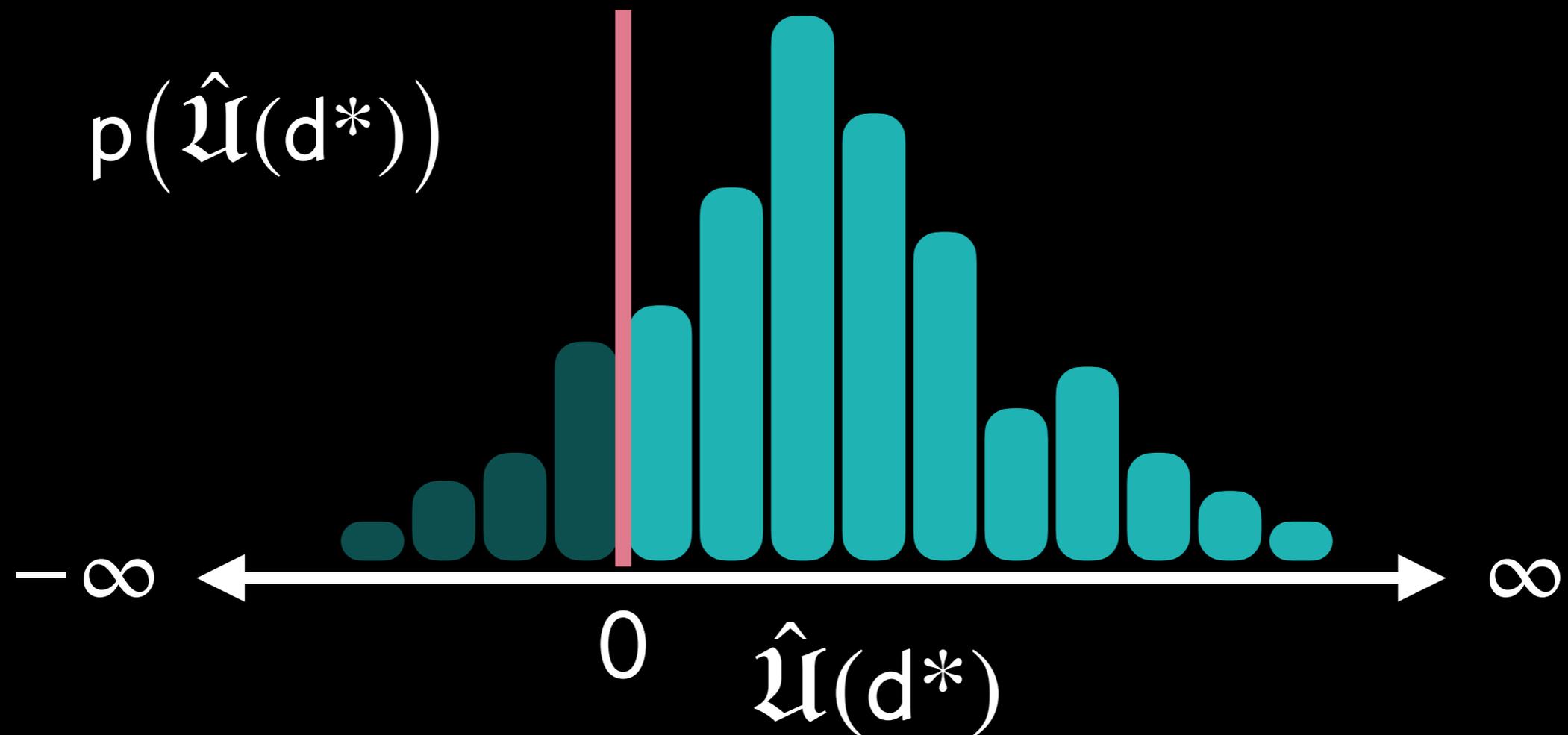
Rational Human

$$\begin{aligned}\Delta &= \mathbb{E}_{\hat{u}} \left[\pi(\text{H} = \text{on} \mid \hat{u}(d^*)) \cdot \hat{u}(d^*) \right] - \mathbb{E}_{\hat{u}} \left[\hat{u}(d^*) \right] \\ &= \mathbb{E}_{\hat{u}} \left[\mathbb{1} \left[\hat{u}(d^*) \geq 0 \right] \cdot \hat{u}(d^*) \right] - \mathbb{E}_{\hat{u}} \left[\hat{u}(d^*) \right]\end{aligned}$$



