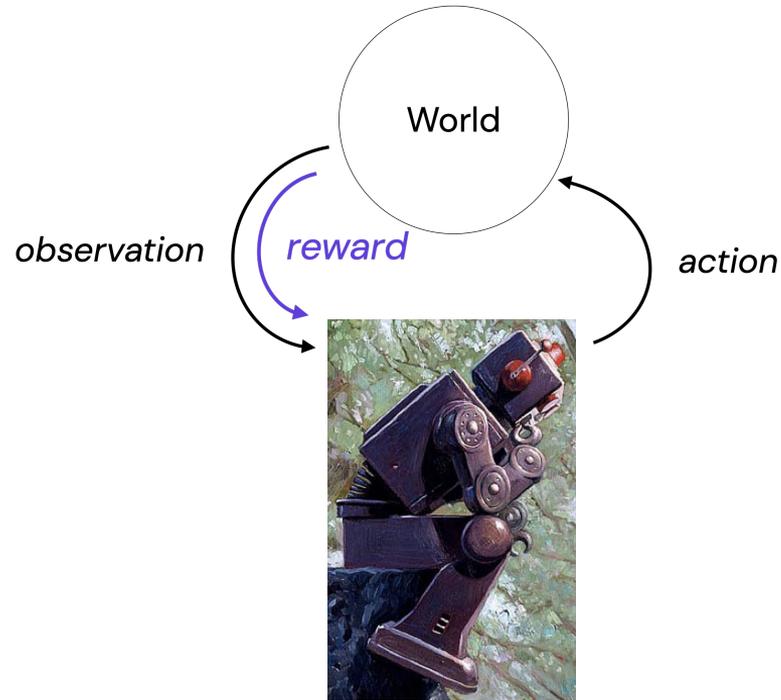


On Tasks and Rewards



Overview



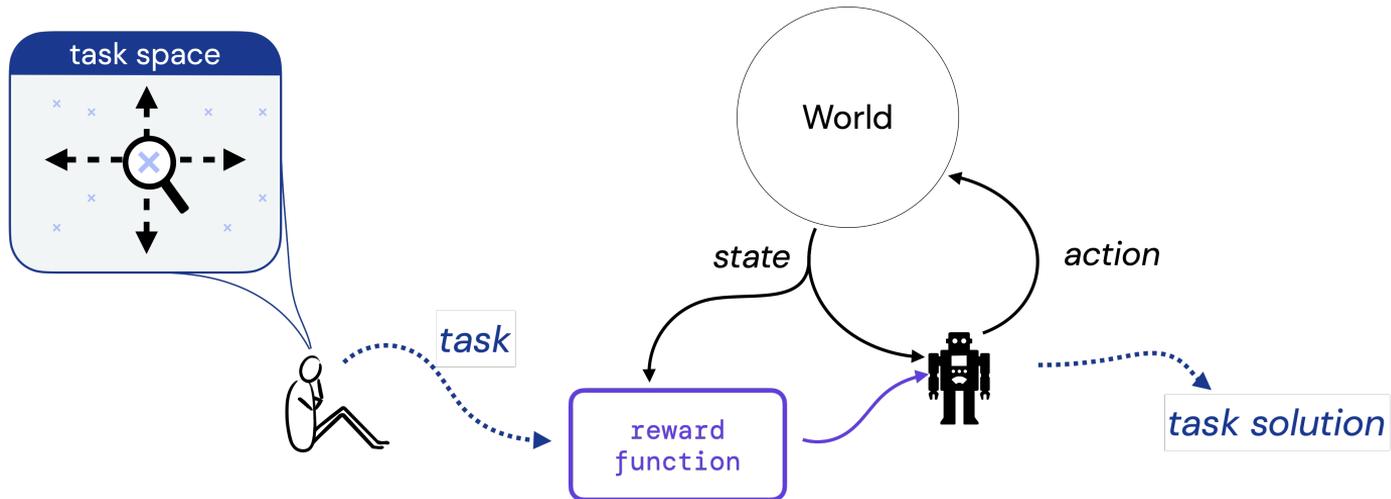
“Part of the appeal of reinforcement learning is that it is in a sense the whole AI problem in a microcosm.”

– [Sutton, 1992](#)

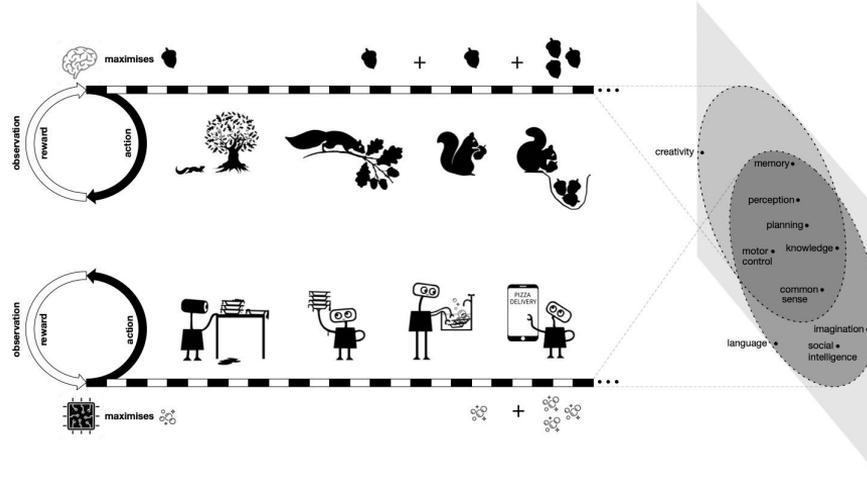
The Reward Hypothesis

“...all of what we mean by goals and purposes can be well thought of as maximization of the expected value of the cumulative sum of a received scalar signal (reward)”

