Modeling Human Beliefs about Al Behavior for Scalable Oversight

Leon Lang @uva.nl
University of Amsterdam

Patrick Forré
University of Amsterdam

p.d. for re@uva.nl

Abstract

Contemporary work in AI alignment often relies on human feedback to teach AI systems human preferences and values. Yet as AI systems grow more capable, human feedback becomes increasingly unreliable. This raises the problem of scalable oversight: How can we supervise AI systems that exceed human capabilities? In this work, we propose to model the human evaluator's beliefs about the AI system's behavior to better interpret the human's feedback. We formalize human belief models and theoretically analyze their role in inferring human values. We then characterize the remaining ambiguity in this inference and conditions for which the ambiguity disappears. To mitigate reliance on exact belief models, we then introduce the relaxation of human belief model covering. Finally, we propose using foundation models to construct covering belief models, providing a new potential approach to scalable oversight.

Contents

1	Inti	roduction	2
2	Human belief models		3
	2.1	Conventions and preliminaries for linear algebra	4
	2.2	Preliminaries on Markov Decision Processes	5
	2.3	The human's ontology and reward object	5
	2.4	The human's feature belief and observation return function	6
	2.5	The definition of human belief models	7
	2.6	Complete belief models and the ambiguity	9
	2.7	Faithful belief models	12
	2.8	Conceptual examples	13
3	Hu	man belief model covering	15
	3.1	Human belief model covering and its implications	16
	3.2	Morphisms of human belief models and ontology translations	17
	3.3	An example of symmetry-invariant features and reward functions	19
	3.4	A proposal for belief model covering in practice	23

4	Discussion					
	4.1	Summary	26			
	4.2	Related work	26			
	4.3	Future work	28			
	4.4	Conclusion	29			
Re	References					
A Preliminary results on linear algebra						
B Balanced human belief models and choices						
4.4 Conclusion References A Preliminary results on linear algebra B Balanced human belief models and choices C A diagram in the category of human belief models D Details on the example with invariant features			43			
D	Det	ails on the example with invariant features	48			
E	Mat	thematical interpretations of related work in our framework	51			

1 Introduction

In recent years, reinforcement learning from human feedback (RLHF) (Christiano et al., 2017) and variations like direct preference optimization (Rafailov et al., 2023) have become a practical approach for aligning language models (Ziegler et al., 2020; Stiennon et al., 2022; Bai et al., 2022; Ouyang et al., 2022). These techniques have then been used in the alignment of large-scale foundation models like ChatGPT (OpenAI, 2022) and GPT-4 (OpenAI et al., 2024), Gemini (Gemini Team, 2023), Llama (Touvron et al., 2023; Grattafiori et al., 2024), and Claude (Bai et al., 2022; Anthropic, 2023a;b).

RLHF uses feedback of human evaluators on AI behavior to give information about desired behavior. But as AI systems become more capable, they may eventually outstrip our ability to evaluate them. The problem of scalable oversight therefore asks how to effectively teach our preferences to AI systems when they become more capable than ourselves (Amodei et al., 2016; Christiano et al., 2018; Irving et al., 2018; Leike et al., 2018; Bowman et al., 2022). Conceptually, if human evaluators lack the capacity to fully understand the AI's behavior, then this may incentivize the AI to perform behavior that the human evaluator believes to be good, possibly at the expense of actual value.

Recent work has in fact revealed problems with RLHF that stem from erroneous human beliefs about AI behavior. In an early example (Amodei et al., 2017), an AI was supposed to grasp a ball with a robot hand in a simulation. The evaluators, who looked at a video of the behavior, sometimes believed that the hand was grasping the ball when it was in fact only in front of it, leading to positive feedback for the wrong behavior. In Cloud et al. (2024), a synthetic overseer in a baseline setting is unaware of whether a reached goal is a diamond (positive) or ghost (negative), leading the policy to learn to approach the ghost. Denison et al. (2024) shows a language model modifying a file unbeknownst to a synthetic evaluator, successfully deceiving it into believing a checklist has been accomplished. Similar models have also been shown to mislead humans who are limited in their time or competence to evaluate question answering or programming tasks, leading them to believe that the performance is better than it is (Wen et al., 2024).

Given the crucial role of evaluator beliefs in erroneous feedback, in this theoretical work, we are exploring the idea to *model* human beliefs about AI behavior. Our idea is that if we knew what the human *believes* the AI has done in its environment, then we could properly assign the human's feedback to that believed behavior instead of to the actual behavior. For example, imagine the robot hand from Amodei et al. (2017) is in front of the ball and the human evaluator gives positive feedback. Additionally, assume we know the

human believes the ball was *actually grasped*. If so, then we know that the human thinks grasping the ball is positive, and can incentivize this behavior from now on.

But what are beliefs? Prior work (Lang et al., 2024) takes the view that a belief is a probability distribution over state sequences in the AI's environment. They show in theoretical toy examples that modeling such beliefs can help to learn effectively from feedback. However, in reality it is unrealistic that humans form probability distributions over entire state sequences, and it would be prohibitive to specify such a belief explicitly. We think it is more realistic that humans think in terms of **features**, which we conceptualize as higher-level and independent properties of trajectories in an environment. We then model a human belief about an AI's behavior as a vector of feature strengths. This leads to our notion of a **human belief model**, which we introduce in Section 2. This model, together with feedback on observations, then theoretically allows to infer the return function over trajectories up to an inherent **ambiguity**, which we characterize in terms of the human belief model. The ambiguity disappears when the model is **complete**, which holds, in particular, when the human's beliefs of all observations linearly cover the feature combinations that are possible in the environment (Theorem 2.11). We then analyze conceptual toy examples of non-complete and complete models and consequences for the resulting return function inference.

We would like to use human belief modeling in practice to make concrete progress on scalable oversight. However, realistically, we cannot model human' beliefs precisely; after all, we do not even have an explicit specification of the human's feature set! In Section 3, we thus relax the requirement of exact modeling by investigating belief models that can **represent** all return functions and feedback functions that are compatible with the *true* belief model. We then say that these models **cover** the true belief model. When such a model is complete, it can replace the true model for the return function inference (Theorem 3.2). In a conceptual example, we analyze a human with symmetry-invariant features, which we can cover with a complete model that assumes symmetry-invariant reward functions.

We also characterize model coverage by the existence of what we call **model morphisms** and **linear belief-compatible ontology translations** (Theorem 3.5). Intuitively, this means that if one can *linearly translate* from the specified model's features to the human's features in a way that is compatible with their beliefs about observations, then coverage holds. The notion of a linear ontology translation is strikingly similar to work in the field of mechanistic interpretability on sparse autoencoders (Cunningham et al., 2023; Bricken et al., 2023), which linearly map from the internal representation space of language models to a space of human-interpretable features. Inspired by this connection, in Section 3.4, we make a preliminary proposal to use adapted foundation models to construct a belief model that covers the true human belief model, which we hope motivates future work that makes empirical progress on scalable oversight.

Finally, in our discussion in Section 4, we summarize our work, survey related work, motivate theoretical and empirical future work, and conclude.

2 Human belief models

In this section we define our generalized notion of human belief models that can help to infer return functions from the human's feedback.

In Section 2.1 we start by explaining our conventions and some preliminaries for linear algebra. We recommend also experienced readers to briefly skim this section to understand our notation. In Section 2.2, we briefly define Markov decision processes and slightly adapt them: We do *not* make the assumption that the return function over trajectories comes from a reward function. This allows us to be slightly more general. In Section 2.3, we then introduce the notion of the human's *ontology*, which is a map that assigns a vector of feature strengths to each trajectory. Our crucial assumption is that the return function evaluates a trajectory by summing up the rewards of each feature, weighted by the feature strengths. This allows us to recover classical reward functions as a special case.

In reinforcement learning, the goal is to maximize the return function, but we assume to be this function unknown: It represents the implicit values of a human evaluator and needs to be inferred from the human's feedback. In Section 2.4, we explain the human to give feedback based on forming a belief over features based on observations. We call the resulting feedback the observation return function. The fundamental question

is how to use it to infer the true return function. To do this, we assume that the whole human belief model, which we introduce in Section 2.5, is known; in other words, all aspects of the human's feedback process are known except the return parameters. Under this assumption, one can infer the return function from the human's feedback up to an ambiguity, which we characterize in Section 2.6. We also define complete human belief models for which the ambiguity vanishes and characterize them in Theorem 2.11. In Section 2.7, we introduce the dual notion of faithfulness. In Section 2.8, we study conceptual examples that give an intuition for when the ambiguity does, or does not, vanish.

2.1 Conventions and preliminaries for linear algebra

Let X be a set. For $x \in X$, we define the delta function $\delta_x : X \to \mathbb{R}$ as usual by

$$\delta_x(x') = \begin{cases} 1, & x' = x, \\ 0, & \text{else} \end{cases}$$

Let \mathbb{R} be the real numbers. Then \mathbb{R}^X denotes the vector space of functions from X to \mathbb{R} , which one can also view as column-vectors indexed by $x \in X$. For a function $v \in \mathbb{R}^X$, we write v(x) for the entry of v at position $x \in X$. The standard basis functions of \mathbb{R}^X are given by $\{e_x\}_{x \in X}$ with $e_x = \delta_x$. Then every function $v \in \mathbb{R}^X$ can be written as $v = \sum_{x \in X} v(x) \cdot e_x$. We define the standard scalar product $\langle \cdot, \cdot \rangle : \mathbb{R}^X \times \mathbb{R}^X \to \mathbb{R}$ by

$$\langle v, w \rangle \coloneqq \sum_{x \in X} v(x) w(x).$$

For a family of functions $v^i \in \mathbb{R}^X$ indexed by indices $i \in I$, we denote the span by

$$\mathbb{R}\left\langle v^i\mid i\in I\right\rangle \coloneqq \left\{\sum_{i\in I}\lambda_iv^i \ \middle| \ \lambda_i\in\mathbb{R} \text{ for all } i\in\mathbb{R}\right\}.$$

Let $A: \mathbb{R}^X \to \mathbb{R}^Y$ be a linear map. Then we view A also as a matrix with matrix elements $A_{yx} := [A(e_x)](y) \in \mathbb{R}$ for $x \in X, \ y \in Y$. We write $A_y \in \mathbb{R}^X$ for the row of A at index $y \in Y$, which is the function with entries $A_y(x) = A_{yx}$. Consequently, if $v \in \mathbb{R}^X$ and $y \in Y$, then we obtain

$$[A(v)](y) = \sum_{x \in X} v(x) [A(e_x)](y) = \sum_{x \in X} A_{yx} v(x) = \langle A_y, v \rangle.$$

This corresponds to the typical way that linear functions can be represented as matrix-vector products. If $\mathcal{V} \subseteq \mathbb{R}^X$ is a vector subspace, then we write $A|_{\mathcal{V}}: \mathcal{V} \to \mathbb{R}^Y$ for the restriction of A to \mathcal{V} , which is simply given by $(A|_{\mathcal{V}})(v) = A(v)$ for all $v \in \mathcal{V}$. We denote the kernel and image of $A: \mathbb{R}^X \to \mathbb{R}^Y$ by

$$\begin{split} \ker(A) &\coloneqq \left\{ v \in \mathbb{R}^X \ \middle| \ A(v) = 0 \right\} \subseteq \mathbb{R}^X, \\ \operatorname{im}(A) &\coloneqq \left\{ w \in \mathbb{R}^Y \ \middle| \ \exists v \in \mathbb{R}^X \colon A(v) = w \right\} \subseteq \mathbb{R}^Y. \end{split}$$

I.e., they are the set of functions in \mathbb{R}^X that are sent to zero by A, and the set of functions in \mathbb{R}^Y that are hit by A, respectively. Both are vector subspaces of their respective surrounding vector spaces. We explain basic properties of kernels and images in Appendix A.2 and will refer to those results when needed.

Sometimes, our linear functions "come from" a deterministic function in the other direction. I.e., if $h: Y \to X$ is a function, then we define the linear dual function $h^*: \mathbb{R}^X \to \mathbb{R}^Y$ for $v \in \mathbb{R}^X$ and $y \in Y$ by

$$[h^*(v)](y) := v(h(y)). \tag{1}$$

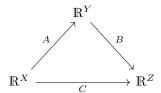
The matrix elements are then given by

$$h_{yx}^* = [h^*(e_x)](y) = e_x(h(y)) = \begin{cases} 1, & h(y) = x, \\ 0, & \text{else.} \end{cases}$$
 (2)

For two linear maps $A: \mathbb{R}^X \to \mathbb{R}^Y$ and $B: \mathbb{R}^Y \to \mathbb{R}^Z$, we write their composition as $B \circ A: \mathbb{R}^X \to \mathbb{R}^Z$, which has matrix elements $(B \circ A)_{zx} = \sum_{y \in Y} B_{zy} A_{yx}$ for $x \in X$, $z \in Z$. This also implies

$$(B \circ A)_z = \sum_{y \in Y} B_{zy} A_y. \tag{3}$$

Finally, whenever we draw a diagram of (usually linear) functions in this paper, the diagram *commutes*, meaning that all directed pathways from one node to another node are the same function. For example, in a diagram of the form



we have $C = B \circ A$. The only exceptions to this convention are Equations (20) and (26) in the appendix.

2.2 Preliminaries on Markov Decision Processes

We work in the following setting: We have an MDP (S, A, T, P_0, T, G) , with a finite set of states S and actions A, a transition kernel $T: S \times A \to \Delta(S)$, an initial state distribution $P_0 \in \Delta(S)$, a finite time horizon $T \in \{0, 1, 2, 3, ...\}$, and the human's implicit return function G, which we clarify after defining trajectories below. We fix this generic MDP for the rest of the paper.

We now define trajectories. In the rest of the paper, whenever we have a state-action sequence $\xi \in (\mathcal{S} \times \mathcal{A})^T \times \mathcal{S}$ present in some context and then write about states and actions $s_0, a_0, \ldots, s_{T-1}, a_{T-1}, s_T$, then they implicitly refer to the states and actions in ξ . Then the set of **trajectories** be given by

$$\Xi = \left\{ \xi \in (\mathcal{S} \times \mathcal{A})^T \times \mathcal{S} \ \middle| \ \xi \text{ is possible} \right\} \subseteq (\mathcal{S} \times \mathcal{A})^T \times \mathcal{S},$$

where we call a state-action sequence ξ possible if $P_0(s_0) > 0$ and $T(s_t \mid s_{t-1}, a_{t-1}) > 0$ for all $t \ge 1$. Then the return function is a function $G \in \mathbb{R}^{\Xi}$.

Note that this formalism does *not* assume that G decomposes into a reward function $R: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathbb{R}$. In other words, the formalism allows for human preferences that do not satisfy the reward hypothesis (Bowling et al., 2022) and is thus slightly more general than typical reinforcement learning settings.

Policies are functions $\pi: \mathcal{S} \to \Delta(\mathcal{A})$. A policy π together with the MDP induces a distribution $P^{\pi} \in \Delta(\Xi)$ over trajectories by sampling initial states from P_0 , actions from π , and transitions from \mathcal{T} . The policy evaluation function is then given by

$$J(\pi) := \underset{\xi \sim P^{\pi}(\cdot)}{\mathbf{E}} \left[G(\xi) \right]. \tag{4}$$

The goal in reinforcement learning is to find a policy π that maximizes this evaluation function.

2.3 The human's ontology and reward object

We assume that the return function $G:\Xi\to\mathbb{R}$ encodes what a given human cares about. We imagine that G measures the quality of each trajectory linearly in certain feature strengths associated to each trajectory; here, features are high-level return-relevant concepts.