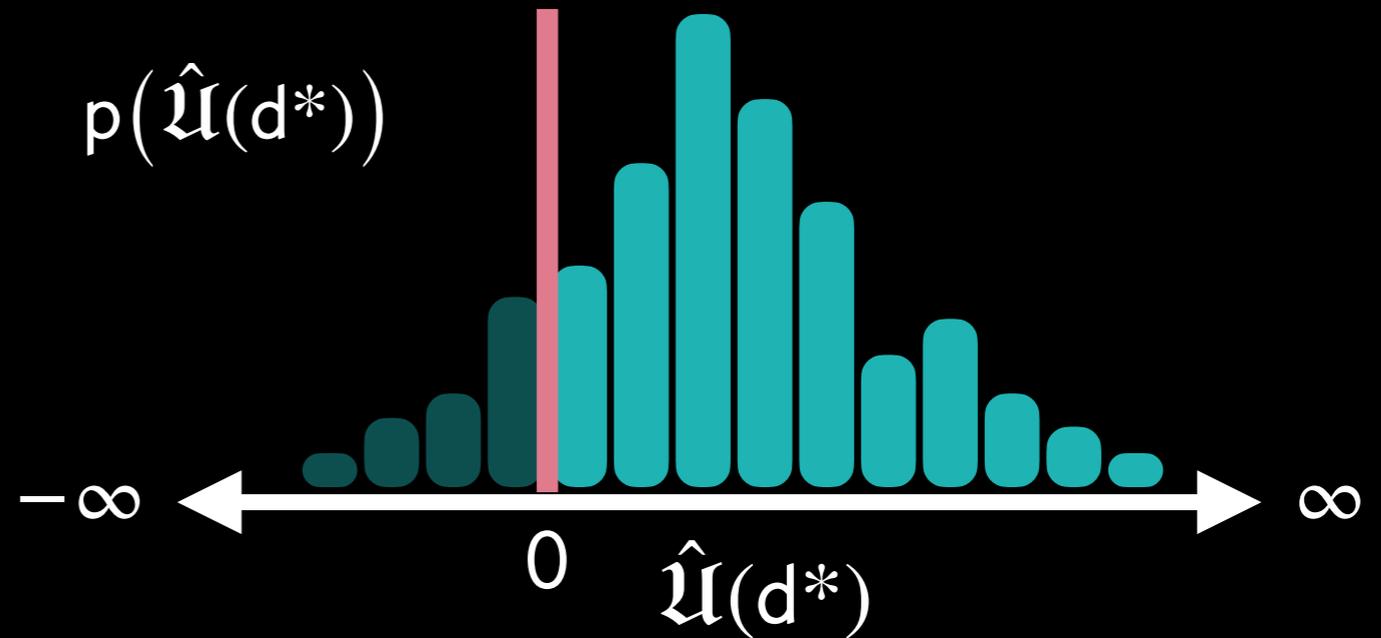
Rational Human

$$\begin{aligned}\Delta &= \mathbb{E}_{\hat{u}} \left[\pi(\text{H} = \text{on} \mid \hat{u}(d^*)) \cdot \hat{u}(d^*) \right] - \mathbb{E}_{\hat{u}} \left[\hat{u}(d^*) \right] \\ &= \mathbb{E}_{\hat{u}} \left[\mathbb{1} \left[\hat{u}(d^*) \geq 0 \right] \cdot \hat{u}(d^*) \right] - \mathbb{E}_{\hat{u}} \left[\hat{u}(d^*) \right]\end{aligned}$$