-- [Sutton \(2004\)](#), [Littman \(2017\)](#)



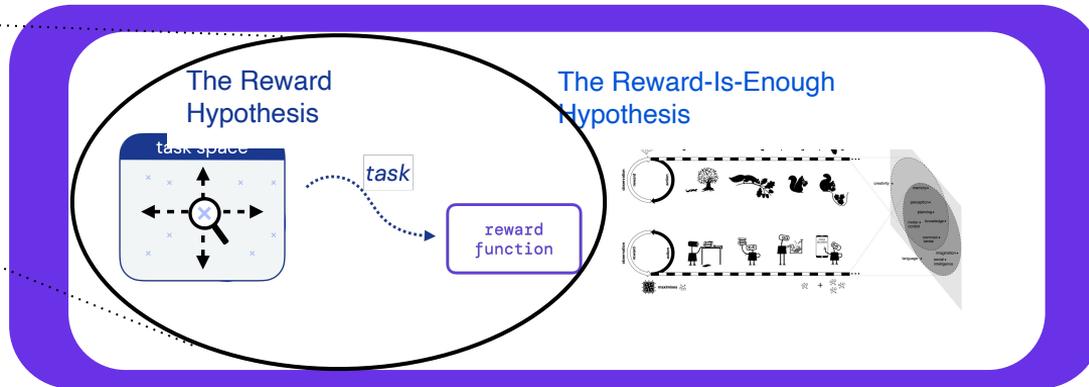
Reward is Enough



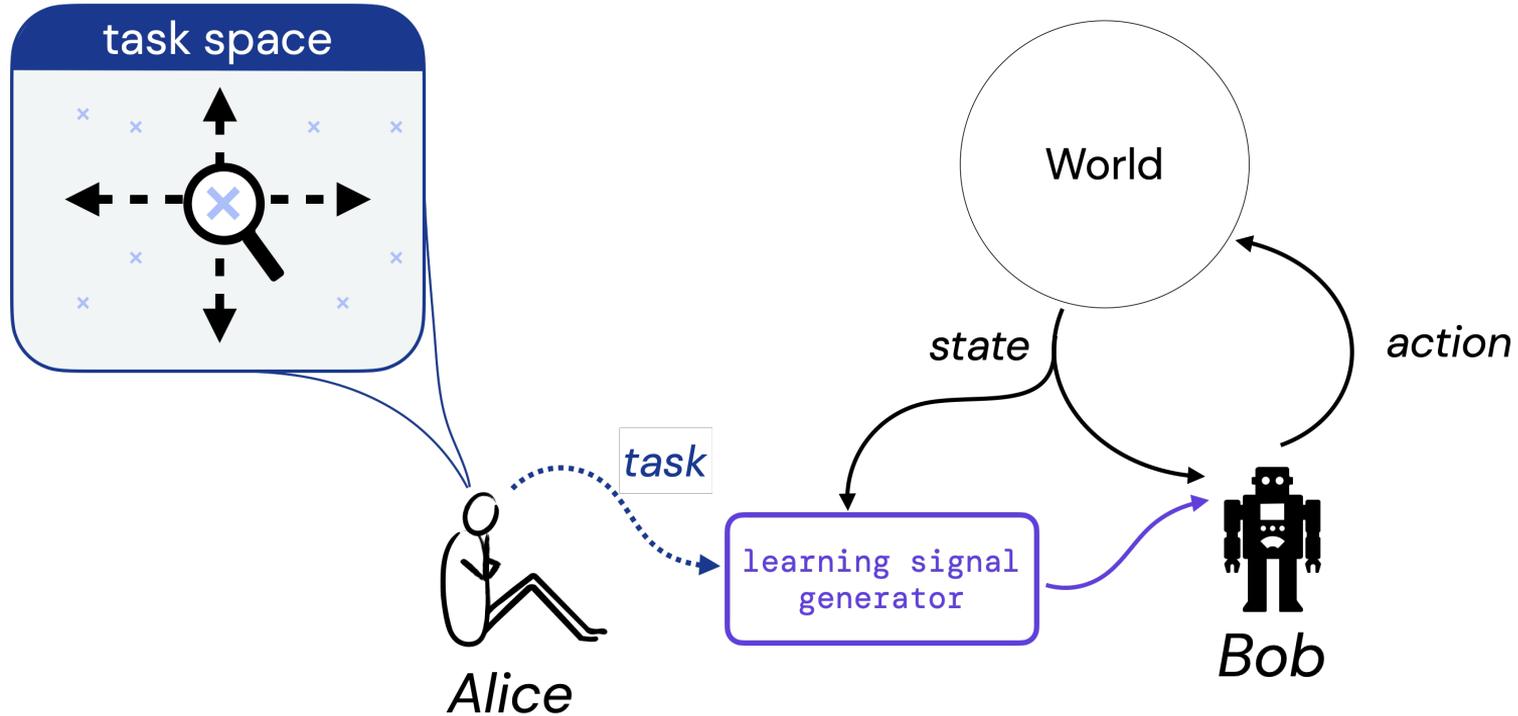
“Intelligence, and its associated abilities, can be understood as subserving the maximisation of reward by an agent acting in its environment”

-- [Silver, Singh, Precup, Sutton \(2021\)](#)

On the Expressivity of Markov Reward (NeurIPS'2021)