More precisely, we assume the human comes equipped with a finite **feature set** Ω . The human's **ontology** is a function $\lambda : \Xi \to \mathbb{R}^{\Omega}$ that encodes for each $\xi \in \Xi$ and $\omega \in \Omega$ the extent $[\lambda(\xi)](\omega)$ to which the feature ω is present in the trajectory ξ . As we discuss in greater detail in Section 2.4, these feature strengths $\lambda(\xi)$

are idealized, i.e., the human may not be able to compute them. In the general theory, we allow them to be negative, though it may help to imagine them to be non-negative. The human's **reward object** is a fixed function $R_{\Omega} \in \mathbb{R}^{\Omega}$ that assigns to each feature $\omega \in \Omega$ the degree $R_{\Omega}(\omega)$ to which the human likes this feature. The return function $G: \Xi \to \mathbb{R}$ evaluates a trajectory by summing up the quality of all features, weighted by the extent to which they appear in the trajectory:

$$G(\xi) = \sum_{\omega \in \Omega} \left[\lambda(\xi) \right] (\omega) \cdot R_{\Omega}(\omega) = \left\langle \lambda(\xi), R_{\Omega} \right\rangle.$$

Let $\Lambda: \mathbb{R}^{\Omega} \to \mathbb{R}^{\Xi}$ be given by $[\Lambda(\tilde{R}_{\Omega})](\xi) \coloneqq \langle \lambda(\xi), \tilde{R}_{\Omega} \rangle$. Then Λ is a linear function from which we can recover λ using the matrix elements of $\Lambda: [\lambda(\xi)](\omega) = \Lambda_{\xi\omega}$ (Proposition A.1). We slightly abuse the terminology by referring to both $\lambda: \Xi \to \mathbb{R}^{\Omega}$ and $\Lambda: \mathbb{R}^{\Omega} \to \mathbb{R}^{\Xi}$ as the human's ontology. This representation of the ontology satisfies $\Lambda(R_{\Omega}) = G$.

Example 2.1 (Classical reward functions). Assume that the human's feature set is $\Omega = \mathcal{S} \times \mathcal{A} \times \mathcal{S}$, i.e., the set of all state-action-state transitions. Then the reward object $R_{\Omega} = R \in \mathbb{R}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}}$ is a reward function in the classical sense. Let $\gamma \in [0,1]$ be a discount factor and define the ontology $\mathbf{\Lambda} = \mathbf{\Gamma} : \mathbb{R}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}} \to \mathbb{R}^{\Xi}$ by

$$\left[\mathbf{\Gamma}(\tilde{R})\right](\xi) = \sum_{t=0}^{T-1} \gamma^t \tilde{R}(s_t, a_t, s_{t+1}).$$

 $G = \Gamma(R)$ is then given by the discounted sum of rewards of individual transitions in a trajectory, as is typical in reinforcement learning. The matrix elements of the ontology Γ are given by

$$\mathbf{\Gamma}_{\xi,(s,a,s')} = \left[\mathbf{\Gamma}(e_{(s,a,s')})\right](\xi) = \sum_{t=0}^{T-1} \gamma^t e_{(s,a,s')}(s_t, a_t, a_{t+1}) = \sum_{t=0}^{T-1} \gamma^t \delta_{(s,a,s')}(s_t, a_t, s_{t+1}),$$

In other words, the extent to which the "feature" (s, a, s') is present in the trajectory ξ is simply the discounted number of times that it appears.

2.4 The human's feature belief and observation return function

In standard reinforcement learning, it is typically assumed that we know the return function G to be maximized by a policy. However, G is the implicit return function of a given human. We must thus rely on feedback by the human to infer information about G. A naive idea would be to show the human each trajectory ξ , ask them to evaluate it by computing $G(\xi)$, and to use these returns to train a policy. There are three challenges to this plan:

- (i) The human might not have access to the full trajectory ξ or all trajectories $\xi \in \Xi$. For example, in complex environments, they would usually only receive partial observations.
- (ii) Even if ξ were fully accessible, the human might not have the *capacity* to compute the features $\lambda(\xi)$. For example, if ξ encodes a proof-attempt of a mathematical conjecture and ω encodes "correctness", then the human may not have the competence to determine the extent $[\lambda(\xi)](\omega)$ to which ξ is correct.
- (iii) Even if the human could fully compute $\lambda(\xi)$, they may not have full access to their reward object R_{Ω} in order to then compute $G(\xi)$ as an explicit number.

To address (i), we assume a set of **observations** \mathcal{O} that the human can receive. They may be come from an observation function $O: \Xi \to \mathcal{O}$, as in some examples from Amodei et al. (2017); Lang et al. (2024); Cloud et al. (2024); Denison et al. (2024). For example, $O(\xi)$ could be a sequence of observation segments, one for each time-step. Alternatively, \mathcal{O} could be a set of simulations to probe the human's opinion on specific situations that may be present in real trajectories. We can also imagine \mathcal{O} to be a subset of Ξ . For example, it can make sense to only show trajectories to the human that they are able to correctly evaluate, as in easy-to-hard generalization (Sun et al., 2024; Hase et al., 2024; Ding et al., 2024). In any case, we

assume \mathcal{O} to be a fixed set of observations where, for the most part, we are agnostic about the process that generates them.

To address (ii), we assume the human then has a **feature belief function**, i.e. a function $\epsilon : \mathcal{O} \to \mathbb{R}^{\Omega}$ that encodes for each $o \in \mathcal{O}$ the extent $[\epsilon(o)](\omega)$ to which the human *believes* the feature ω to be present in the observation o or an associated unknown trajectory ξ . Even if $\mathcal{O} = \Xi$, we can have $\epsilon \neq \lambda$, which reflects that the human may believe the features of a trajectory to be different from what they actually are.

The human then judges observations $o \in \mathcal{O}$ according to the **observation return function** $G_{\mathcal{O}} \in \mathbb{R}^{\mathcal{O}}$. It evaluates a trajectory by summing up the quality of all features, weighted by the extent to which the human believes them to be present:

$$G_{\mathcal{O}}(o) := \sum_{\omega \in \Omega} \left[\epsilon(o) \right] (\omega) \cdot R_{\Omega}(\omega) = \left\langle \epsilon(\xi), R_{\Omega} \right\rangle.$$

Let $\mathcal{E}: \mathbb{R}^{\Omega} \to \mathbb{R}^{\Xi}$ be given by $\left[\mathcal{E}(\tilde{R}_{\Omega})\right](\xi) := \left\langle \epsilon(o), \tilde{R}_{\Omega} \right\rangle$. Then as for Λ , \mathcal{E} is a linear function from which we can recover ϵ using the matrix elements of $\mathcal{E}: \left[\epsilon(o)\right](\omega) = \mathcal{E}_{o\omega}$ (Proposition A.1). We again slightly abuse the terminology by referring to $both \ \epsilon: \mathcal{O} \to \mathbb{R}^{\Omega}$ and $\mathcal{E}: \mathbb{R}^{\Omega} \to \mathbb{R}^{\mathcal{O}}$ as the human's feature belief function. This representation of the featre belief function satisfies $\mathcal{E}(R_{\Omega}) = G_{\mathcal{O}}$.

To address (iii), in RLHF it is typically assumed that the human can make *choices* between observations (Christiano et al., 2017) that *implicitly* reflect the underlying return function (or, on our case, observation return function). We discuss this setting in Appendix B for the special case that all relevant linear functions are row-constant. However, in the main paper, we do not want to overcomplicate the theory and assume that the human can *directly compute* $G_{\mathcal{O}}$. In principle, if we would show each $o \in \mathcal{O}$ to the human, we could then record the entire observation return function $G_{\mathcal{O}}$. Thus, we assume $G_{\mathcal{O}}$, which represents the human's feedback, to be fully known.

2.5 The definition of human belief models

How can we infer the return function G from the human's feedback, represented by $G_{\mathcal{O}}$? In order to do this, we assume the features Ω , ontology Λ , and feature belief function \mathcal{E} are all known; additionally, we assume we may have a priori knowledge of a vector subspace of **valid reward objects** $\mathcal{V} \subseteq \mathbb{R}^{\Omega}$ in which R_{Ω} resides: $R_{\Omega} \in \mathcal{V}$. This may come from any a priori knowledge, e.g. of certain symmetries in the environment that leave rewards invariant (cf. Example 2.17 and Section 3.3). This leads to the following notion:

Definition 2.2 (Human belief model, representing). Let Ξ be the set of trajectories in an MDP and \mathcal{O} the fixed set of observations that a human evaluator receives. Then a human belief model (or belief model, or model, for short) is a tuple $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$, where:

- Ω is a set called feature set;
- $\Lambda : \mathbb{R}^{\Omega} \to \mathbb{R}^{\Xi}$ is a linear function, called ontology;
- $\mathcal{E}: \mathbb{R}^{\Omega} \to \mathbb{R}^{\mathcal{O}}$ is a linear function, called feature belief function;
- and $\mathcal{V} \subseteq \mathbb{R}^{\Omega}$ is a vector subspace, called space of valid reward objects.

A human belief model \mathcal{M} represents the true return function $G:\Xi\to\mathbb{R}$ and observation return function $G_{\mathcal{O}}:\mathcal{O}\to\mathbb{R}$ if there is a reward object $R_{\Omega}\in\mathcal{V}$ with $G=\Lambda(R_{\Omega})$ and $G_{\mathcal{O}}=\mathcal{E}(R_{\Omega})$.

When appropriate, we will also use the representation $\lambda : \Xi \to \mathbb{R}^{\Omega}$ with $[\lambda(\xi)](\omega) = \mathbf{\Lambda}_{\xi\omega}$ and $\epsilon : \mathcal{O} \to \mathbb{R}^{\Omega}$ with $[\epsilon(o)](\omega) = \mathcal{E}_{o\omega}$ of the ontology and feature belief function, respectively (cf. Proposition A.1). We

visualize a model $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$ as



where we use the latter version in the special case that $\mathcal{V} = \mathbb{R}^{\Omega}$.

Remark 2.3 (Which human belief model?). A priori, there can be many human belief models $\widehat{\mathcal{M}} = (\widehat{\Omega}, \widehat{\Lambda}, \widehat{\mathcal{E}}, \widehat{\mathcal{V}})$ that can represent the true return function G and the human's feedback $G_{\mathcal{O}}$. For example, in Section 3.3 we will discuss a situation with three different models, and in Appendix C one with many more. In fact, there is always a trivial belief model that can represent G and $G_{\mathcal{O}}$: View \mathbb{R} as the \mathbb{R} -vectorspace $\mathbb{R}^{\{\star\}}$ of functions from a one-element set $\{\star\}$ to \mathbb{R} . Associate to $G: \Xi \to \mathbb{R} = \mathbb{R}^{\{\star\}}$ and $G_{\mathcal{O}}: \mathcal{O} \to \mathbb{R} = \mathbb{R}^{\{\star\}}$ via Proposition A.1 the linear functions $\lim_{X \to \mathbb{R}} (G) : \mathbb{R}^{\{\star\}} \to \mathbb{R}^{\mathbb{Z}}$ and $\lim_{X \to \mathbb{R}} (G_{\mathcal{O}}) : \mathbb{R}^{\{\star\}} \to \mathbb{R}^{\mathcal{O}}$ given by

$$\big[\, \mathrm{lin}(G) \big] (r \cdot e_\star) \coloneqq r \cdot G, \quad \big[\, \mathrm{lin}(G_{\mathcal{O}}) \big] (r \cdot e_\star) \coloneqq r \cdot G_{\mathcal{O}}$$

for $r \in \mathbb{R}$ and the only e_{\star} . Then, define the human belief model

$$\mathcal{M} \coloneqq (\{\star\}, \ \operatorname{lin}(G), \ \operatorname{lin}(G_{\mathcal{O}}), \ \mathbb{R}^{\{\star\}})$$

with the feature set $\{\star\}$. Then this represents G and $G_{\mathcal{O}}$ with the reward object $e_{\star} \in \mathbb{R}^{\{\star\}}$ since $\left[\operatorname{lin}(G) \right](e_{\star}) = G$ and $\left[\operatorname{lin}(G_{\mathcal{O}}) \right](e_{\star}) = G_{\mathcal{O}}$. Intuitively, the feature \star then means "goodness", the "ontology" $G : \Xi \to \mathbb{R}^{\{\star\}}$ measures the extent to which "goodness" is present in a trajectory, and the "feature belief function" $G_{\mathcal{O}} : \mathcal{O} \to \mathbb{R}^{\{\star\}}$ measures how much goodness the human believes to be present in an observation. This trivial model is not very interesting for our purposes since assuming that the human belief model is known would then presuppose knowledge of G itself — which is the return function that we want to infer based on feedback $G_{\mathcal{O}}$.

A belief model $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$ is more interesting if Ω consists of a variety of meaningful concepts such that G and $G_{\mathcal{O}}$ decompose linearly over these features, in such a way that:

- (i) It is easier to know the model \mathcal{M} than to know the return function G.
- (ii) Given \mathcal{M} and $G_{\mathcal{O}}$, it is easy to determine G.

For the rest of this section, we imagine to know the "true" belief model \mathcal{M} , which we assume to have these properties. Note that these are philosophical assumptions that guide how we talk about \mathcal{M} and how we imagine its application. But mathematically, all of our claims here and in the upcoming sections are correct whenever \mathcal{M} represents G and $G_{\mathcal{O}}$.

In Section 3, we then relax the requirement of "knowledge" of \mathcal{M} , and specifically in Section 3.4 we make a proposal for model specification using foundation models. We then also show how G could in principle be learned by learning $G_{\mathcal{O}}$ via supervised learning. We think that section will make it plausible that properties (i) and (ii) can hold in practice.

Example 2.4. Assume the setting from Example 2.1, in which $\Omega = S \times A \times S$ and $\Lambda = \Gamma : \mathbb{R}^{S \times A \times S} \to \mathbb{R}^{\Xi}$. We now explain the belief model from Lang et al. (2024) in our framework, thus showing that we generalize their work.

Namely, assume the human comes equipped with a belief function $b: \mathcal{O} \to \mathbb{R}^{\Xi}$ that assigns a probability distribution over trajectories $b(o) \in \Delta(\Xi) \subseteq \mathbb{R}^{\Xi}$ to each observation $o \in \mathcal{O}$. For example, this could be

a posterior belief if o is sampled by first sampling a trajectory ξ according to some distribution and then computing the observation as $o = O(\xi)$ for an observation function $O : \Xi \to \mathcal{O}$.

Then we assume that $\epsilon: \mathcal{O} \to \mathbb{R}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}}$ computes, for each $o \in \mathcal{O}$ and $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$, the expected discounted number of times that the transition (s, a, s') appears in the trajectory ξ expected probabilistically from observing o:

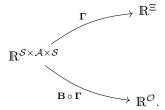
$$\left[\epsilon(o)\right](s,a,s') \coloneqq \sum_{\xi \in \Xi} \left[b(o)\right](\xi) \cdot \mathbf{\Gamma}_{\xi,(s,a,s')}.$$

Let $\mathbf{B}: \mathbb{R}^{\Xi} \to \mathbb{R}^{\mathcal{O}}$ be the linear function given by $[\mathbf{B}(G)](o) := \langle b(o), G \rangle$, which has matrix elements $\mathbf{B}_{o\xi} = [b(o)](\xi)$ by Proposition A.1. We obtain

$$\mathcal{E}_{o,(s,a,s')} = [\epsilon(o)](s,a,s') = \sum_{\xi \in \Xi} \mathbf{B}_{o\xi} \, \mathbf{\Gamma}_{\xi,(s,a,s')} = (\mathbf{B} \circ \mathbf{\Gamma})_{o,(s,a,s')}.$$

Thus, $\mathcal{E} = \mathbf{B} \circ \mathbf{\Gamma}$.