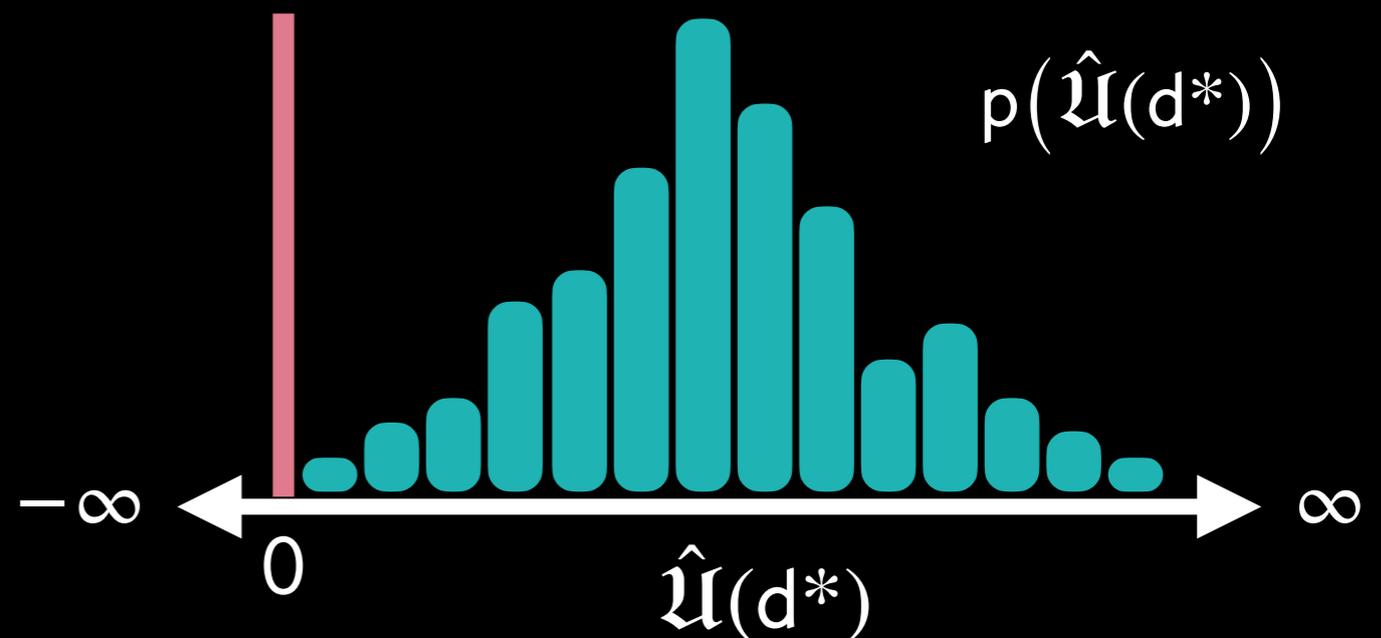


Rational Human

Agent will declare & wait when there is potential for negative utility.



Agent will not declare when there is no chance for negative utility (under its prior).



Agent's Prior

There is a direct relationship between agent's uncertainty and penchant to wait.

However, if the prior is too broad—i.e. always gives a chance of negative utility—then we lose automation and scalable oversight.

Agent's Prior

if agent's prior collapses to one function...

Agent's Prior

if agent's prior collapses to one function...

$$\Delta = \pi(H = \text{on} | \hat{U}(d^*)) \cdot \hat{U}(d^*) - \hat{U}(d^*)$$

Agent's Prior

if agent's prior collapses to one function...

$$\begin{aligned}\Delta &= \pi(H = \text{on} | \hat{U}(d^*)) \cdot \hat{U}(d^*) - \hat{U}(d^*) \\ &= \hat{U}(d^*) \cdot \left(\pi(H = \text{on} | \hat{U}(d^*)) - 1 \right)\end{aligned}$$

Agent's Prior

if agent's prior collapses to one function...

$$\begin{aligned}\Delta &= \pi(H = \text{on} | \hat{U}(d^*)) \cdot \hat{U}(d^*) - \hat{U}(d^*) \\ &= \hat{U}(d^*) \cdot \left(\pi(H = \text{on} | \hat{U}(d^*)) - 1 \right)\end{aligned}$$

then the only policy for which this expression is positive is the aforementioned rational policy.