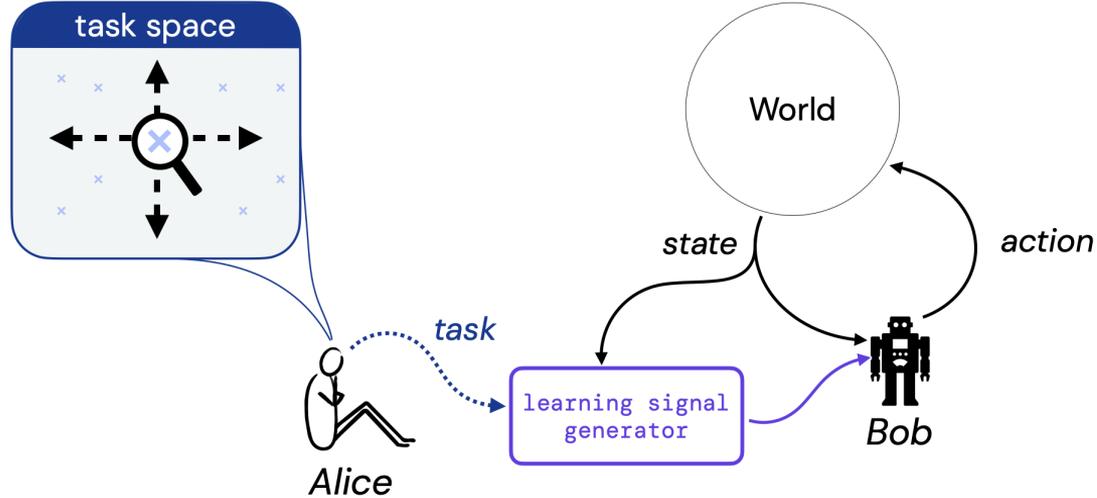


Formalizing the Reward Hypothesis



The Two Question View

Expression Question: Which signal can be used as a mechanism for expressing a given task?

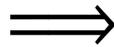


The Two Question View

Expression Question: Which signal can be used as a mechanism for expressing a given task?



$\mathcal{T} = ?$



The Reward Hypothesis (formalized)

Given *any* task \mathcal{T} and *any* environment E there is a reward function that realizes \mathcal{T} in E

Task Question: What is a task?

The Two Question View

Expression Question: Which signal can be used as a mechanism for expressing a given task?

$\mathcal{T} = ?$



The Reward Hypothesis (formalized)

Given *any* task \mathcal{T} and *any* environment E there is a reward function that realizes \mathcal{T} in E

Task Question: What is a task?

Assumption. All environments are finite Controlled Markov Processes (CMPs).

$$E = (\mathcal{S}, \mathcal{A}, T, \gamma, s_0)$$

$R(s), R(s, a), R(s, a, s'), R(s')$

Related Work: Other Perspectives on Reward

Safety

[Everitt et al. 2017](#), [Ortega et al. 2018](#), [Kumar et al. 2020](#)
[Uesato et al. 2020](#)

Preferences

[MacGlashan et al. 2016](#), [Wirth et al. 2017](#), [Christiano et al. 2017](#), [Xu et al. 2020](#) ★

Reward Learning & Design

[Ackley & Littman 1992](#), [Singh et al. 2010](#), [Sorg 2011](#), [Zheng et al. 2020](#), [Jeon et al. 2020](#) ★

Constrained MDPs

★
[Mannor & Shimkin 2004](#), [Szepesvári 2020](#), [Roijers et al. 2020](#), [Zahavy et al. 2021](#)

Teaching

[Goldman & Kearns 1995](#), [Simard et al. 2017](#), [Ho et al. 2019](#)

Logical tasks in RL

[Littman et al. 2017](#), [Li et al. 2017](#), [Jothimurugan et al. 2020](#), [Tasse et al. 2020](#)

Expectations, Discount, and Rationality

[Mitten 1974](#), [Sobel 1975](#), [Weng 2011](#), [Pitis 2019](#), [Gottipati et al. 2020](#) ★

Target Distribution

[Akshay et al. 2013](#), [Hafner et al. 2020](#)

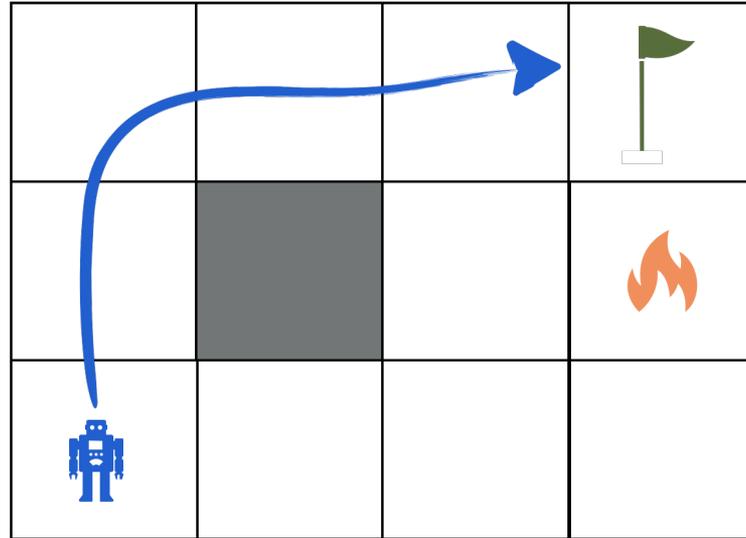
IRL, CIRL, Assistive Learning

[Syed et al. 2008](#), [Hadfield-Menell et al. 2016](#), [Amin et al. 2017](#), [Shah et al. 2020](#)