We finally assume that the set of valid reward objects is simply given by $\mathcal{V} = \mathbb{R}^{S \times A \times S}$. Thus, in total, our model is given by $(\Omega, \Lambda, \mathcal{E}, \mathcal{V}) = (S \times A \times S, \Gamma, \mathbf{B} \circ \Gamma, \mathbb{R}^{S \times A \times S})$:



2.6 Complete belief models and the ambiguity

In this whole subsection, we fix an MDP with trajectories Ξ , observations \mathcal{O} , and corresponding return function $G \in \mathbb{R}^{\Xi}$ and observation return function $G_{\mathcal{O}} \in \mathbb{R}^{\mathcal{O}}$. We also fix a human belief model $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$ that represents G and $G_{\mathcal{O}}$ with an implicit reward object $R_{\Omega} \in \mathcal{V}$: $\Lambda(R_{\Omega}) = G$ and $\mathcal{E}(R_{\Omega}) = G_{\mathcal{O}}$. We can visualize the model together with the reward object and return functions as follows:



We are concerned with the following question:

Question 2.5. Assume the human belief model \mathcal{M} and the human feedback, in the form of the observation return function $G_{\mathcal{O}}$, are known. Under what conditions, or to what extent, can we infer the return function G?

In Section 3, we will relax the assumption that the belief model is known precisely.

The idea for the answer is as follows: When trying to infer G from $G_{\mathcal{O}}$, we make use of the knowledge that they are connected by the unknown reward object $R_{\Omega} \in \mathcal{V}$. Thus, one can first try to determine a reward object $\tilde{R}_{\Omega} \in \mathcal{V}$ with the correct observation return function $\mathcal{E}(\tilde{R}_{\Omega}) = G_{\mathcal{O}}$. Then the question is whether the corresponding return function $\Lambda(\tilde{R}_{\Omega})$ equals the true return function G. They differ by the return function $G' = \Lambda(\tilde{R}_{\Omega}) - G$. The ambiguity will then be defined as the set of all these differences, and the question will be how to characterize it and when it vanishes, leading to the notion of a complete human belief model.

Definition 2.6 (Feedback-compatible, ambiguity). We define the set of return functions that are **feedback-compatible** with $G_{\mathcal{O}}$ as

$$\mathrm{FC}^{\mathcal{M}}(G_{\mathcal{O}}) \coloneqq \Big\{ \tilde{G} \in \mathbb{R}^\Xi \ \big| \ \exists \tilde{R}_{\Omega} \in \mathcal{V} \colon \mathcal{E}(\tilde{R}_{\Omega}) = G_{\mathcal{O}} \ and \ \mathbf{\Lambda}(\tilde{R}_{\Omega}) = \tilde{G} \Big\}.$$

We define the **ambiguity** left in the return function G after the observation return function $G_{\mathcal{O}}$ is known by

$$\operatorname{Amb}^{\mathcal{M}}(G, G_{\mathcal{O}}) := \left\{ G' \in \mathbb{R}^{\Xi} \mid G' = \tilde{G} - G \text{ for } \tilde{G} \in \operatorname{FC}^{\mathcal{M}}(G_{\mathcal{O}}) \right\}.$$

Then clearly, we have

$$FC^{\mathcal{M}}(G_{\mathcal{O}}) = G + Amb^{\mathcal{M}}(G, G_{\mathcal{O}}).$$

This leaves open the question of how to characterize the ambiguity and when it vanishes.

Proposition 2.7 (Ambiguity characterization). We have

$$Amb^{\mathcal{M}}(G, G_{\mathcal{O}}) = \mathbf{\Lambda}(\ker(\mathcal{E}) \cap \mathcal{V}).$$

Proof. For one direction, let $G' \in \text{Amb}^{\mathcal{M}}(G, G_{\mathcal{O}})$. Then $G' = \Lambda(\tilde{R}_{\Omega}) - G$ for $\tilde{R}_{\Omega} \in \mathcal{V}$ with $\mathcal{E}(\tilde{R}_{\Omega}) = G_{\mathcal{O}}$. By the linearity of \mathcal{E} we obtain

$$\mathcal{E}(\tilde{R}_{\Omega} - R_{\Omega}) = \mathcal{E}(\tilde{R}_{\Omega}) - G_{\mathcal{O}} = 0$$

and thus $\tilde{R}_{\Omega} - R_{\Omega} \in \ker(\mathcal{E}) \cap \mathcal{V}$. Since Λ is linear, we obtain

$$G' = \Lambda(\tilde{R}_{\Omega}) - G = \Lambda(\tilde{R}_{\Omega}) - \Lambda(R_{\Omega}) = \Lambda(\tilde{R}_{\Omega} - R_{\Omega}) \in \Lambda(\ker(\mathcal{E}) \cap \mathcal{V}).$$

For the other direction, let $G' \in \Lambda(\ker(\mathcal{E}) \cap \mathcal{V})$. Then $G' = \Lambda(R'_{\Omega})$ for $R'_{\Omega} \in \ker(\mathcal{E}) \cap \mathcal{V}$. Define $\tilde{R}_{\Omega} := R_{\Omega} + R'_{\Omega} \in \mathcal{V}$. By the linearity of \mathcal{E} and by $R'_{\Omega} \in \ker(\mathcal{E})$ we obtain

$$\mathcal{E}(\tilde{R}_{\Omega}) = \mathcal{E}(R_{\Omega}) + \mathcal{E}(R'_{\Omega}) = G_{\mathcal{O}} + \mathcal{E}(R'_{\Omega}) = G_{\mathcal{O}}.$$

Finally, the linearity of Λ shows

$$G' = \Lambda(R'_{\Omega}) = \Lambda(\tilde{R}_{\Omega} - R_{\Omega}) = \Lambda(\tilde{R}_{\Omega}) - \Lambda(R_{\Omega}) = \Lambda(\tilde{R}_{\Omega}) - G.$$

This shows $G' \in Amb^{\mathcal{M}}(G, G_{\mathcal{O}})$.

See Proposition B.6 for a version of the preceding proposition where the feedback is given by a choice probability function instead of $G_{\mathcal{O}}$. See Appendix C.4 for the ambiguities of many human belief models.

Remark 2.8. In light of the previous proposition, it turns out that the ambiguity does not depend on the true return function and observation return function, and we can thus simply write it as $Amb^{\mathcal{M}}$.

Example 2.9. We continue Example 2.4, with the model given by $\mathcal{M} = (\mathcal{S} \times \mathcal{A} \times \mathcal{S}, \Gamma, \mathbf{B} \circ \Gamma, \mathbb{R}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}})$. Proposition 2.7 implies that the ambiguity is given by $\mathrm{Amb}^{\mathcal{M}} = \Gamma(\ker(\mathbf{B} \circ \Gamma)) = \ker(\mathbf{B}) \cap \mathrm{im}(\Gamma)$. Lang et al. (2024) contains characterizations of this ambiguity in examples and special cases.

It is important to know when the ambiguity disappears, which is captured as *completeness* in the following definition.

Definition 2.10 (Completeness). We call the human belief model $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$ complete if $\ker(\mathcal{E}) \cap \mathcal{V} \subseteq \ker(\Lambda) \cap \mathcal{V}$, i.e. if $\operatorname{Amb}^{\mathcal{M}} = \Lambda(\ker(\mathcal{E}) \cap \mathcal{V}) = 0$ (cf. Proposition 2.7).

We will explain in Interpretation 2.12 why we call this property completeness. We find equivalent and sufficient conditions of completeness in the following proposition:

Theorem 2.11 (Completeness characterization). Remember the representations $\lambda : \Xi \to \mathbb{R}^{\Omega}$ and $\epsilon : \mathcal{O} \to \mathbb{R}^{\Omega}$ of the ontology $\mathbf{\Lambda} : \mathbb{R}^{\Omega} \to \mathbb{R}^{\Xi}$ and feature belief function $\mathcal{E} : \mathbb{R}^{\Omega} \to \mathbb{R}^{\mathcal{O}}$, respectively. Consider the following four statements:

- 1. M is complete.
- 2. There exists a linear function $Z: \mathbb{R}^{\mathcal{O}} \to \mathbb{R}^{\Xi}$ with $\Lambda|_{\mathcal{V}} = Z \circ \mathcal{E}|_{\mathcal{V}}$.
- 3. There exists a linear function $Z: \mathbb{R}^{\mathcal{O}} \to \mathbb{R}^{\Xi}$ with $\Lambda = Z \circ \mathcal{E}$.
- 4. For all $\xi \in \Xi$, $\lambda(\xi) \in \mathbb{R}^{\Omega}$ is contained in the span of the image of ϵ : $\lambda(\xi) \in \mathbb{R} \langle \epsilon(o) \mid o \in \mathcal{O} \rangle$.

Then 1 and 2 are equivalent, and 3 and 4 are equivalent and imply 1 and 2:

$$1 \bigcirc 2 \longleftarrow 3 \bigcirc 4$$

Furthermore, the linear function Z from 2 and 3 satisfies $G = Z(G_{\mathcal{O}})$. Finally, if $\mathcal{V} = \mathbb{R}^{\Omega}$, then all four statements are equivalent.

Proof. Note that $\ker(\mathbf{\Lambda}) \cap \mathcal{V} = \ker(\mathbf{\Lambda}|_{\mathcal{V}})$ and $\ker(\mathcal{E}) \cap \mathcal{V} = \ker(\mathcal{E}|_{\mathcal{V}})$. Thus, the equivalence of statements 1 and 2 immediately follow from Proposition A.3. That 3 implies 2 is obvious.

Now assume statement 4. For all $\xi \in \Xi$ and $o \in \mathcal{O}$, let $Z_{\xi o} \in \mathbb{R}$ be coefficients with

$$\lambda(\xi) = \sum_{o \in \mathcal{O}} Z_{\xi o} \epsilon(o).$$

This implies

$$\Lambda_{\xi\omega} = \sum_{o \in \mathcal{O}} Z_{\xi o} \mathcal{E}_{o\omega}.$$

Define the linear function $Z: \mathbb{R}^{\mathcal{O}} \to \mathbb{R}^{\Xi}$ to have the matrix elements $Z_{\xi o}$. Then the previous equation implies $\Lambda = Z \circ \mathcal{E}$, proving statement 3. That 3 implies 4 follows by reversing these arguments.

Assume statement 2. Then

$$G = \Lambda|_{\mathcal{V}}(R_{\Omega}) = (Z \circ \mathcal{E}|_{\mathcal{V}})(R_{\Omega}) = Z(\mathcal{E}(R_{\Omega})) = Z(G_{\mathcal{O}}).$$

That all statements are equivalent if $\mathcal{V} = \mathbb{R}^{\Omega}$ is clear.

Interestingly, for complete models, the preceding proposition shows that $G = Z(G_{\mathcal{O}})$ for a linear function Z that only depends on the human belief model \mathcal{M} , thus showing once more that we can infer G from $G_{\mathcal{O}}$ if the model is complete. Since Z may be impractical to find, it may however be more useful to determine G as the unique feedback-compatible return function by first finding \tilde{R}_{Ω} with $\mathcal{E}(\tilde{R}_{\Omega}) = G_{\mathcal{O}}$.

Interpretation 2.12. Overall, we have thus answered Question 2.5: When the human belief model is known, then G can be inferred from $G_{\mathcal{O}}$ up to the ambiguity $\mathrm{Amb}^{\mathcal{M}} = \Lambda(\ker(\mathcal{E}) \cap \mathcal{V})$, which vanishes if the model \mathcal{M} is complete. Looking back at the definition of feedback-compatible return functions, G is then determined as $G = \Lambda(\tilde{R}_{\Omega})$ for any $\tilde{R}_{\Omega} \in \mathcal{V}$ that gives rise to the human's feedback: $\mathcal{E}(\tilde{R}_{\Omega}) = G_{\mathcal{O}}$. For completeness, we then found further equivalent and sufficient conditions in Theorem 2.11.

We now explain why we chose the name completeness. By Theorem 2.11, statement 3, a sufficient condition for completeness is the existence of coefficients $Z_{\xi o} \in \mathbb{R}$ such that for all $\xi \in \Xi$, we can write $\lambda(\xi)$ as a linear combination of the $\epsilon(o)$:

$$\lambda(\xi) = \sum_{o \in \mathcal{O}} Z_{\xi o} \epsilon(o).$$

This means that the observation-dependent feature beliefs $\epsilon(o)$ span the true feature strengths of trajectories. In this sense, the belief model is "complete": The human's interpretations of observations sufficiently entail what happens in the MDP. It is thus not a surprise that completeness implies that the ambiguity disappears, and thus that we can infer the return function from the observation return function $G_{\mathcal{O}}$.

2.7 Faithful belief models

Let again $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$ be a human belief model for an MDP together with trajectories Ξ and observations \mathcal{O} , and let $R_{\Omega} \in \mathcal{V}$ with $G = \Lambda(R_{\Omega})$, $G_{\mathcal{O}} = \mathcal{E}(R_{\Omega})$. Here, we briefly study the notion of faithfulness of human belief models, which is dual to completeness:

Definition 2.13 (Faithfulness). We call the human belief model \mathcal{M} faithful if $\ker(\Lambda) \cap \mathcal{V} \subseteq \ker(\mathcal{E}) \cap \mathcal{V}$.

Many human belief models we study in this paper are faithful. The model from Example 2.9 is faithful since $\ker(\Gamma) \subseteq \ker(\mathbf{B} \circ \Gamma)$. All models from Section 3.3 are faithful. In Appendix C.2 we study several more faithful human belief models. We can find a characterization of faithfulness that is completely dual to Theorem 2.11:

Proposition 2.14 (Faithfulness characterization). Remember the representations $\lambda : \Xi \to \mathbb{R}^{\Omega}$ and $\epsilon : \mathcal{O} \to \mathbb{R}^{\Omega}$ of the ontology $\mathbf{\Lambda} : \mathbb{R}^{\Omega} \to \mathbb{R}^{\Xi}$ and feature belief function $\mathcal{E} : \mathbb{R}^{\Omega} \to \mathbb{R}^{\mathcal{O}}$, respectively. Consider the following four statements:

- 1. M is faithful.
- 2. There exists a linear function $Y: \mathbb{R}^{\Xi} \to \mathbb{R}^{\mathcal{O}}$ with $\mathcal{E}|_{\mathcal{V}} = Y \circ \mathbf{\Lambda}|_{\mathcal{V}}$.
- 3. There exists a linear function $Y: \mathbb{R}^{\Xi} \to \mathbb{R}^{\mathcal{O}}$ with $\mathcal{E} = Y \circ \Lambda$.
- 4. For all $o \in \mathcal{O}$, $\epsilon(o) \in \mathbb{R}^{\Omega}$ is contained in the span of the image of λ : $\epsilon(o) \in \mathbb{R}\langle \lambda(\xi) \mid \xi \in \Xi \rangle$.

Then 1 and 2 are equivalent, and 3 and 4 are equivalent and imply 1 and 2:

$$1 \bigcirc 2 \longleftarrow 3 \bigcirc 4$$

Furthermore, the linear function Y from 2 and 3 satisfies $G_{\mathcal{O}} = Y(G)$. Finally, if $\mathcal{V} = \mathbb{R}^{\Omega}$, then all four statements are equivalent.

Proof. The proof is analogous to the one for Theorem 2.11.