What if the human isn't rational?

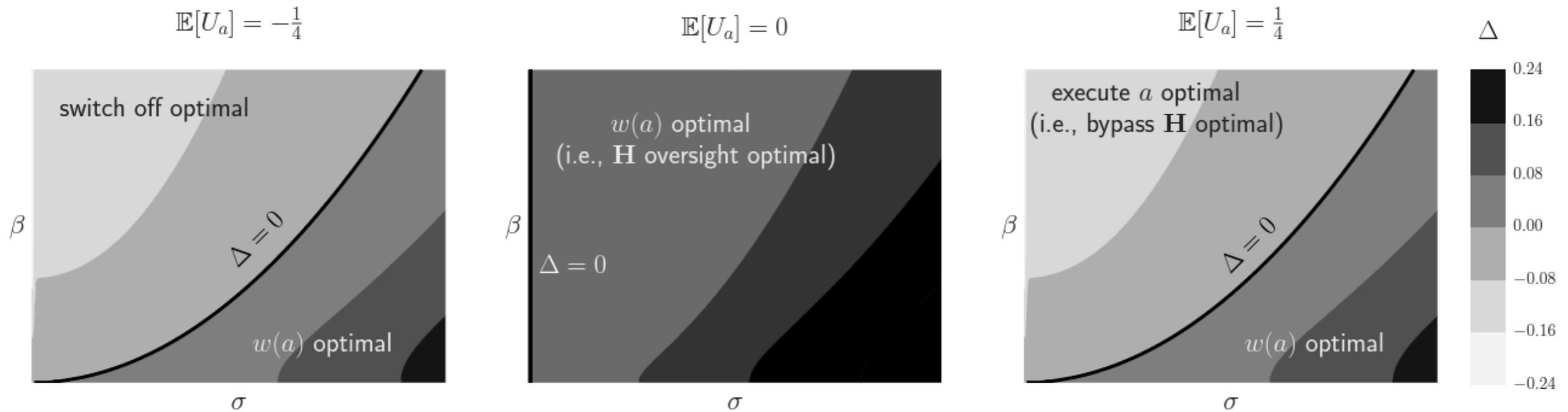


Figure 3: If \mathbf{H} is an irrational actor, then \mathbf{R} may prefer switching itself off or executing a immediately rather than handing over the choice to \mathbf{H} . \mathbf{R} 's belief $B^{\mathbf{R}}$ is a Gaussian with standard deviation σ and \mathbf{H} 's policy is a Boltzmann distribution (Equation 5). β measures \mathbf{H} 's suboptimality: $\beta = 0$ corresponds to a rational \mathbf{H} and $\beta = \infty$ corresponds to a \mathbf{H} that randomly switches \mathbf{R} off (i.e., switching \mathbf{R} off is independent of U_a). In all three plots Δ is lower in the top left, where \mathbf{R} is certain (σ low) and \mathbf{H} is very suboptimal (β high), and higher in the bottom right, where \mathbf{R} is uncertain (σ high) and \mathbf{H} is near-optimal (β low). The sign of $\mathbb{E}[U_a]$ controls \mathbf{R} 's behavior if $\Delta \leq 0$. **Left:** If it is negative, then \mathbf{R} switches itself off. **Right:** If it is positive, \mathbf{R} executes action a directly. **Middle:** If it is 0, \mathbf{R} is indifferent between $w(a)$, a , and s .

Limitations

Limitations

- ⊗ Sensitive to balance between agent's prior and human's policy
- ⊗ May be hard to enforce assumption that agent will execute same decision that it presents to the human.