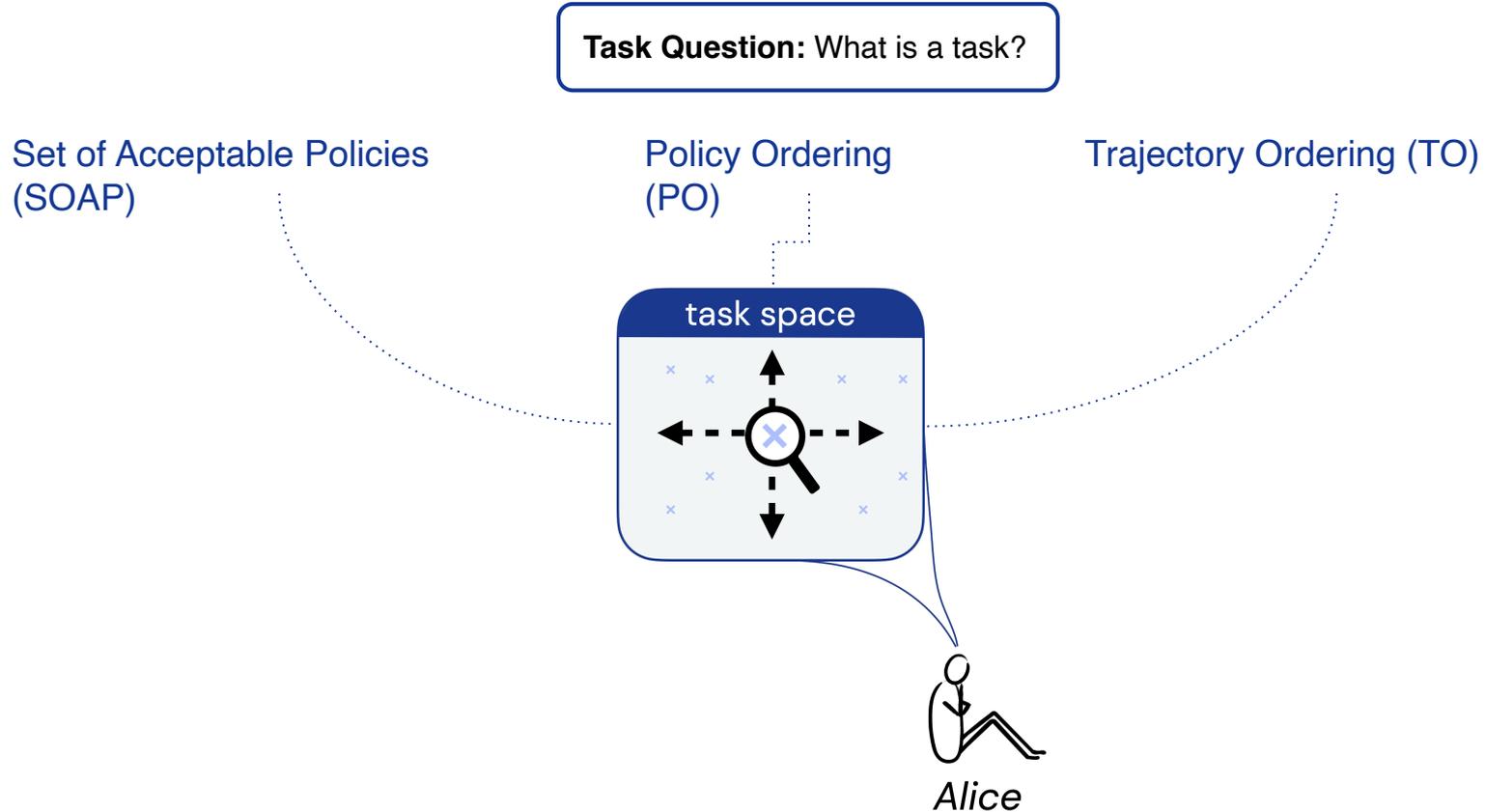
Natural Language

[MacGlashan et al. 2015](#), [Williams et al. 2017](#)

What is a Task?



Task Types: SOAPs, POs, TOs



Task Types: SOAPs, POs, TOs

Task Question: What is a task?

Set of Acceptable Policies
(SOAP)

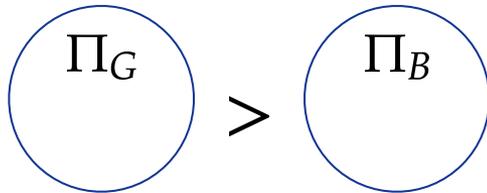
$$\Pi_G \subseteq \Pi$$

Policy Ordering
(PO)

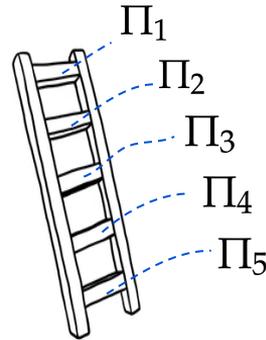
$$L_\Pi$$

Trajectory Ordering (TO)

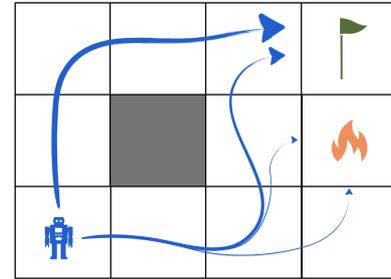
$$L_{\tau, N}$$



Example:
“Reach the goal in less than 10
steps in expectation.”