Interpretation 2.15. By Proposition 2.14, a sufficient condition for faithfulness is the existence of coefficient $Y_{o\xi} \in \mathbb{R}$ such that for all $o \in \mathcal{O}$, we can write $\epsilon(o)$ as a linear combination of the $\lambda(\xi)$:

$$\epsilon(o) = \sum_{\xi \in \Xi} Y_{o\xi} \lambda(\xi).$$

We can suggestively interpret $Y_{o\xi}$ as the human's "belief" that the true trajectory ξ was responsible for the observation o.¹ Under this interpretation, the feature strengths $\epsilon(o)$ are like an "expected value" of true feature strengths, given the human's beliefs. Thus, the feature beliefs encoded by ϵ are "faithful" to a belief over the actual MDP. This interpretation is especially valid for Example 2.4, where Y is given by \mathbf{B} , and thus $Y_{o\xi} = \mathbf{B}_{o\xi} = [b(o)](\xi)$ is the probability of the trajectory ξ , given o.

Faithfulness is also philosophically reasonable for a related reason. Assume R_{Ω} , $\tilde{R}_{\Omega} \in \mathcal{V}$ are two reward objects with $\Lambda(R_{\Omega}) = \Lambda(\tilde{R}_{\Omega})$. Then the resulting return functions coincide in their evaluation of all trajectories. It would then be reasonable to assume that they also coincide in their evaluation of observations, insofar as observations give information about a state of affairs in the actual MDP: $\mathcal{E}(R_{\Omega}) = \mathcal{E}(\tilde{R}_{\Omega})$. This requires that $\Lambda(R_{\Omega} - \tilde{R}_{\Omega}) = 0$ implies $\mathcal{E}(R_{\Omega} - \tilde{R}_{\Omega}) = 0$, which exactly means that $\ker(\Lambda) \cap \mathcal{V} \subseteq \ker(\mathcal{E}) \cap \mathcal{V}$, i.e. the faithfulness of the model. However, since human evaluators do not necessarily have feature beliefs that are "this rational", and since this property is not needed in the rest of our theory, we do not assume faithfulness of our human belief models in the general theory.

¹Though note that Y is usually not unique and $Y_o \in \mathbb{R}^{\Xi}$ is usually not an actual probability distribution over ξ .

2.8 Conceptual examples

We now present some simple conceptual examples to illustrate the theory with respect to the ambiguity. We will not define entire MDPs but only those parts of the formalism that are necessary to reason about the ambiguity. In all our examples, we assume $\mathcal{O} \subseteq \Xi$, i.e., observations are entire trajectories: They are fully observed. This brings other factors than observability into the focus, like the human's *understanding* of the trajectories, and whether there is enough data. We refer to Section 3.3 and Lang et al. (2024) for concrete examples where the MDP is defined entirely and a possibly remaining ambiguity stems from partial observability.

In all examples, we have a feature set $\Omega = \{\omega_1, \omega_2, \omega_3, \omega_4\}$ whose meaning will vary in each scenario. We also have observations o_1, o_2, o_3, o_4, o_5 , though the observation set $\mathcal{O} \subseteq \Xi$ is sometimes only a subset of $\{o_1, o_2, o_3, o_4, o_5\}$. We will make use of the feature belief function $\epsilon : \mathcal{O} \to \mathbb{R}^{\Omega} \cong \mathbb{R}^4$. We represent each $\epsilon(o)$ as a row-vector with entries $\epsilon(o)_i = [\epsilon(o)](\omega_i)$ as follows:

$$\epsilon(o_1) = (1, 0, 0, 1)
\epsilon(o_2) = (0, 1, 0, 1)
\epsilon(o_3) = (0, 0, 1, 1)
\epsilon(o_4) = (0, 3, 0, 2)
\epsilon(o_5) = (0, 0, 0, 1).$$
(5)

In all examples, we assume Ξ , \mathcal{V} , and Λ (and thus overall $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$) to also be known, though we do not make them explicit.

Example 2.16 (Simple Completeness). The goal is to create children stories $\xi \in \Xi$ that appeal to a specific child Alice. We assume that it is known that children care about the presence or absence of exactly the following four features $\Omega = \{\omega_1, \omega_2, \omega_3, \omega_4\}$:

- ω_1 : Rule-breaking.
- ω_2 : Non-linear story-telling.
- ω_3 : Scariness.
- ω_4 : Moral lessons.

We show Alice four stories $\mathcal{O} = \{o_1, o_2, o_3, o_4\} \subseteq \Xi$ that give rise to this feature belief function (cf. Equation (5)):

$$\mathcal{E}: \mathbb{R}^{\Omega} \to \mathbb{R}^{\mathcal{O}}, \quad \mathcal{E} = \begin{pmatrix} 1 & 0 & 0 & 1\\ 0 & 1 & 0 & 1\\ 0 & 0 & 1 & 1\\ 0 & 3 & 0 & 2 \end{pmatrix}. \tag{6}$$

Conceptually, we assume these are the true feature strengths, meaning that $\lambda(o) = \epsilon(o)$ for all $o \in \mathcal{O}$, though this detail does not matter for the ambiguity analysis. The matrix elements mean that o_1 contains rule-breaking and moral lessons, but the story-telling is linear and the story is not scary. o_4 is very non-linear and contains quite some moral lessons, but shows no rule-breaking or scariness. Since the rows of Equation (6) form a basis of \mathbb{R}^{Ω} , we obtain

$$\mathbb{R}\langle \epsilon(o) \mid o \in \mathcal{O} \rangle = \mathbb{R}^{\Omega},$$

and so by Theorem 2.11 it follows that the model is complete. Thus, the ambiguity vanishes and we can infer the return function from Alice's feedback $G_{\mathcal{O}}$ on the stories in \mathcal{O} (cf. Interpretation 2.12).

We demonstrate this now explicitly. Assume R_{Ω} is Alice's (a priori unknown) reward object, and that the resulting observation return function $G_{\mathcal{O}} = \mathcal{E}(R_{\Omega}) \in \mathbb{R}^{\mathcal{O}}$ is given as follows:

$$G_{\mathcal{O}}(o_1) = -3, \quad G_{\mathcal{O}}(o_2) = 1, \quad G_{\mathcal{O}}(o_3) = 0, \quad G_{\mathcal{O}}(o_4) = 4.$$
 (7)

Representing R_{Ω} and $G_{\mathcal{O}}$ as column-vectors, this means:

$$\begin{pmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 3 & 0 & 2 \end{pmatrix} \cdot R_{\Omega} = \mathcal{E}(R_{\Omega}) = G_{\mathcal{O}} = \begin{pmatrix} -3 \\ 1 \\ 0 \\ 4 \end{pmatrix}.$$

This results in

$$R_{\Omega} = \begin{pmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 3 & 0 & 2 \end{pmatrix}^{-1} \cdot \begin{pmatrix} -3 \\ 1 \\ 0 \\ 4 \end{pmatrix} = \begin{pmatrix} 1 & -3 & 0 & 1 \\ 0 & -2 & 0 & 1 \\ 0 & -3 & 1 & 1 \\ 0 & 3 & 0 & -1 \end{pmatrix} \cdot \begin{pmatrix} -3 \\ 1 \\ 0 \\ 4 \end{pmatrix} = \begin{pmatrix} -2 \\ 2 \\ 1 \\ -1 \end{pmatrix}.$$

Thus, we recovered Alice's reward object: She positively values non-linear story-telling (2) and scariness (1) but does not like moral lessons (-1) or rule-breaking (-2). Together with the knowledge of the ontology Λ to associate features to stories, we can create more enjoyable stories for Alice.

This example shows that knowledge of return-relevant features and a successful modeling of the feature belief function can go a long way to determine the correct return function: We essentially need only four "data points" (in the form of observation returns $G_{\mathcal{O}}(o)$) to determine G.

Example 2.17 (Completeness from symmetries). We assume the same situation as in Example 2.16 except this time, we assume that we only have three observations $\mathcal{O} = \{o_1, o_2, o_4\}$:

$$\mathcal{E}: \mathbb{R}^{\Omega} \to \mathbb{R}^{\mathcal{O}}, \quad \mathcal{E} = \begin{pmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 3 & 0 & 2 \end{pmatrix}.$$
 (8)

Note that the third row now represents $\epsilon(o_4)$ since o_3 is not in \mathcal{O} . Assume we have a priori information that children who like scariness do not like moral lessons and vice versa; i.e., it is known that Alice's reward object satisfies the following equation:

$$R_{\Omega}(\omega_3) = -R_{\Omega}(\omega_4).$$

The set of all reward functions with this property forms a vector subspace $\mathcal{V} \subseteq \mathbb{R}^{\Omega}$. It turns out that this implies $\ker(\mathcal{E}) \cap \mathcal{V} = \{0\}$. Indeed, let $R'_{\Omega} \in \ker(\mathcal{E}) \cap \mathcal{V}$. Then $\mathcal{E}(R'_{\Omega}) = 0$ implies $R'_{\Omega}(\omega_1) = R'_{\Omega}(\omega_2) = R'_{\Omega}(\omega_4) = 0$, and $R'_{\Omega} \in \mathcal{V}$ implies $R'_{\Omega}(\omega_3) = -R'_{\Omega}(\omega_4) = 0$. With $\ker(\mathcal{E}) \cap \mathcal{V} = 0$, the ambiguity vanishes: Amb^M = $\Lambda(\ker(\mathcal{E}) \cap \mathcal{V}) = 0$. Thus, the model is complete and the return function can be inferred from Alice's feedback $G_{\mathcal{O}}$ on $\mathcal{O} = \{o_1, o_2, o_4\}$ (cf. Interpretation 2.12).

Example 2.18 (Ambiguity from undetected vulnerabilities). The goal is to correctly evaluate coding-agents to produce valid and safe code blocks $\xi \in \Xi$. We assume that exactly the following four features are relevant for evaluating code:

• ω_1 : Code-vulnerability.

• ω_2 : Syntax error.

• ω_3 : Simplicity.

ω₄: Validity.

We show the human evaluator four code-blocks $\mathcal{O} = \{o_2, o_3, o_4, o_5\} \subseteq \Xi$ that give rise to this feature belief function (cf. Equation (5)):

$$\mathcal{E} \colon \mathbb{R}^{\Omega} \to \mathbb{R}^{\mathcal{O}}, \quad \mathcal{E} = \begin{pmatrix} 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 3 & 0 & 2 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

Crucially, we can assume that o_5 does contain a code-vulnerability, but the human does not detect it. In other words, the ontology $\lambda : \Xi \to \mathbb{R}^{\Omega}$ assigns the feature strengths $\lambda(o_5) = (1,0,0,1) \neq (0,0,0,1) = \epsilon(o_5)$. If the human had detected the code-vulnerability, then o_5 would be mathematically equivalent to o_1 in Example 2.16, the model would be complete, and the ambiguity would disappear.

However, this is not the case. Assuming $V = \mathbb{R}^{\Omega}$, we obtain that $\ker(\mathcal{E}) \cap V \neq \{0\}$ contains the reward object $R'_{\Omega} \in \mathbb{R}^{\Omega}$ with $R'_{\Omega}(\omega_1) = 1$ and $R'_{\Omega}(\omega) = 0$ for all $\omega \neq \omega_1$. With R_{Ω} being the unknown true reward object and $\tilde{R}_{\Omega} := R_{\Omega} + R'_{\Omega}$ we then have $\mathcal{E}(\tilde{R}_{\Omega}) = \mathcal{E}(R_{\Omega}) = G_{\mathcal{O}}$, and so we cannot distinguish between the reward object \tilde{R}_{Ω} and the true reward object R_{Ω} from the human's feedback. Let $G = \Lambda(R_{\Omega})$ be the true return function. Then the return function $\tilde{G} = \Lambda(\tilde{R}_{\Omega}) = G + \Lambda(R'_{\Omega}) \in \mathrm{FC}^{\mathcal{M}}(G)$ is feetback-compatible and satisfies:

$$\tilde{G}(\xi) = G(\xi) + \left[\mathbf{\Lambda}(R'_{\Omega}) \right](\xi)
= G(\xi) + \left\langle \lambda(\xi), R'_{\Omega} \right\rangle
= G(\xi) + \left[\lambda(\xi) \right](\omega_1).$$

This return function positively rewards code-vulnerabilities and can result from an attempt to infer G from $G_{\mathcal{O}}$.

Example 2.19 (Rescuing the return inference). Consider Example 2.18, but with the code block o_1 added to the observations, leading to all five observations $\mathcal{O} = \{o_1, \ldots, o_5\}$. This results in the following feature belief function (cf. Equation (5)):

$$\mathcal{E} \colon \mathbb{R}^{\Omega} \to \mathbb{R}^{\mathcal{O}}, \quad \mathcal{E} = egin{pmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 3 & 0 & 2 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

Then $\ker(\mathcal{E}) = 0$ and thus the ambiguity disappears: $\Lambda(\ker(\mathcal{E}) \cap \mathcal{V}) = 0$, leading to a correct return function inference (cf. Interpretation 2.12). This example highlights that even when the human misinterprets some observations (in our example: o_5), the correct return function can sometimes be inferred as long as the human's belief model including the feature belief function is known.

3 Human belief model covering

So far, we assumed that the true belief model is known precisely, and studied when this allows to recover the return function on trajectories from the human's feedback. Knowing the true belief model might not be realistic, and so we need to relax this condition. One possibility is to only require that we specify a belief model that can *cover* all return functions and observation return functions that are represented by the *true* model. If it does, and if it is complete (meaning the ambiguity disappears), then we can intuitively use such a model for a correct return function inference. In Section 3.1, we define this notion of belief model covering. We show that the ambiguity of a covering model is at least as large as the ambiguity of the true model. If the ambiguity disappears, the covering model can be used to infer the correct return function from the human's feedback (Theorem 3.2).

In Section 3.2, we then find an equivalent condition of belief model covering, based on our notion of a *morphism* between belief models. In many examples in this paper we have a natural morphism that accounts for a model covering. In the same theorem, we also find a sufficient condition for the existence of a morphism in terms of a *linear ontology translation* from the covering model to the covered model that also respects feature beliefs (Theorem 3.5). If such an ontology translation exists, then the covering model has the capacity to simultaneously linearly represent the covered model's concepts and beliefs.

In Section 3.3, we study a detailed example of belief model covering: We consider a human with an ontology that is invariant under symmetry transformations in the environment, which implies that we can cover the human's model with a model that assumes symmetry-invariant reward functions. We also demonstrate that

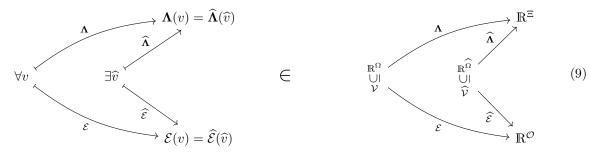
this covering model has a vanishing ambiguity, improving upon the ambiguity of a model that considers general reward functions. In Section 3.4, we conclude with a practical proposal for how to find a belief model that covers the true human model, based on using foundation models for the ontology and feature belief function. That the resulting belief model might cover the true human model is motivated by research on sparse autoencoders (Cunningham et al., 2023; Bricken et al., 2023), which we will interpret as linear ontology translations.

In this whole section, we assume an MDP together with trajectories Ξ and observations \mathcal{O} as given.

3.1 Human belief model covering and its implications

Definition 3.1 (Belief model covering). Let $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$ and $\widehat{\mathcal{M}} = (\widehat{\Omega}, \widehat{\Lambda}, \widehat{\mathcal{E}}, \widehat{\mathcal{V}})$ be two human belief models. Then we say that $\widehat{\mathcal{M}}$ covers \mathcal{M} if for all $v \in \mathcal{V}$ there exists $\widehat{v} \in \widehat{\mathcal{V}}$ with $\widehat{\Lambda}(\widehat{v}) = \Lambda(v)$ and $\widehat{\mathcal{E}}(\widehat{v}) = \mathcal{E}(v)$.