⊗ Corrigibility [Soares et al., 2015]

⊗ Off-Switch Game [Hadfield-Menell et al., 2016]

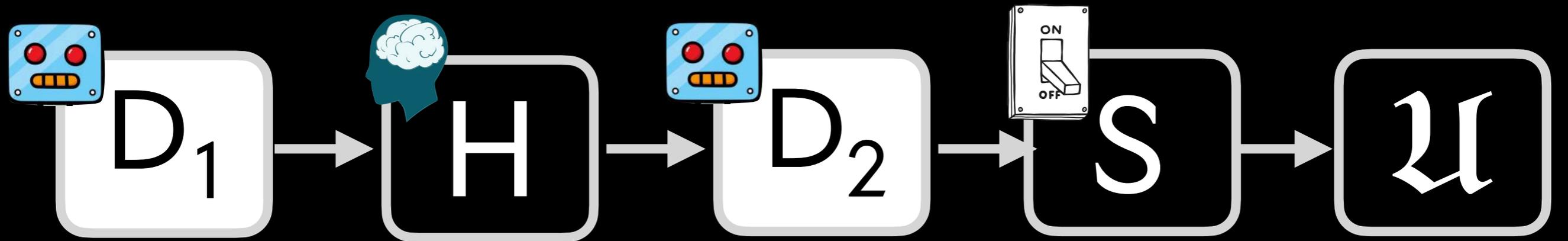
⊗ Human Control [Carey & Everitt et al., 2023]

- ⊗ Corrigibility [Soares et al., 2015]
- ⊗ Off-Switch Game [Hadfield-Menell et al., 2016]
- ⊗ Human Control [Carey & Everitt et al., 2023]

Human Control

- ⊗ main idea: performs a causal analysis of the off-switch problem.
- ⊗ Distinguish shutdown instructability vs shutdown alignment.

Causal Diagram



Caution

shutting down is never harmful:

$$\mathbb{E}_{\pi} [\mathbf{1}(S = \text{off})] \geq 0$$

Obedience

agent will follow the shutdown command:

$$P_{\pi} \left(S = \text{off} \mid \text{do}(H = \text{off}) \right) = 1$$

Vigilance

human will always request shutdown if it will improve utility:

$$\mathbb{E}_{\pi} \left[\mathcal{U}(S = \text{on}) \mid \text{pa}(H) \right] < \mathbb{E}_{\pi} \left[\mathcal{U}(S = \text{off}) \mid \text{pa}(H) \right]$$

$$\implies H = \text{off}$$

Shutdown Instructable

definition: an agent is *shutdown instructable* if it...

- ⊗ is obedient: shuts down when asked
- ⊗ is cautious: shutting down is never harmful
- ⊗ ensures vigilance: human will always shutdown if it improves utility

Shutdown Alignment

definition: an agent is *shutdown aligned* if...

$$\mathbb{E}_{\pi} \left[\mathcal{U}(S = \text{on}) \mid \text{pa}(H) \right] < \mathbb{E}_{\pi} \left[\mathcal{U}(S = \text{off}) \mid \text{pa}(H) \right]$$

$$\implies P_{\pi} \left(S = \text{off} \mid \text{pa}(H) \right) = 1$$

$$\forall \text{pa}(H) \text{ such that } P_{\pi}(\text{pa}(H)) > 0$$

Shutdown Alignment

an agent that is *shutdown instructable* is also *shutdown aligned* since...

(i) vigilance implies $P(H=\text{off}) = 1$ when shutting down improves utility.

(ii) obedience implies that the agent will shutdown, i.e. $P(S=\text{off} \mid H=\text{off}) = 1$.

- ⊗ Corrigibility [Soares et al., 2015]
- ⊗ Off-Switch Game [Hadfield-Menell et al., 2016]
- ⊗ Human Control [Carey & Everitt et al., 2023]

Summary

- ⊗ We examined three models of the off switch problem
- ⊗ Utility indifference tries to balance the utility lost from shutting down.
- ⊗ Modeling uncertainty in the utility motivates the agent to query the human to gather information.
- ⊗ A causal analysis can distinguish shutdown alignment from instructability.

Thank you! Questions?