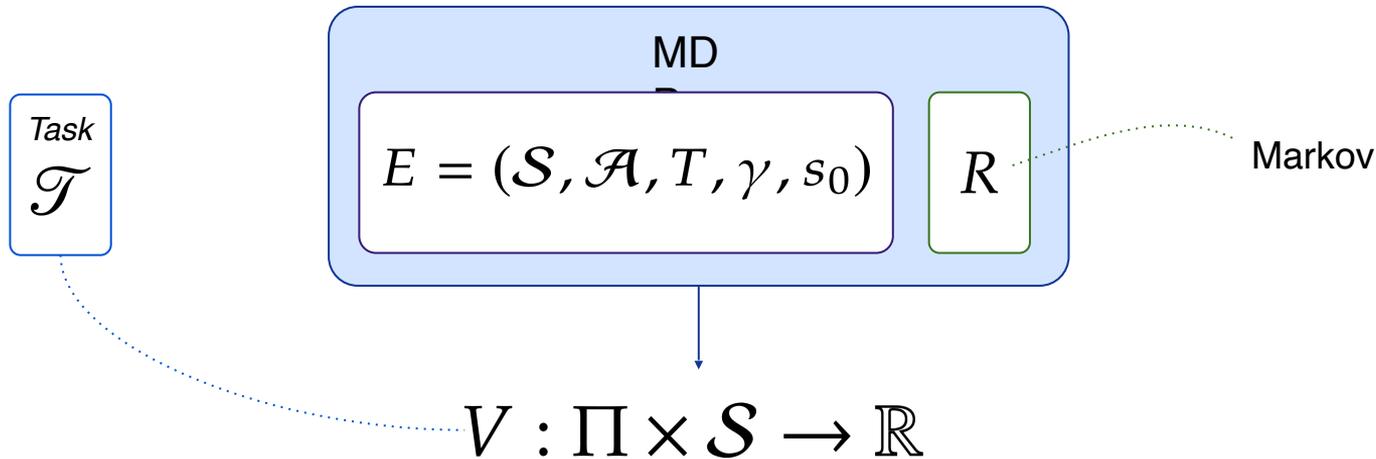
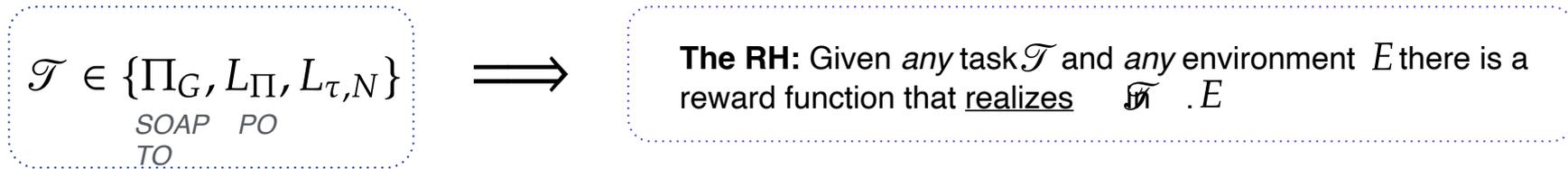


Example:
“I prefer you reach the goal in 5 steps,
else within 10, else don't bother.”



Example:
I prefer safely reaching the goal
and avoid lava at all costs.

Task Realization



Task Realization

$$\mathcal{T} \in \{\underbrace{\Pi_G}_{\text{SOAP}}, \underbrace{L_\Pi}_{\text{PO}}, \underbrace{L_{\tau,N}}_{\text{TO}}\}$$



The RH: Given any task \mathcal{T} and any environment E there is a reward function that realizes \mathcal{T} in E .

Set of Acceptable Policies (SOAP)

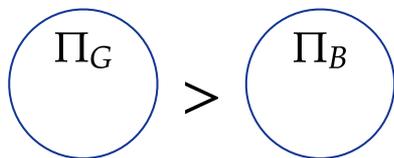
$$\Pi_G \subseteq \Pi$$

Policy Ordering (PO)

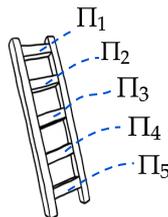
$$L_\Pi$$

Trajectory Ordering (TO)

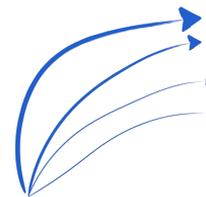
$$L_{\tau,N}$$



$$V^{\pi_g}(s_0) > V^{\pi_b}(s_0)$$

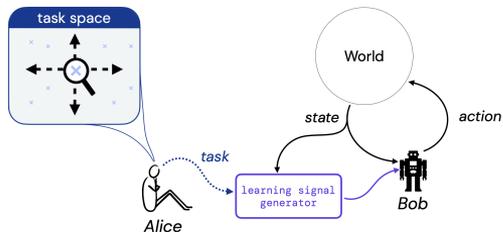


$$V^{\pi_1}(s_0) > V^{\pi_2}(s_0) \dots$$



$$G(\tau_1; s_0) > G(\tau_2; s_0) \dots$$

Recap



Which signal can be used to express any task?

The RH: Reward

What is a task?

$\mathcal{T} \in \{\Pi_G, L_{\Pi}, L_{\tau, N}\}$
SOAP PO
TO

MAIN QUESTION

Given *any* task \mathcal{T} and *any* environment $E = (\mathcal{S}, \mathcal{A}, T, \gamma, s_0)$,
is there a Markov reward function that realizes \mathcal{T} in E

Question 1: What Can Reward Express?