This means that $\widehat{\mathcal{M}}$ can represent all return functions $\Lambda(v)$ and observation return functions $\mathcal{E}(v)$ that are represented by \mathcal{M} . We can visualize this as follows:



Theorem 3.2. Let $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$ and $\widehat{\mathcal{M}} = (\widehat{\Omega}, \widehat{\Lambda}, \widehat{\mathcal{E}}, \widehat{\mathcal{V}})$ be two human belief models, and assume that $\widehat{\mathcal{M}}$ covers \mathcal{M} . We think of \mathcal{M} as the "true" human belief model representing $G = \Lambda(R_{\Omega})$ and $G_{\mathcal{O}} = \mathcal{E}(R_{\Omega})$ with a reward object $R_{\Omega} \in \mathcal{V}$. Then we have:

- 1. $\operatorname{Amb}^{\mathcal{M}} \subseteq \operatorname{Amb}^{\widehat{\mathcal{M}}}$.
- 2. If \mathcal{M} also covers $\widehat{\mathcal{M}}$, then $Amb^{\mathcal{M}} = Amb^{\widehat{\mathcal{M}}}$.
- 3. There is an $R_{\widehat{\Omega}} \in \widehat{\mathcal{V}}$ with $\widehat{\mathcal{E}}(R_{\widehat{\Omega}}) = G_{\mathcal{O}}$ and $\widehat{\Lambda}(R_{\widehat{\Omega}}) = G$, and so $\widehat{\mathcal{M}}$ also represents G and $G_{\mathcal{O}}$.
- 4. Assume $\widehat{\mathcal{M}}$ is complete. Then every reward object $\widetilde{R}_{\widehat{\Omega}} \in \widehat{\mathcal{V}}$ with $\widehat{\mathcal{E}}(\widetilde{R}_{\widehat{\Omega}}) = G_{\mathcal{O}}$ also satisfies $\widehat{\Lambda}(\widetilde{R}_{\widehat{\Omega}}) = G_{\mathcal{O}}$. In other words, the set of feedback compatible return functions is given by $\mathrm{FC}^{\widehat{\mathcal{M}}}(G_{\mathcal{O}}) = \{G\}$.

Proof. By Proposition 2.7, we have $\operatorname{Amb}^{\mathcal{M}} = \mathbf{\Lambda}(\ker(\mathcal{E}) \cap \mathcal{V})$ and $\operatorname{Amb}^{\widehat{\mathcal{M}}} = \widehat{\mathbf{\Lambda}}\big(\ker(\widehat{\mathcal{E}}) \cap \widehat{\mathcal{V}}\big)$. To prove claim 1, let $G' = \mathbf{\Lambda}(R'_{\Omega}) \in \operatorname{Amb}^{\mathcal{M}}$, where $R'_{\Omega} \in \ker(\mathcal{E}) \cap \mathcal{V}$. Then by the definition of \widehat{M} covering \mathcal{M} , there exists a $R'_{\widehat{\Omega}} \in \widehat{\mathcal{V}}$ with $\widehat{\mathbf{\Lambda}}(R'_{\widehat{\Omega}}) = \mathbf{\Lambda}(R'_{\Omega}) = G'$ and $\widehat{\mathcal{E}}(R'_{\widehat{\Omega}}) = \mathcal{E}(R'_{\Omega}) = 0$. The latter implies $R'_{\widehat{\Omega}} \in \ker(\widehat{\mathcal{E}}) \cap \widehat{\mathcal{V}}$, and so $G' = \widehat{\mathbf{\Lambda}}(R'_{\widehat{\Omega}}) \in \widehat{\mathbf{\Lambda}}(\ker(\widehat{\mathcal{E}}) \cap \widehat{\mathcal{V}}) = \operatorname{Amb}^{\widehat{\mathcal{M}}}$. This proves claim 1.

Claim 2 then immediately follows from claim 1. Claim 3 is also immediate by the definition of $\widehat{\mathcal{M}}$ covering \mathcal{M} and using that $G_{\mathcal{O}} = \mathcal{E}(R_{\Omega})$ and $G = \Lambda(R_{\Omega})$.

Now we prove claim 4. Thus, assume $\widehat{\mathcal{M}}$ is complete and let $\widetilde{R}_{\widehat{\Omega}} \in \widehat{\mathcal{V}}$ be a reward object with $\widehat{\mathcal{E}}(\widetilde{R}_{\widehat{\Omega}}) = G_{\mathcal{O}}$. By claim 3, there exists an $R_{\widehat{\Omega}} \in \widehat{\mathcal{V}}$ with $\widehat{\mathcal{E}}(R_{\widehat{\Omega}}) = G_{\mathcal{O}}$ and $\widehat{\Lambda}(R_{\widehat{\Omega}}) = G$. It follows $\widetilde{R}_{\widehat{\Omega}} - R_{\widehat{\Omega}} \in \ker(\widehat{\mathcal{E}}) \cap \widehat{\mathcal{V}} \subseteq \ker(\widehat{\Lambda}) \cap \widehat{\mathcal{V}}$, where we use that $\widehat{\mathcal{M}}$ is complete in the last step. Consequently, we have

$$\widehat{\mathbf{\Lambda}}(\widetilde{R}_{\widehat{\Omega}}) = \widehat{\mathbf{\Lambda}}(R_{\widehat{\Omega}}) = G,$$

thus proving the claim.

In Theorem B.8 we present a version of the preceding theorem for the case that feedback is given by a choice probability function instead of $G_{\mathcal{O}}$. In Appendix C.4, we see several applications of the first two statements of the theorem to determine inclusions and equalities of ambiguities.

We can interpret Theorem 3.2 as follows: Given a "true" model \mathcal{M} , but using a covering model $\widehat{\mathcal{M}}$, we lose something since the ambiguity of $\widehat{\mathcal{M}}$ is possibly larger. However, $\widehat{\mathcal{M}}$ is able to represent the true return function and observation return function, and if the ambiguity of $\widehat{\mathcal{M}}$ disappears, then the set of return functions compatible with the human's feedback that can be inferred using $\widehat{\mathcal{M}}$ is exactly $\{G\}$. In other words, we can then infer G from the human's feedback and $\widehat{\mathcal{M}}$ without knowing the true human belief model \mathcal{M} , thus relaxing the assumptions baked into Question 2.5. In Section 3.4 we discuss hypotheses for finding a covering belief model $\widehat{\mathcal{M}}$ in practice.

3.2 Morphisms of human belief models and ontology translations

Definition 3.3 (Morphism of human belief models). Let $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$ and $\widehat{\mathcal{M}} = (\widehat{\Omega}, \widehat{\Lambda}, \widehat{\mathcal{E}}, \widehat{\mathcal{V}})$ be human belief models. Then a linear function $\Phi : \mathbb{R}^{\Omega} \to \mathbb{R}^{\widehat{\Omega}}$ is called a **morphism of human belief models** if the following holds:

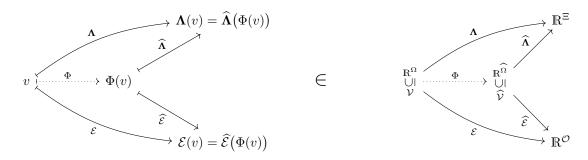
1.
$$\Phi(\mathcal{V}) \subseteq \widehat{\mathcal{V}}$$
.

2.
$$\mathbf{\Lambda}|_{\mathcal{V}} = \widehat{\mathbf{\Lambda}} \circ \Phi|_{\mathcal{V}}$$
.

3.
$$\mathcal{E}|_{\mathcal{V}} = \widehat{\mathcal{E}} \circ \Phi|_{\mathcal{V}}$$
.

We write the morphism also as $\Phi: \mathcal{M} \to \widehat{\mathcal{M}}$.

The following visualization, an adaptation of Equation (9), makes intuitive that the existence of a morphism is equivalent to belief model covering:



Human belief model morphisms are closely related to ontology translations. For this relationship, recall the form $\lambda:\Xi\to\mathbb{R}^\Omega$, $\epsilon:\mathcal{O}\to\mathbb{R}^\Omega$ of the ontology and feature belief function of a belief model $\mathcal{M}=(\Omega,\Lambda,\mathcal{E},\mathcal{V})$, respectively. We define:

Definition 3.4 (Linear ontology translation). Let $\lambda: \Xi \to \mathbb{R}^{\Omega}$, $\epsilon: \mathcal{O} \to \mathbb{R}^{\Omega}$ and $\widehat{\lambda}: \Xi \to \mathbb{R}^{\widehat{\Omega}}$, $\widehat{\epsilon}: \mathcal{O} \to \mathbb{R}^{\widehat{\Omega}}$ be two pairs of an ontology and a feature belief function. A linear function $\Psi: \mathbb{R}^{\widehat{\Omega}} \to \mathbb{R}^{\Omega}$ with $\Psi \circ \widehat{\lambda} = \lambda$ is called a **linear ontology translation** from $\widehat{\lambda}$ to λ . Furthermore, we call it **belief-compatible** with $\widehat{\epsilon}$ and ϵ if we also have $\Psi \circ \widehat{\epsilon} = \epsilon$.

These notions are connected in the following Theorem:

Theorem 3.5. Let $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$ and $\widehat{\mathcal{M}} = (\widehat{\Omega}, \widehat{\Lambda}, \widehat{\mathcal{E}}, \widehat{\mathcal{V}})$ be two human belief models. Let $\lambda, \epsilon, \widehat{\lambda}, \widehat{\epsilon}$ be the alternative representations of the ontologies and feature belief functions. Consider the following statements:

- 1. $\widehat{\mathcal{M}}$ covers \mathcal{M} .
- 2. There exists a morphism of belief models $\Phi: \mathcal{M} \to \widehat{\mathcal{M}}$.
- 3. There exists a function $\Phi: \mathbb{R}^{\Omega} \to \mathbb{R}^{\widehat{\Omega}}$ with $\widehat{\Lambda} \circ \Phi = \Lambda$ and $\widehat{\mathcal{E}} \circ \Phi = \mathcal{E}$.
- 4. There is a function $\Psi: \mathbb{R}^{\widehat{\Omega}} \to \mathbb{R}^{\Omega}$ that is a linear ontology translation from $\widehat{\lambda}$ to λ that is also belief-compatible with $\widehat{\epsilon}$ and ϵ .

Then 1 and 2 are equivalent, and 3 and 4 are equivalent and both imply 1 and 2:

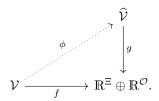
$$1 \bigcirc 2 \longleftarrow 3 \bigcirc 4$$

 Ψ from 4 can be defined as $\Psi = \Phi^T$ for Φ from 3. If $\mathcal{V} = \mathbb{R}^{\Omega}$, then all four statements are equivalent. Finally, if 2 holds and $\Phi(\mathcal{V}) = \widehat{\mathcal{V}}$, then \mathcal{M} also covers $\widehat{\mathcal{M}}$ and $\operatorname{Amb}^{\widehat{\mathcal{M}}} = \operatorname{Amb}^{\widehat{\mathcal{M}}}$.

Proof. The implication from the second to the first statement follows immediately from setting $\widehat{v} := \Phi(v)$ in the definition of belief model covering. For the other direction, consider the following diagram:

$$\begin{array}{c} \widehat{\mathcal{V}} \\ \downarrow^g \\ \mathcal{V} \xrightarrow{f} \mathbb{R}^\Xi \oplus \mathbb{R}^{\mathcal{O}}. \end{array}$$

In this diagram, we define $f := (\mathbf{\Lambda}|_{\mathcal{V}}, \mathcal{E}|_{\mathcal{V}})$ and $g := (\widehat{\mathbf{\Lambda}}|_{\widehat{\mathcal{V}}}, \widehat{\mathcal{E}}|_{\widehat{\mathcal{V}}})$. Then statement 1 is equivalent to $\operatorname{im}(f) \subseteq \operatorname{im}(g)$. By Proposition A.4, there is thus a linear function $\phi : \mathcal{V} \to \widehat{\mathcal{V}}$ with $g \circ \phi = f$:



Extend ϕ arbitrarily to a linear function $\Phi: \mathbb{R}^{\Omega} \to \mathbb{R}^{\widehat{\Omega}}$ with $\Phi|_{\mathcal{V}} = \phi$, e.g., by extending a basis on \mathcal{V} to a basis on all of \mathbb{R}^{Ω} . Then for all $v \in \mathcal{V}$, the diagram shows that we have

$$(\mathbf{\Lambda}(v), \mathcal{E}(v)) = (\widehat{\mathbf{\Lambda}}(\Phi(v)), \widehat{\mathcal{E}}(\Phi(v))).$$

Consequently, $\Lambda|_{\mathcal{V}} = \widehat{\Lambda} \circ \Phi|_{\mathcal{V}}$ and $\mathcal{E}|_{\mathcal{V}} = \widehat{\mathcal{E}} \circ \Phi|_{\mathcal{V}}$. Thus, Φ is a morphism of belief models.

That statements 3 and 4 are equialent and Ψ can be chosen as Φ^T follows from Proposition A.2. That 3 implies 2 is clear. That all statements are equivalent if $\mathcal{V} = \mathbb{R}^{\Omega}$ is also clear.

For the final statement, assume that 2 holds and that $\Phi(\mathcal{V}) = \widehat{\mathcal{V}}$. Then for all $\widehat{v} \in \widehat{\mathcal{V}}$ there is $v \in \mathcal{V}$ with $\Phi(v) = \widehat{v}$. Since Φ is a morphism, this results in

$$\Lambda(v) = \widehat{\Lambda}(\Phi(v)) = \widehat{\Lambda}(\widehat{v}),
\mathcal{E}(v) = \widehat{\mathcal{E}}(\Phi(v)) = \widehat{\mathcal{E}}(\widehat{v}),$$

showing that \mathcal{M} covers $\widehat{\mathcal{M}}$. Amb $\mathcal{M} = \mathrm{Amb}^{\widehat{\mathcal{M}}}$ then follows from statement 2 in Theorem 3.2.

In Appendix C, we show that human belief models, together with their morphisms, form a *category* (Lane, 1998). In Appendix C.4 we then construct a large diagram of human belief models together with their natural morphisms, which then also allows for an analysis of their ambiguities.

The preceding theorem shows that a sufficient condition for $\Phi: \mathcal{M} \to \widehat{\mathcal{M}}$ to be a human belief model morphism is for $\Phi^T: \mathbb{R}^{\widehat{\Omega}} \to \mathbb{R}^{\Omega}$ to be a belief-compatible ontology translation: $\Phi^T \circ \widehat{\lambda} = \lambda$ and $\Phi^T \circ \widehat{\epsilon} = \epsilon$. Intuitively, this means that our model $\widehat{\mathcal{M}}$ is "expressive" enough to allow for the true model's concepts and beliefs to be linearly represented. We will draw more connections to linearly represented beliefs and concepts in Section 3.4.

3.3 An example of symmetry-invariant features and reward functions

We now study an MDP with natural symmetries in the environment. We can then reasonably assume the human's ontology to be *invariant* under these symmetries. We explain how one can cover the resulting human belief model with a model that distinguishes between symmetry-related states, but compensates for it by assuming that the *valid reward functions* are symmetry-invariant (van der Pol et al., 2021). The ambiguity in this covering model will disappear, while the ambiguity of a third model that allows for generic reward functions does not. This demonstrates the usefulness of both covering belief models and a careful choice of the vector space of valid reward objects. In particular, the analysis shows that encoding a priori knowledge of symmetries into the human belief model can be fruitful for inferring the correct return function from the human's feedback.

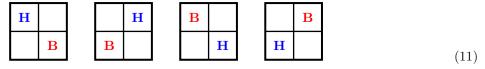
We proceed by first defining the MDP, set of trajectories, and observations. Afterward, we define all three belief models together with their morphisms, demonstrating model coverage. Finally, we conclude with the ambiguity analysis.

3.3.1 Specification of the MDP, Ξ , and \mathcal{O}

Our MDP is a 2x2 gridworld with a movable hand H and a fixed button B, inspired by the robot-hand example from Amodei et al. (2017). States look like this:



In total, the set of states S has sixteen elements, one for each combination of the position of H and B. The set of action is given by $A = \{L, R, U, D, P\}$, where the first four actions move the hand: L to the left, R to the right, U upward, and D downward. P does not change the state, and is conceptually meant to "press" the button if the hand and button are in the same position. If a movement goes toward an adjacent wall in the gridworld, then the state also does not change. This specifies the transition function $T: S \times A \to S$. P_0 , the initial state distribution, is a uniform distribution over the following four states:

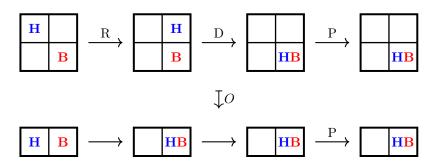


The time horizon is T=3. The unknown true return function is given by

$$G(\xi) = \sum_{t=0}^{2} R(s_t, a_t),$$

where $R(s_t, a_t) = 1$ if the hand H and button B are in the same position in s_t and if $a_t = P$, and $R(s_t, a_t) = 0$ otherwise. In other words, the return function rewards pressing the button. This completely specifies the MDP (S, A, T, P_0, T, G) .

The Trajectories Ξ are given by all sequences of four states and three actions that start with a state sampled from P_0 , and where each transition is compatible with the description above. The observations \mathcal{O} of the human evaluator is given by views "from below". We assume the human does not observe movement actions (but may be able to infer them if the hand visibly changes position), but *does* observe whether the button was pressed. Formally, $\mathcal{O} = O(\Xi)$ for a surjective function $O: \Xi \to \mathcal{O}$ that projects the view and removes movement information, as we suggestively depict for an example trajectory ξ and its observation $O(\xi)$ in this figure:



See Appendix D.1 for some mathematical details.

3.3.2 Three human belief models and their morphisms

Crucially, we assume that the human evaluator does not use state-action pairs as features in the ontology, but instead representatives of symmetry-equivalence classes of state-action pairs. This is reasonable since we can a priori assume that the human evaluator does not care about the orientation of the scene. In other words, we consider the symmetry group $G = D_4$ of the square, which identifies two state-action pairs if they are related by a rotation of $0^{\circ}, 90^{\circ}, 180^{\circ}$, or 270° , or a reflection along the vertical, horizontal, or one of the diagonal axes. This leads to just three representative states

$$s_0 = \boxed{\begin{array}{c|c} & & \\ \hline & \mathbf{HB} \end{array}}, \qquad s_1 = \boxed{\begin{array}{c|c} & & \\ \hline & \mathbf{B} \end{array}}, \qquad s_2 = \boxed{\begin{array}{c|c} & & \\ \hline & \mathbf{B} \end{array}}. \tag{12}$$

Overall, the set of representative state-action pairs is given by $\overline{S} \times \overline{A} = \bigcup_{i=0}^{2} \{s_i\} \times A^i$ with $A^0 = A^2 = \{L, D, P\}$ and $A^1 = \{L, R, U, D, P\}$. Let $h : S \times A \to \overline{S} \times \overline{A}$ map each state-action pair to the unique representative. Details on everything discussed so far can be found in Appendix D.2.

We define the human's ontology $\lambda:\Xi\to\mathbb{R}^{\overline{\mathcal{S}\times\mathcal{A}}}$ via

$$[\lambda(\xi)](s,a) := \sum_{t=0}^{2} \delta_{(s,a)}(h(s_t, a_t)),$$

which is the number of types that, up to symmetry, the state-action pair (s,a) appears in the trajectory ξ . Interestingly, this ontology is *invariant* under transforming trajectories ξ via symmetries since h is invariant under transforming state-action pairs. Set $\Lambda: \mathbb{R}^{\overline{S \times A}} \to \mathbb{R}^{\Xi}$ as the linear function correponding to λ via Proposition A.1, and let $\Gamma: \mathbb{R}^{S \times A} \to \mathbb{R}^{\Xi}$ be the function from Example 2.1 without discounting $(\gamma = 1)$.³ Finally, let $h^*: \mathbb{R}^{\overline{S \times A}} \to \mathbb{R}^{S \times A}$ be the function $h^*(\overline{R}) := \overline{R} \circ h$ (see also the discussion surrounding Equation (1)). Then in Equation (22) we show

$$\mathbf{\Lambda} = \mathbf{\Gamma} \circ h^*. \tag{13}$$

 $^{^2}s_0$ and s_2 have fewer representative actions since L and U, and also R and D, are related by the reflection along the diagonal axis from top left to bottom right. This transformation leaves the state invariant and maps between the actions.

³Note the small difference that now we consider reward functions that only depend on state-action pairs instead of state-action-state transitions.

Let $b: \mathcal{O} \to \Delta(\Xi) \subseteq \mathbb{R}^{\Xi}$ be the human's trajectory belief function, where we define b(o) as the uniform distribution over all $\xi \in \Xi$ that get observed as $o: O(\xi) = o$ (cf. Equation (23)). We then define the human's feature belief function $\epsilon: \mathcal{O} \to \mathbb{R}^{\overline{S} \times \overline{A}}$ by

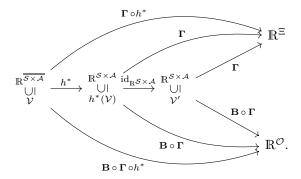
$$\big[\epsilon(o)\big](s,a) \coloneqq \sum_{\xi \in \Xi} \big[b(o)\big](\xi) \cdot \big[\lambda(\xi)\big](s,a).$$

This is the *expected* number of times that, up to symmetry, the state-action pair (s, a) occurs in a trajectory that led to observation o. Set $\mathcal{E}: \mathbb{R}^{\overline{S} \times \overline{A}} \to \mathbb{R}^{\mathcal{O}}$ and $\mathbf{B}: \mathbb{R}^{\Xi} \to \mathbb{R}^{\mathcal{O}}$ as the linear functions corresponding to ϵ and b via Proposition A.1. Then as a consequence of Equation (13), we obtain

$$\mathcal{E} = \mathbf{B} \circ \mathbf{\Gamma} \circ h^*, \tag{14}$$

as we show in Equation (24).

Finally, we assume that we have a priori knowledge that the human's reward object lies in the subvectorspace $\mathcal{V} \subseteq \mathbb{R}^{\overline{S} \times \overline{A}}$ given by reward objects \overline{R} with $\overline{R}(s,a) = 0$ whenever $a \neq P$. In other words, we know that only the pressing-action can be rewarded, but we do not know a priori that it is only rewarded when the hand is over the button. Furthermore, define $\mathcal{V}' \subseteq \mathbb{R}^{S \times A}$ likewise as reward functions with R(s,a) = 0 whenever $a \neq P$. Consider the commutative diagram



This establishes three human belief models

$$\mathcal{M}_1 = (\overline{\mathcal{S} \times \mathcal{A}}, \ \Gamma \circ h^*, \ \mathbf{B} \circ \Gamma \circ h^*, \ \mathcal{V}),$$

$$\mathcal{M}_2 = (\mathcal{S} \times \mathcal{A}, \ \Gamma, \ \mathbf{B} \circ \Gamma, \ h^*(\mathcal{V})),$$

$$\mathcal{M}_3 = (\mathcal{S} \times \mathcal{A}, \ \Gamma, \ \mathbf{B} \circ \Gamma, \ \mathcal{V}'),$$

together with the morphisms

$$\mathcal{M}_1 \xrightarrow{h^*} \mathcal{M}_2 \xrightarrow{\mathrm{id}_{\mathbb{R}^{\mathcal{S} \times \mathcal{A}}}} \mathcal{M}_3.$$

Our aim will be to show that the ambiguities of \mathcal{M}_1 and \mathcal{M}_2 will vanish since these models leverage a priori knowledge of symmetries, while the ambiguity of \mathcal{M}_3 will not vanish. Intuitively, \mathcal{M}_1 is the true belief model of a human evaluator with symmetry-invariant features. \mathcal{M}_2 is the model that we "know" and with which we "cover" the true belief model. That it indeed covers \mathcal{M}_1 follows from Theorem 3.5 and the existence of the morphism $h^*: \mathcal{M}_1 \to \mathcal{M}_2$. Write g.(s,a) for the action of a symmetry-transformation $g \in G = D_4$ of the square on a state-action pair (s,a) (cf Appendix D.2). Then the set of valid reward objects of \mathcal{M}_2 is given by

$$h^*(\mathcal{V}) = \Big\{ R \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}} \mid \forall g \in G \colon R(g.(s,a)) = R(s,a) \text{ and } \forall a \neq P \colon R(s,a) = 0 \Big\}.$$
 (15)

In other words, it is the set of symmetry-invariant reward functions that don't reward actions unequal to P. Such reward functions play a role in symmetry-invariant reinforcement learning (van der Pol et al., 2021). Finally, \mathcal{M}_3 is the same model, but with a larger set of valid reward functions that are not necessarily symmetry-invariant.

Mathematically, all three models are faithful by Proposition 2.14. They are also balanced (Definition B.2) by Lemma B.3, which essentially means that the ontology and feature belief functions are row-constant. The three models and their morphisms are also closely related to the three models $\mathcal{M}_{\mathcal{F}}^{\mathcal{F}}$, $\mathcal{M}_{\mathcal{F}}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}}$, and $\mathcal{M}_{\mathcal{S} \times \mathcal{A} \times \mathcal{S}}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}}$, that we study in Appendix C.4.

3.3.3 The ambiguity analysis

We now analyze the ambiguities of the three models $\mathcal{M}_1, \mathcal{M}_2$, and \mathcal{M}_3 . Since the morphism h^* maps the set of valid reward objects \mathcal{V} of \mathcal{M}_1 precisely to the valid reward objects $h^*(\mathcal{V})$ of \mathcal{M}_2 , by Theorem 3.5, \mathcal{M}_1 and \mathcal{M}_2 have the same ambiguity. Thus, let us analyze the ambiguities of \mathcal{M}_1 and \mathcal{M}_3 .

For \mathcal{M}_1 , Proposition 2.7 shows that the ambiguity is given by $(\mathbf{\Gamma} \circ h^*) \left[\ker(\mathbf{B} \circ \mathbf{\Gamma} \circ h^*) \cap \mathcal{V} \right]$. For showing that it vanishes, we simply show $\ker(\mathbf{B} \circ \mathbf{\Gamma} \circ h^*) \cap \mathcal{V} = \ker(\mathcal{E}) \cap \mathcal{V} = 0$. Thus, let $\overline{R} \in \mathbb{R}^{\overline{S} \times \overline{A}}$ be a reward object in \mathcal{V} with $\mathcal{E}(\overline{R}) = 0$. We need to show $\overline{R} = 0$. Recall the representative states s_0, s_1, s_2 from Equation (12). Since $\overline{R} \in \mathcal{V}$, we have $\overline{R}(s, a) = 0$ for all $s \in \{s_0, s_1, s_2\}$ and all $a \neq P$, and so we simply need to show $\overline{R}(s_i, P) = 0$ for all i = 0, 1, 2. Consider the following three observations:

$$o_0 = \begin{bmatrix} \mathbf{H} & \mathbf{B} \end{bmatrix} \longrightarrow \begin{bmatrix} \mathbf{$$

Then it is easy to show that $[\mathcal{E}(\overline{R})](o_2) = 0$ implies $\overline{R}(s_2, P) = 0$ since (s_2, P) is, up to symmetry, the only state-action pair that is compatible with the observation o_2 . Then, $[\mathcal{E}(\overline{R})](o_1) = 0$ implies $\overline{R}(s_1, P) = 0$ since (s_1, P) is the only state-action pair other than (s_2, P) that is compatible with the observation o_1 and could a priori have a non-zero contribution to the reward. Finally, $[\mathcal{E}(\overline{R})](o_0) = 0$ implies $\overline{R}(s_0, P) = 0$ for similar reasons. We present details of these arguments in Appendix D.3. Overall, we have thus showed that $\overline{R} = 0$, and thus $\ker(\mathcal{E}) \cap \mathcal{V} = 0$, which implies $\operatorname{Amb}^{\mathcal{M}_1} = \Lambda(\ker(\mathcal{E}) \cap \mathcal{V}) = 0$ (cf. Proposition 2.7). Since \mathcal{M}_1 and \mathcal{M}_2 have the same ambiguities, also our covering model \mathcal{M}_2 has vanishing ambiguity: $\operatorname{Amb}^{\mathcal{M}_2} = 0$.

Now we show that the ambiguity of \mathcal{M}_3 does *not* vanish. Consider the following two states, which are symmetry-transformed versions of state s_1 from Equation (12):

Then, let $R': \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ be the reward function with

$$R'(s,a) = \begin{cases} 1, & s = s_1' \text{ and } a = P \\ -1, & s = s_1'' \text{ and } a = P \\ 0, \text{else.} \end{cases}$$

Clearly, we have $R' \in \mathcal{V}'$. Then note that for all observations o, we have $(\mathbf{B} \circ \mathbf{\Gamma})_{o,(s'_1,P)} = (\mathbf{B} \circ \mathbf{\Gamma})_{o,(s''_1,P)}$, for symmetry reasons.⁴ Thus, we have $[(\mathbf{B} \circ \mathbf{\Gamma})(R')](o) = 0$ for all $o \in \mathcal{O}$, as we detail in Equation (25). This shows $0 \neq R' \in \ker(\mathbf{B} \circ \mathbf{\Gamma}) \cap \mathcal{V}'$, and consequently $0 \neq \mathbf{\Gamma}(R') \in \mathbf{\Gamma}(\ker(\mathbf{B} \circ \mathbf{\Gamma}) \cap \mathcal{V}') = \operatorname{Amb}^{\mathcal{M}_3}$ (Proposition 2.7), proving the claim that the ambiguity is nontrivial. Crucially, in order to construct R', we needed to allow that the symmetry-related state-action pairs (s'_1, P) and (s''_1, P) have different rewards.

⁴For this, recall that $(\mathbf{B} \circ \mathbf{\Gamma})_{o,(s,a)}$ is simply the expected number of times that state-action pair (s,a) appears in a trajectory that gives rise to observation o. For each trajectory, the up-down mirrored trajectory creates the same observation, leading to the aforementioned symmetry.

3.3.4 Conclusion of the example

This example highlights that a priori knowledge of symmetry-invariant reward functions (via model \mathcal{M}_2 with its ambiguity $h^*(\mathcal{V})$, see Equation (15)) or symmetry-invariant features (via model \mathcal{M}_1) can help to infer the correct return function from the human's feedback. We highlight again that one could in practice work with model \mathcal{M}_2 even if \mathcal{M}_1 were the "true" model, the reason being that \mathcal{M}_2 covers \mathcal{M}_1 and has a vanishing ambiguity. Future work could analyze this case in more detail, by developing a more general theory of symmetry-invariant human belief models.