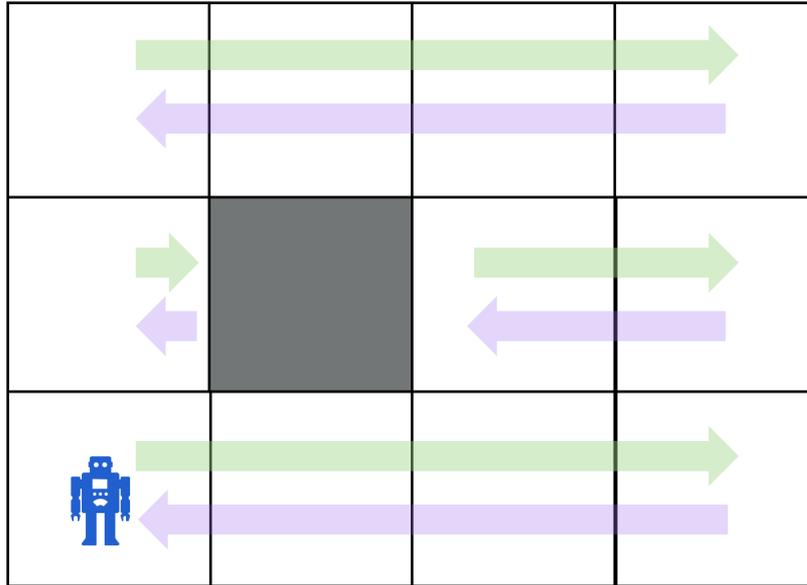
Theorem 1. *For each of SOAP, PO, and TO, there exist (E, \mathcal{T}) pairs for which no reward function realizes \mathcal{T} in E .*

MAIN QUESTION

Given *any* task \mathcal{T} and *any* environment $E = (\mathcal{S}, \mathcal{A}, T, \gamma, s_0)$,
is there a Markov reward function that realizes \mathcal{T} in E ?

Expressivity Example 1

What kinds of SOAPs are not expressible?

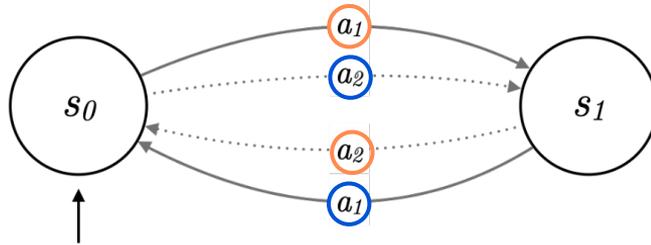


$$\Pi_G = \{\pi_{\leftarrow}, \pi_{\rightarrow}, \dots\}$$

SOAP = “Always go in the same direction”

Expressivity Example 2

What kinds of SOAPs are not expressible?



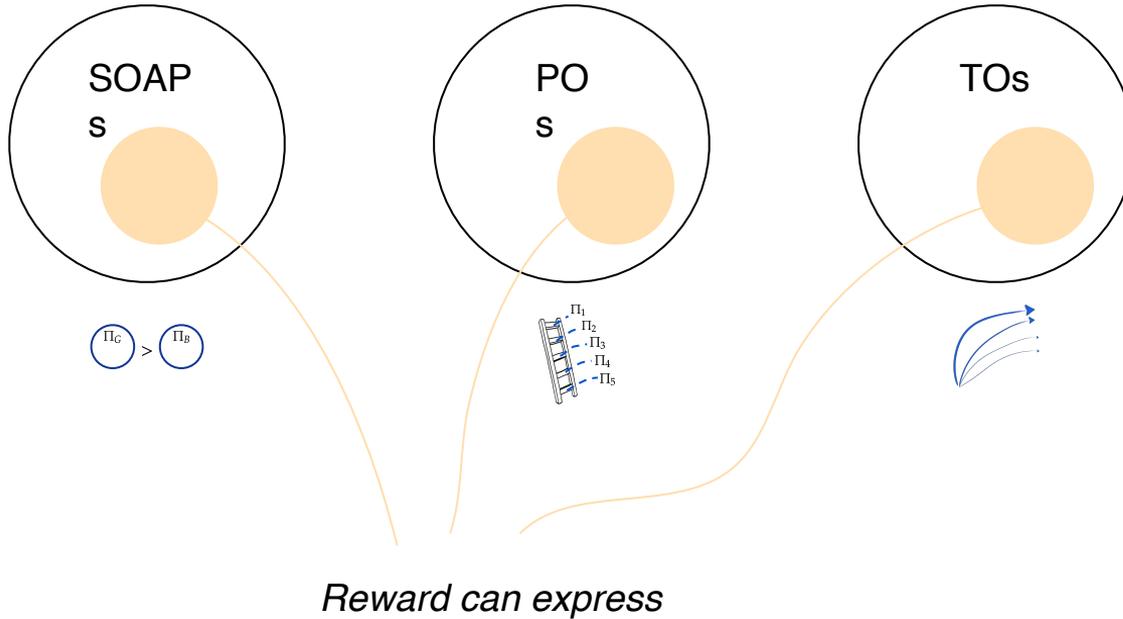
$$\Pi_G = \{\pi_{21}, \pi_{12}\}$$

XOR Problem

...Other types?

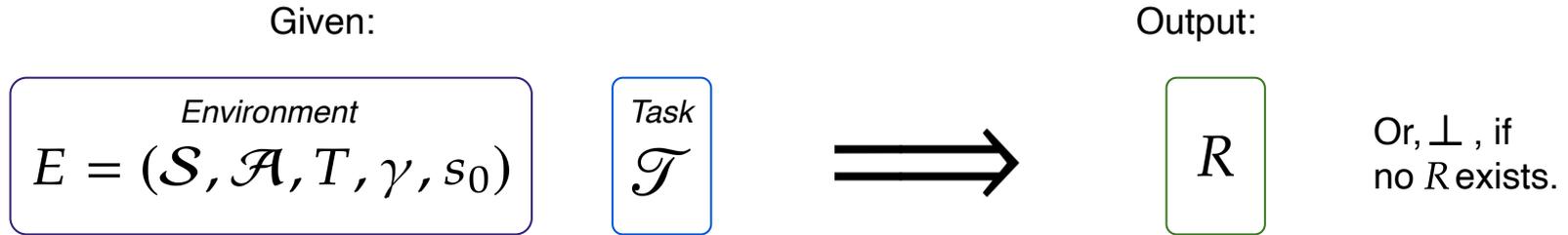
Question 2: Can We Find the Realizing Rewards?

Definition 1. The *REWARDDESIGN* problem is: **Given** $E = (\mathcal{S}, \mathcal{A}, T, \gamma, s_0)$, and a \mathcal{T} , **output** a reward function R_{alice} that ensures \mathcal{T} is realized in $M = (E, R_{alice})$.



Main Result 2: Reward Design

Definition 1. The REWARDDESIGN problem is: **Given** $E = (\mathcal{S}, \mathcal{A}, T, \gamma, s_0)$, and a \mathcal{T} , **output** a reward function R_{alice} that ensures \mathcal{T} is realized in $M = (E, R_{alice})$.



Theorem 2. The REWARDDESIGN problem can be solved in polynomial time, for any finite E , and any \mathcal{T} .

Corollary 1. Given \mathcal{T} and E , deciding whether \mathcal{T} is expressible in E is solvable in polynomial time for any finite E .

Algorithm: SOAP Reward Design

Algorithm 1 SOAP Reward Design

INPUT: $E = (\mathcal{S}, \mathcal{A}, T, \gamma, s_0), \Pi_G$.

OUTPUT: R , or \perp .

```
1:  $\Pi_{\text{fringe}} = \text{compute\_fringe}(\Pi_G)$ 
2: for  $\pi_{g,i} \in \Pi_G$  do                                ▶ Compute state-visitation distributions.
3:    $\rho_{g,i} = \text{compute\_exp\_visit}(\pi_{g,i}, E)$ 

4: for  $\pi_{f,i} \in \Pi_{\text{fringe}}$  do
5:    $\rho_{f,i} = \text{compute\_exp\_visit}(\pi_{f,i}, E)$ 

6:  $C_{\text{eq}} = \{\}$                                           ▶ Make Equality Constraints.
7: for  $\pi_{g,i} \in \Pi_G$  do
8:    $C_{\text{eq}}.\text{add}(\rho_{g,0}(s_0) \cdot X = \rho_{g,i}(s_0) \cdot X)$ 

9:  $C_{\text{ineq}} = \{\}$                                         ▶ Make Inequality Constraints.
10: for  $\pi_{f,j} \in \Pi_{\text{fringe}}$  do
11:    $C_{\text{ineq}}.\text{add}(\rho_{f,j}(s_0) \cdot X + \epsilon \leq \rho_{g,0}(s_0) \cdot X)$ 

12:  $R_{\text{out}}, \epsilon_{\text{out}} = \text{linear\_programming}(\text{obj.} = \max \epsilon, \text{constraints} = C_{\text{ineq}}, C_{\text{eq}})$     ▶ Solve LP.

13: if  $\epsilon_{\text{out}} > 0$  then                                ▶ Check if successful.
    return  $R_{\text{out}}$ 
14: else
    return  $\perp$ 
```

Recap

MAIN QUESTION

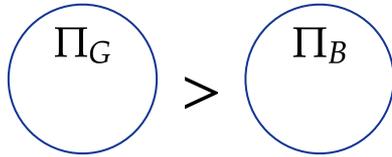
Given *any* task \mathcal{T} and *any* environment $E = (\mathcal{S}, \mathcal{A}, T, \gamma, s_0)$,
is there a Markov reward function that realizes \mathcal{T} in E

Theorem 1. *For each of SOAP, PO, and TO, there exist (E, \mathcal{T}) pairs for which no reward function realizes \mathcal{T} in E .*

Theorem 2. *The REWARDDESIGN problem can be solved in polynomial time, for any finite E , and any \mathcal{T} .*

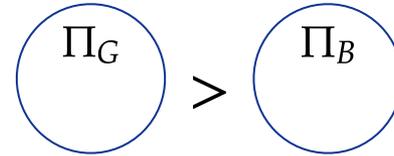
Other Analysis: Two Kinds of SOAP

“range” SOAP



$$V^{\pi_g}(s_0) > V^{\pi_b}(s_0)$$

“equal” SOAP



$$V^{\pi_g}(s_0) = V^{\pi'_g}(s_0)$$

AND

$$V^{\pi_g}(s_0) > V^{\pi_b}(s_0)$$

Proposition 1. *The “range” realization of SOAP is strictly more general than the “equal” realization.*

Other Analysis

Extensions of Main Results

Theorem 3. *There exist choices of $E_{\neg} = (\mathcal{S}, \mathcal{A}, s_0)$ and \mathcal{T} , such that there is no (T, R, γ) that realizes \mathcal{T} in E_{\neg} .*

Theorem 4. *The FINITE-REWARDDECISION problem is NP-hard.*

Multi-Environment

Theorem 5. *Given a task \mathcal{T} and a finite set of CMPs, $\mathcal{E} = \{E_1, \dots, E_n\}$, with shared state-action space, there exists a polynomial time algorithm that outputs one reward function that realizes the task (when possible) in all CMPs in \mathcal{E} .*

Theorem 6. *Task realization is not closed under sets of CMPs. That is, there exist choices of \mathcal{T} and $\mathcal{E} = \{E_1, \dots, E_n\}$ such that \mathcal{T} is realizable in each $E_i \in \mathcal{E}$ independently, but there is not a single reward function that realizes \mathcal{T} in all $E_i \in \mathcal{E}$ simultaneously.*

Limitations & Assumptions

Environment.

- > Finite CMPs.
- > γ is part of E .

Task.

- > Tasks of interest are SOAPs, POs, and TOs.

Task Realization.