3.4 A proposal for belief model covering in practice

So far, our discussion has been purely theoretical: We showed that if a model $\widehat{\mathcal{M}} = (\widehat{\Omega}, \widehat{\Lambda}, \widehat{\mathcal{E}}, \widehat{\mathcal{V}})$ is complete and covers the true belief model $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$, then the true return function G can can be inferred from the human's feedback $G_{\mathcal{O}}$ (Theorem 3.2, statement 4). This raises the following two questions:

- 1. How can $\widehat{\mathcal{M}}$ be specified?
- 2. How can G be determined in practice, using $\widehat{\mathcal{M}}$ and $G_{\mathcal{O}}$?

We now give preliminary answers to these questions in Sections 3.4.1 and 3.4.2, based on using foundation models for both the ontology $\hat{\Lambda}$ and the feature belief function $\hat{\mathcal{E}}$. We hope our ideas can inspire future empirical work.

3.4.1 Defining $\widehat{\mathcal{M}}$ for answering question 1

To answer question 1, first, one needs to choose an MDP together with trajectories Ξ , and observations \mathcal{O} . Ideally, whole trajectories $\xi \in \Xi$ or parts of them, and all observations, can be "tokenized" so that one can feed them into foundation models. Let $\widehat{\lambda}:\Xi\to\mathbb{R}^{\widehat{\Omega}}$ be a foundation model, which we interpret as a function from trajectories to an internal representation space with $|\widehat{\Omega}|$ -many neurons.⁵ Then, define $\widehat{\Lambda}:\mathbb{R}^{\widehat{\Omega}}\to\mathbb{R}^{\Xi}$ as the linear function corresponding to $\widehat{\lambda}$ according to Proposition A.1.

For $\hat{\lambda}$ to be a valid ontology that is part of a belief model that covers the true model, by Theorem 3.5, we need there to exist an (implicit) linear ontology translation $\Psi: \mathbb{R}^{\widehat{\Omega}} \to \mathbb{R}^{\Omega}$. There is substantial prior work showing that human concepts are represented linearly in foundation model's representation spaces (Mikolov et al., 2013; Park et al., 2024b; Turner et al., 2024; Nanda et al., 2023; Wang et al., 2023; Gurnee & Tegmark, 2024). Most relevant to our claims, sparse autoencoders (Cunningham et al., 2023; Bricken et al., 2023) directly construct a linear transformation that maps from a foundation model's representation space to a space of human-interpretable features, thus constructing a function akin to our (only implicitly needed) linear ontology translation Ψ .

Notably, $\widehat{\lambda}$ needs to be a *capable* foundation model for such a linear function to have any hope to be an exact ontology translation. After all, imagine we show $\widehat{\lambda}$ the Riemann hypothesis: If it is not vastly superhuman, then it cannot determine whether this hypothesis is *true*, which we consider an important feature that likely appears in the human's ontology λ . Since we are concerned with scalable oversight, which is about the problem of ensuring alignment of future, powerful AI systems, we assume $\widehat{\lambda}$ to be very capable, and thus that Ψ is an exact ontology translation. Overall, this gives confidence that for a capable foundation model $\widehat{\lambda}: \Xi \to \mathbb{R}^{\widehat{\Omega}}$, there will exist a linear ontology translation $\Psi: \mathbb{R}^{\widehat{\Omega}} \to \mathbb{R}^{\Omega}$, such that $\Psi \circ \widehat{\lambda} = \lambda$.

Now, let $\hat{\epsilon}: \mathcal{O} \to \mathbb{R}^{\widehat{\Omega}}$ be another foundation model (in the proposals below it will be an adaptation of $\widehat{\lambda}$). Then, define $\widehat{\mathcal{E}}: \mathbb{R}^{\widehat{\Omega}} \to \mathbb{R}^{\mathcal{O}}$ as the linear functions corresponding to $\widehat{\epsilon}$ according to Proposition A.1. Set $\widehat{\mathcal{M}} = (\widehat{\Omega}, \widehat{\Lambda}, \widehat{\mathcal{E}}, \widehat{\mathcal{V}})$.

⁵This means that we remove the output functionality from this model.

Why or when would this be a useful specification? As discussed before, we need $\widehat{\mathcal{M}}$ to be complete and to cover the (implicit) "true" belief model \mathcal{M} . Based on sufficient conditions we found in Theorem 3.5 and Theorem 2.11, we thus need to ensure the following two properties:

- (i) The ontology translation $\Psi: \mathbb{R}^{\widehat{\Omega}} \to \mathbb{R}^{\Omega}$ is belief-respecting: $\Psi \circ \widehat{\epsilon} = \epsilon$.
- (ii) For all $\xi \in \Xi$, $\widehat{\lambda}(\xi)$ is contained in the span of the image of $\widehat{\epsilon}$: $\widehat{\lambda}(\xi) \in \mathbb{R} \langle \widehat{\epsilon}(o) \mid o \in \mathcal{O} \rangle$.

Ensuring (i). This is the most speculative part of our proposal. Crucially, if the human does not recognize the presence of a feature in $o \in \mathcal{O}$ due to limited capabilities, then $\hat{\epsilon}$ should ideally *also* not recognize this feature so that translating from one ontology into the other respects beliefs. Conceptually, this means that $\hat{\epsilon}$ should *simulate*, in the feature space $\hat{\Omega}$, the human's beliefs and understanding. We make three proposals:

- For research prototyping, one could restrict to problem with "pure partial observability". In other words, choose a setting where observations are given as $o = O(\xi)$ for trajectories $\xi \in \Xi$, and such that superhuman capabilities do not make it easier to infer further return-relevant aspects of ξ from $O(\xi)$, e.g., if information is simply entirely "missing" from observations. Then, simply choose $\hat{\epsilon} := \hat{\lambda}$, applied to observations instead of trajectories. In this case, the fact that crucial information is equally obstructed to the human with feature belief function ϵ and to the AI with $\hat{\epsilon}$ should ensure that the ontology translation also correctly translates feature beliefs: $\Psi \circ \hat{\epsilon} = \epsilon$.
- Now consider a setting that may go beyond "pure partial observability". One speculative idea for how to construct $\hat{\epsilon}$ is to prepend a "belief prompt" bp to inputs of $\hat{\lambda}$ that nudges the model to think more in the way a human evaluator would think (Park et al., 2024a):

$$\widehat{\epsilon}(o) \coloneqq \widehat{\lambda}_{\mathrm{bp}}(o) \coloneqq \widehat{\lambda}(\mathrm{bp}, o).$$

An example of such a prompt would be

bp = "Think about the following input like a typical human evaluator:"

Alternatively, one could potentially achieve this by finetuning $\hat{\lambda}$ to obtain $\hat{\epsilon}$. Unfortunately, while foundation models can simulate the behavior of specific people in their *outputs*, it is unclear whether this also reflects in their *internal representations*. For example, prior work shows that the truth-value of statements can sometimes be linearly predicted from internal representations even when the model outputs the non-truth, leading to the potential for AI lie detectors (Azaria & Mitchell, 2023; Burns et al., 2023b). However, other work trains a foundation model to predict human behavior in experiments and finds it to have internal representations that can predict human neural activity when engaging in the same task (Binz et al., 2024).

• Alternatively, one could also consider defining $\widehat{\epsilon} := \widehat{\lambda}_{\text{early}}$ as an earlier training checkpoint of $\widehat{\lambda}$. Being an earlier training checkpoint, $\widehat{\epsilon}$ would then be less capable than $\widehat{\lambda}$, leading to the potential that it has the same blindspots in understanding observations $o \in \mathcal{O}$ as the human evaluator.

Ensuring (ii). To ensure that for all $\xi \in \Xi$ we have $\widehat{\lambda}(\xi) \in \mathbb{R}\langle \widehat{\epsilon}(o) \mid o \in \mathcal{O} \rangle$, it is important to ensure that the image of $\widehat{\epsilon} : \mathcal{O} \to \mathbb{R}^{\widehat{\Omega}}$ is "large", i.e., spans as much as possible of the representation space. To form more intuitions on this, note that by Theorem 3.2, for $\widehat{\mathcal{M}}$ to be complete (which would be implied by property (ii)) requires also the true belief model \mathcal{M} to be complete. Again, by Theorem 2.11, a sufficient condition for this is that for all $\xi \in \Xi$, $\lambda(\xi) \in \mathbb{R}\langle \epsilon(o) \mid o \in \mathcal{O} \rangle$. In particular, any "bad" feature $\omega \in \Omega$ that can ever be present in a trajectory (meaning $[\lambda(\xi)](\omega) \neq 0$) needs to also sometimes be believed to be present by the human (meaning there exists an $o \in \mathcal{O}$ with $[\epsilon(o)](\omega) \neq 0$). This is a relaxation from the requirement that the human understands all observations perfectly, but it does mean that there needs to be a variety of observations $o \in \mathcal{O}$ that is understandable enough for the human to sometimes detect any possible problem. It will depend on the specific MDP and setup of an experiment to reason about how to ensure or purposefully violate this property.

3.4.2 Learning G using $\widehat{\mathcal{M}}$ and $G_{\mathcal{O}}$ for answering question 2

Now that we have discussed preliminary approaches for how to specify a complete model $\widehat{\mathcal{M}} = (\widehat{\Omega}, \widehat{\mathbf{\Lambda}}, \widehat{\mathcal{E}}, \widehat{\mathcal{V}})$ that covers the true belief model \mathcal{M} , we can turn to the question of how to use \mathcal{M} and the human's feedback $G_{\mathcal{O}}: \mathcal{O} \to \mathbb{R}$ to determine the true return function G. By statement 4 in Theorem 3.2, we can determine G as $G = \widehat{\mathbf{\Lambda}}(R_{\widehat{\Omega}})$ for $R_{\widehat{\Omega}} \in \mathbb{R}^{\widehat{\Omega}}$ with $\widehat{\mathcal{E}}(R_{\widehat{\Omega}}) = G_{\mathcal{O}}$.

We now unpack that. Remembering the relation between $\widehat{\mathbf{\Lambda}}$ and the foundation model $\widehat{\lambda}$, and $\widehat{\mathcal{E}}$ and the foundation model $\widehat{\epsilon}$ via Proposition A.1, we want to determine $R_{\widehat{\mathbf{O}}}$ such that for all $o \in \mathcal{O}$, we have

$$G_{\mathcal{O}}(o) = \left[\widehat{\mathcal{E}}(R_{\widehat{\Omega}})\right](o) = \left\langle \widehat{\epsilon}(o), R_{\widehat{\Omega}} \right\rangle. \tag{16}$$

This can be achieved by attaching $R_{\widehat{\Omega}}$ as a linear reward probe to the representation space of $\widehat{\epsilon}$ and learning the function $G_{\mathcal{O}}$ by supervised learning.

In the case that $G_{\mathcal{O}}$ cannot be directly evaluated but that it is indirectly accessible in the form of *choice probabilities* between observations, one can learn $R_{\widehat{\Omega}}$ by logistic regression as in (Christiano et al., 2017). See Appendix B.2 for a possible correspondence between $\widehat{\mathcal{E}}(R_{\widehat{\Omega}})$ and the resulting choice probabilities. Notably, the approach via logistic regression, however, loses a theoretical guarantee: Namely, $G_{\mathcal{O}}$ can at best be learned up to an additive constant. If $\widehat{\epsilon}$ and $\widehat{\lambda}$ were *balanced* (meaning that the total weight of feature strengths is constant over all observations and trajectories, respectively), this would lead to the inferred G also being correct up to an additive constant by Proposition B.6. Since additive constants in return functions are inconsequential for policy optimization, this would be fine. However, typically the representation spaces of foundation models are not normalized, the balancedness property does not hold, and this guarantee breaks. It is then an empirical question to what extent this breakage is an issue or how to resolve it.

After successful learning, we can then compute the true return function G for $\xi \in \Xi$ as:

$$G(\xi) = \left[\widehat{\mathbf{\Lambda}}(R_{\widehat{\Omega}})\right](\xi) = \left\langle \widehat{\lambda}(\xi), R_{\widehat{\Omega}} \right\rangle. \tag{17}$$

In other words, we attach the learned linear reward probe $R_{\widehat{\Omega}}$ to the representation space of $\widehat{\lambda}$ and use it to compute returns. These can then be used to train a policy to maximize the policy evaluation function Equation (4) via standard reinforcement learning techniques.

3.4.3 Further remarks

Looking at the definition of the human belief model $\widehat{\mathcal{M}}$, we see that it includes the ontology $\widehat{\mathbf{\Lambda}}:\mathbb{R}^{\widehat{\Omega}}\to\mathbb{R}^\Xi$ and feature belief function $\widehat{\mathcal{E}}:\mathbb{R}^{\widehat{\Omega}}\to\mathbb{R}^{\mathcal{O}}$. These can be extremely large matrices. However, since in the process of training $R_{\widehat{\Omega}}$ and computing G, we only need to be able to *query* the resulting observation return function and return function on specific observations and trajectories as in Equations (16) and (17), the matrices never need to be stored or used in their entirety. Thus, the size of the matrices is not a concern.

We also remark on a special case we mentioned before: If we are in a case of "pure" partial observability where observations $o \in \mathcal{O}$ contain no information that $\widehat{\lambda}$ understands better than the human, then we proposed to simply set $\widehat{\epsilon} := \widehat{\lambda}$, applied to observations $o \in \mathcal{O}$. In that case, the procedure we describe is essentially classical RLHF, with the only — crucial — difference that during training of $R_{\widehat{\Omega}}$, we only show observations, instead of full trajectories, to the model. This prevents a model misspecification where the data the model reads differs from what the human sees, and could theoretically allow generalization to data $\xi \in \Xi$. Lang et al. (2024) consider naive RLHF, where the initialized return function reads entire trajectories during training time while the human's view is obstructed, leading to issues of deceptive inflation and overjustification after policy optimization.

3.4.4 Conclusions for the proposal

We have thus explained how to specify $\widehat{\mathcal{M}}$ using (possibly adapted) foundation models, and how to learn G by supervised learning of $G_{\mathcal{O}}$ with an attached reward probe. We note again that these ideas are preliminary and leave open many choices to be made, and that empirical research would need to determine whether the

proposal can be instantiated in a compelling way. In Section 4.3.2, we motivate further future work related to this proposal.

4 Discussion

4.1 Summary

In this work, we have introduced the notion of a human belief model, based on modeling a human's ontology and feature belief function. The goal of such a model is to aid the inference of the human's implicit return function from feedback. In our framework, the feedback, in the form of an observation return function, is viewed as carrying information about the reward of features that the human believes to be associated with an observation. Once the feature rewards, in the form of a reward object within a valid set, are inferred, they can be used together with the human's ontology to infer a return function. We defined and characterized the resulting ambiguity in the return function in terms of the belief model and showed that for complete models, the ambiguity disappears. Complete models have an important sufficient condition in terms of a linear coverage of the features of any trajectory by feature beliefs of observations (Theorem 2.11). This shows that for observations that are varied enough and provide ample information in total, a correct return function inference is possible, which we demonstrated in simple conceptual examples in Section 2.8.

Since the human belief model is not known in practice, we then introduced the notion of belief model coverage. Here, model $\widehat{\mathcal{M}}$ covers another model \mathcal{M} if $\widehat{\mathcal{M}}$ can represent all return functions and observation return functions that can be expressed in \mathcal{M} . If a complete model covers the true belief model, then it can be used for the return function inference just as well (Theorem 3.2). We then characterized belief model coverage in terms of belief model morphisms, and found an important sufficient condition given by the existence of a linear translation from the covering model's ontology to the true model's ontology that is also compatible with the feature belief functions (Theorem 3.5). We then studied a conceptual example of a human with symmetry-invariant feature beliefs, which we could cover with a model whose completeness stems from the valid reward functions being symmetry-invariant.