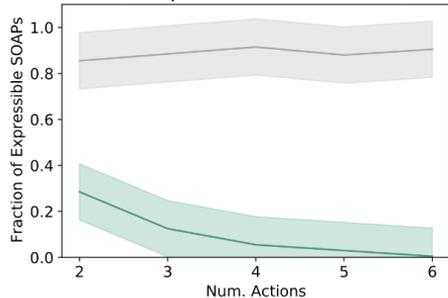
- > Start-state value determines task realization.
- > Ignore learning dynamics.

Reward Functions.

- > Deterministic.
- > Markov.

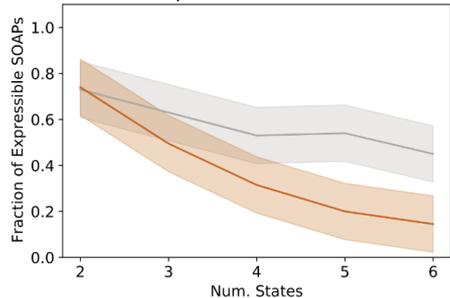
Experiment 1: SOAP Expressivity

Estimate of Expressible SOAPs vs. Num. Actions



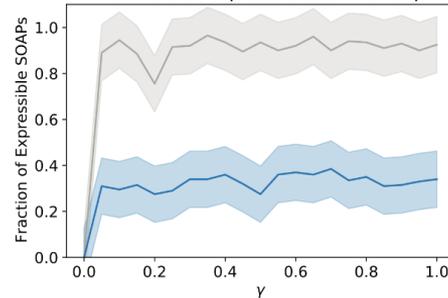
(a) Vary Num. Actions

Estimate of Expressible SOAPs vs. Num. States



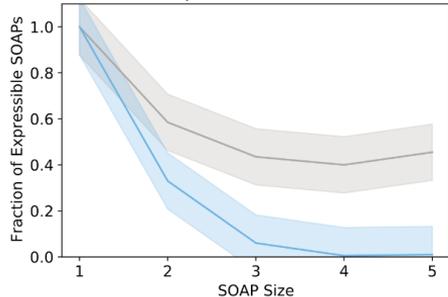
(b) Vary Num. States

Estimate of Expressible SOAPs vs. γ



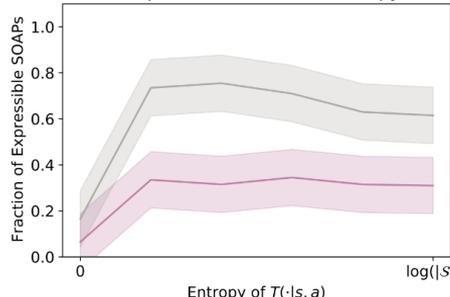
(c) Vary γ

Estimate of Expressible SOAPs vs. SOAP Size



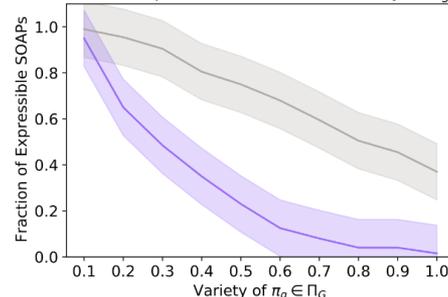
(d) Vary SOAP Size

Estimate of Expressible SOAPs vs. Entropy of $T(\cdot|s, a)$



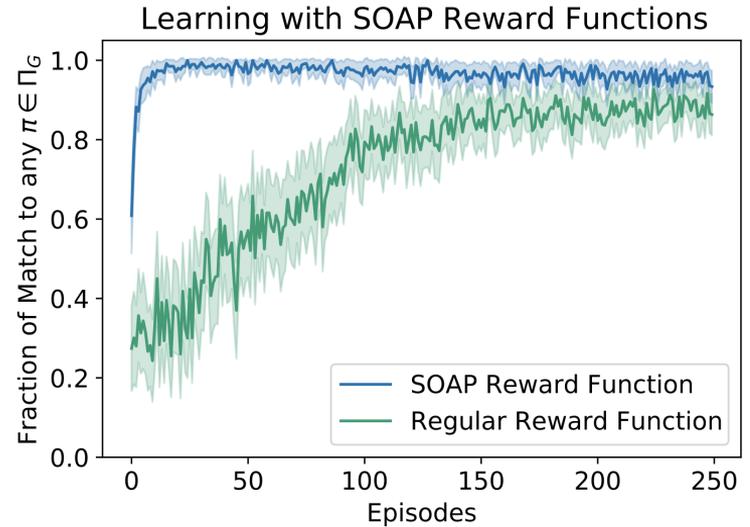
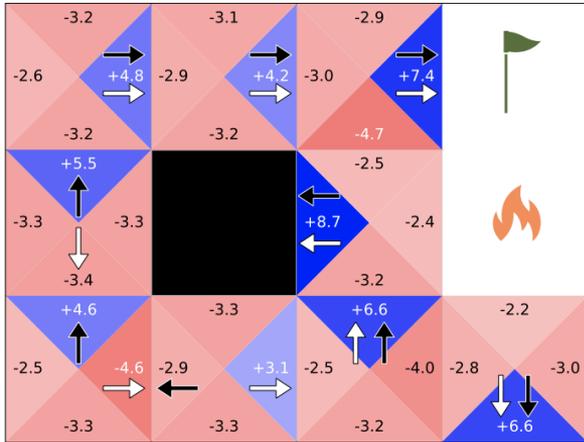
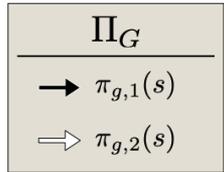
(e) Vary Entropy of T

Estimate of Expressible SOAPs vs. Variety of $\pi_a \in \Pi_G$



(f) Vary the Spread of Π_G

Experiment 3: Learning with SOAP Rewards (Grid)



Main Result Overview