Finally, in Section 3.4 we sketched a proposal for how to find covering human belief models in practice, by using foundation models for both the human's ontology and feature belief function. For the latter, it is important to ensure that the foundation model has a similar understanding of the observtions as the human evaluator, for which we sketched out three proposals. That the resulting belief model might cover the true belief model relies on prior research on the linear representation hypothesis, which provides evidence that a belief-respecting linear ontology translation could indeed exist.⁶ Our hope is that our proposal can help to find covering models that are easier to determine than the return function itself, which might subsequently be trained with modest effort as an approach to scalable oversight (Remark 2.3).

4.2 Related work

Other human modeling approaches. The human belief models we introduce in this work are largely about modeling a human's ontology and feature belief function for learning from the human's feedback. There is also other work that models aspects of humans for better goal inference. For example, reward-rational choice (Jeon et al., 2020) requires modeling human choice probabilities for choices over various options like trajectory-pairs, language-utterances, or initial environment states. This framework can be regarded a special case of assistance problems (Fern et al., 2014; Hadfield-Menell et al., 2016; Shah et al., 2021), which requires a model of the human's action selection in a cooperative two-player game. This framework has recently been generalized to a partially observable setting (Emmons et al., 2024). Much of this work makes specific assumptions about the human's model, like a Boltzmann-rational or optimal selection of choices. (Zhi-Xuan et al., 2024) instantiates a version of assistance games in which a human's utterance is modeled using a language model. Finally, Hatgis-Kessell et al. (2025) researches how to influence human evaluators to conform with a theoretical model of human choices. Compared to all this work, we instead model human

⁶We emphasize again that our theory only requires its existence and no explicit specification of this ontology translation.

beliefs about AI behavior instead of human actions or choices.⁷ In Section 4.3.1 we propose research directions to combine these types of work.

Deception in AI. Our work generalizes the human belief modeling from Lang et al. (2024), which is meant to address deceptive AI behavior that also surfaces in recent empirical work for cases where humans lack evaluative capacity (Cloud et al., 2024; Denison et al., 2024; Wen et al., 2024; Williams et al., 2024). Deception in AI systems can also occur for various other reasons (Park et al., 2024c), e.g., when language agents are put under pressure (Scheurer et al., 2023). An important theoretical concern is deceptive alignment, in which an AI system follows the given goals while actually planning a later takeover (Hubinger et al., 2019). Recent work (Greenblatt et al., 2024) substantiates this concern by showing that Claude plays along with a new harmful goal for the purpose of preventing that the learning process updates its safety behavior in the long-term. Finally, deception has also been formalized for structural causal games (Ward et al., 2023).

knowledge. In Section 3.4, Surfacing latent we already mentioned sparse coders (Cunningham et al., 2023; Bricken et al., 2023), which construct an explicit linear transformation from a foundation model's representation space to human-interpretable features, which is in the spirit of a linear ontology translation (Definition 3.4). Instead of decoding the AI's entire ontology in a human-interpretable way, other paradigms seek to construct a reporter that can be queried for specific information (Christiano et al., 2021). In this direction, recent work builds toward AI lie detectors by finding internal linear representations of truth (Burns et al., 2023b; Azaria & Mitchell, 2023; Marks & Tegmark, 2024), with alternative interpretations of such findings discussed in Liu et al. (2023). Other work linearly predicts concepts like harmfulness (Zou et al., 2024), theft advice (Roger, 2023), the activation of a harmful backdoor (MacDiarmid et al., 2024), and concepts related to honesty and power, among others (Zou et al., 2023).

Amplified oversight. While our work aims to learn from feedback of humans who potentially lack capabilities, work on amplified oversight tries to *increase* the human's evaluation capabilities through AI assistance to achieve scalable oversight. Recursive reward modeling is the general proposal to train AI models by reward modeling and then using their assistance to evaluate the next generation of AIs (Leike et al., 2018). Saunders et al. (2022) shows that model-written critiques of summaries can help humans find flaws that they would have missed on their own. This raises the question why to trust the critiquing AI, which leads to the idea to, in turn, criticize the critic. Recursively, this leads to a debate in which the debaters are trained to produce arguments that are persuasive to a human judge (Irving et al., 2018). This requires for debates to surface true and useful information to the human judges, which has found support for reading comprehension tasks (Michael et al., 2023). Optimizing debaters to be persuasive then increases the human's ability to identify the truth (Khan et al., 2024; Kenton et al., 2024). Finally, some work uses AI to directly give feedback based on a constitution (Bai et al., 2022) or model spec (Guan et al., 2025), thus omitting humans from the evaluation process.

Easy-to-hard and weak-to-strong generalization. Instead of amplifying the evaluator capabilities, other work for scalable oversight relies on easy examples that humans can reliably evaluate, which then requires the reward model to generalize to data that is harder to evaluate (Sun et al., 2024). Language models have also shown to generalize tasks like STEM questions from easy to hard data (Hase et al., 2024). Weak-to-strong generalization differs by trying to generalize from weak evaluation on potentially hard data (Burns et al., 2023a). Our proposal from Section 3.4 can be considered an approach to weak-to-strong generalization since we aim to learn a correct return function G from feedback of an evaluator with potentially faulty beliefs. Since weak supervision is plausibly cheaper to obtain than strong supervision, there is also work that investigates tradeoffs to find the correct allocation of a fixed budget to label data with different data labeling quality (Mallen & Belrose, 2024).

In Appendix E, we briefly interpret amplified oversight, easy-to-hard generalization, and classical RLHF together with weak-to-strong generalization in terms of our theoretical framework by interpreting their

⁷Only in Appendix B do we consider human choices.

underlying evaluator belief modeling assumptions and reasoning about their ambiguities and learned return functions.

4.3 Future work

4.3.1 Theory extensions

Several extensions and generalizations of our theory could be studied in future work. We assumed that we can exactly specify a belief model that covers the true human belief model. Future work could develop an approximate theory, akin to Lang et al. (2024, Theorem 5.4). Furthermore, we assumed that the human's return function is linear in the features of trajectories. One could study non-linear models, which would also theoretically ground the use of non-linear reward probes instead of the linear probes we propose in Section 3.4. We also assumed that learned return functions can read entire trajectories, which might be unrealistic for very complex environments. It would thus be interesting to develop a theory based on a second set of observations \mathcal{O}' for the learned return function and the trained policy. In the practical proposal, this would also mean that we cannot assume to have a foundation model $\hat{\lambda}$ that translates to the human's true ontology—instead, it would represent another feature belief function. Relaxing the capabilities of the foundation model could also help to model a case in which the human evaluator has information that is hidden from the learned return function or resulting policy.

The ambiguity is a measure of the *information* that is available in the feedback, given a belief model, for determining the human's return function. Future work could theoretically study concrete learning procedures, which are about extracting said information. One could, for example, study training distributions over the observations \mathcal{O} to determine sample complexity bounds for the error of the resulting return function, or the regret of the resulting policy. This is theoretically interesting since the usage of the learned return function in policy optimization involves two shifts compared to the return function's learning process: First, it needs to evaluate trajectories instead of observations, making use of an ontology instead of a feature belief function; and second, the policy optimization leads to a further distribution shift over trajectories, which can in turn lead to increased regret even if the return function seemed accurate before (Fluri et al., 2024). Finally, if there is remaining ambiguity, it would also be desirabe to study protocols for deciding for a return function within the ambiguity, possibly using a priori knowledge about "human-like" return functions that is not captured by our notion of valid reward objects.

Going beyond our framework, one could attempt a synthesis with work that models the human's action selection, which we discussed at the start of Section 4.2. For example, one could consider reward-rational choice or assistance games under the assumption that humans form beliefs based on observations. Emmons et al. (2024) does a first step in this direction by assuming humans form rational beliefs based on knowledge of a prior of the agent's policy. In that framework, they then study the notion of information interference, which leads to an increase in the human's uncertainty about the world state. If future work were to study more general, possibly faulty, belief models for humans in partially observable assistance games, this could make it possible to go beyond information interference to study deception.

4.3.2 Empirical work

We would be interested in attempts to instantiate our proposal from Section 3.4. One could test the proposal in settings with synthetic humans with a known ontology and feature belief function, which can for small MDPs allow to compute the ambiguity explicitly. This should make it possible to make concrete predictions about experimental results. It would also be interesting to study settings with partial observability in which capable AI does not have an advantage over humans in understanding the meaning of the observations. As we discussed in our proposal, that should allow to use the same foundation model for the ontology and feature belief function, with the latter only applied to observations. This could then be compared with a baseline of "naive" RLHF in which the initialized return function reads entire trajectories during the learning process, similar to the conceptual examples from Lang et al. (2024). It would be interesting to study different levels of observability to dial the ambiguity up or down. We also encourage to break some of our theoretical assumptions, e.g., by using non-linear reward probes. Finally, it would be desirable to learn how our approach can be combined with approaches in the direction of amplified oversight, easy-to-hard

generalization, or weak-to-strong generalization discussed in Section 4.2. For example, a combination of our work with easy-to-hard generalization could seek to only train the reward pobe on data where the human beliefs are *modeled correctly* by $\hat{\epsilon}$, which should be a superset of easy data.

Instead of studying the whole pipeline, one could also empirically assess the underlying assumptions. The most precarious assumption is that it is possible to construct a feature belief function $\hat{\epsilon}$, for which we discuss proposals in Section 3.4.1. In situations that are not purely based on partial observability, we proposed to prompt the model to step "into the shoes" of a human evaluator, or to use earlier training checkpoints of a model to decrease its capabilities to that of a human evaluator (assuming future models that would otherwise be superhuman). For this to work, there needs to be a linear ontology translation that is compatible with $\hat{\epsilon}$ and the human's true feature belief function ϵ . One could test this empirically by using sparse autoencoders (Cunningham et al., 2023; Bricken et al., 2023) trained on an unobstructed capable foundation model, and evaluating the resulting human-readable features when applied to $\hat{\epsilon}$.

4.4 Conclusion

In conclusion, in this work we theoretically studied models of human beliefs about AI behavior for the goal of scalable oversight, thus contributing to the theory of how to align advanced AI. We hope that future work will build on our theory, empirically study our practical proposal, or engage in other research on how to make AI systems safe and aligned with the goal of decreasing the risks of advanced AI.

Author Contributions

Leon Lang developed the ideas, derived all theoretical results, and wrote the paper. Patrick Forré gave general guidance and feedback.

Acknowledgments and disclosure of funding

We thank Scott Emmons, Davis Foote, Erik Jenner, Micah Carroll, and Alex Cloud for discussions that shaped how we think about human evaluators with limited capabilities. We thank Open Philanthropy for financial support.

References

Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete Problems in AI Safety, 2016. URL https://arxiv.org/abs/1606.06565.

Dario Amodei, Paul Christiano, and Alex Ray. Learning from human preferences. https://openai.com/research/learning-from-human-preferences, 2017. Accessed: 2023-12-13.

Anthropic. Introducing Claude. https://www.anthropic.com/index/introducing-claude, 2023a. Accessed: 2023-09-05.

Anthropic. Claude's Constitution. https://www.anthropic.com/index/claudes-constitution, 2023b. Accessed: 2023-09-05.

Amos Azaria and Tom Mitchell. The Internal State of an LLM Knows When It's Lying, 2023. URL https://arxiv.org/abs/2304.13734.

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback, 2022.

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional AI: Harmlessness from AI Feedback. arXiv e-prints, art. arXiv:2212.08073, December 2022. doi: 10.48550/arXiv.2212.08073.

Marcel Binz, Elif Akata, Matthias Bethge, Franziska Brändle, Fred Callaway, Julian Coda-Forno, Peter Dayan, Can Demircan, Maria K. Eckstein, Noémi Éltető, Thomas L. Griffiths, Susanne Haridi, Akshay K. Jagadish, Li Ji-An, Alexander Kipnis, Sreejan Kumar, Tobias Ludwig, Marvin Mathony, Marcelo Mattar, Alireza Modirshanechi, Surabhi S. Nath, Joshua C. Peterson, Milena Rmus, Evan M. Russek, Tankred Saanum, Natalia Scharfenberg, Johannes A. Schubert, Luca M. Schulze Buschoff, Nishad Singhi, Xin Sui, Mirko Thalmann, Fabian Theis, Vuong Truong, Vishaal Udandarao, Konstantinos Voudouris, Robert Wilson, Kristin Witte, Shuchen Wu, Dirk Wulff, Huadong Xiong, and Eric Schulz. Centaur: a foundation model of human cognition, 2024. URL https://arxiv.org/abs/2410.20268.

Michael Bowling, John D. Martin, David Abel, and Will Dabney. Settling the Reward Hypothesis, 2022.

Samuel R. Bowman, Jeeyoon Hyun, Ethan Perez, Edwin Chen, Craig Pettit, Scott Heiner, Kamilė Lukošiūtė, Amanda Askell, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Christopher Olah, Daniela Amodei, Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-Johnson, Jackson Kernion, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Liane Lovitt, Nelson Elhage, Nicholas Schiefer, Nicholas Joseph, Noemí Mercado, Nova DasSarma, Robin Larson, Sam McCandlish, Sandipan Kundu, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Timothy Telleen-Lawton, Tom Brown, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Ben Mann, and Jared Kaplan. Measuring Progress on Scalable Oversight for Large Language Models, 2022.

Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen, Brayden McLean,

- Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christopher Olah. Towards Monose-manticity: Decomposing Language Models With Dictionary Learning. *Transformer Circuits Thread*, 2023. https://transformer-circuits.pub/2023/monosemantic-features/index.html.
- Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbrenner, Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan Leike, Ilya Sutskever, and Jeff Wu. Weak-to-Strong Generalization: Eliciting Strong Capabilities With Weak Supervision, 2023a. URL https://arxiv.org/abs/2312.09390.
- Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. Discovering Latent Knowledge in Language Models Without Supervision. In *The Eleventh International Conference on Learning Representations*, 2023b. URL https://openreview.net/forum?id=ETKGuby0hcs.
- Paul Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep Reinforcement Learning from Human Preferences. arXiv e-prints, art. arXiv:1706.03741, June 2017. doi: 10.48550/arXiv.1706.03741.
- Paul Christiano, Buck Shlegeris, and Dario Amodei. Supervising strong learners by amplifying weak experts, 2018. URL https://arxiv.org/abs/1810.08575.
- Paul Christiano, Ajeya Cotra, and Mark Xu. Eliciting Latent Knowledge. https://docs.google.com/document/d/1WwsnJQstPq91_Yh-Ch2XRL8H_EpsnjrC1dwZXR37PC8/edit, 2021. Accessed: 2023-04-25.
- Alex Cloud, Jacob Goldman-Wetzler, Evžen Wybitul, Joseph Miller, and Alexander Matt Turner. Gradient Routing: Masking Gradients to Localize Computation in Neural Networks, 2024. URL https://arxiv.org/abs/2410.04332.
- Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse Autoencoders Find Highly Interpretable Features in Language Models, 2023. URL https://arxiv.org/abs/2309.08600.
- Carson Denison, Monte MacDiarmid, Fazl Barez, David Duvenaud, Shauna Kravec, Samuel Marks, Nicholas Schiefer, Ryan Soklaski, Alex Tamkin, Jared Kaplan, Buck Shlegeris, Samuel R. Bowman, Ethan Perez, and Evan Hubinger. Sycophancy to Subterfuge: Investigating Reward-Tampering in Large Language Models, 2024. URL https://arxiv.org/abs/2406.10162.
- Mucong Ding, Chenghao Deng, Jocelyn Choo, Zichu Wu, Aakriti Agrawal, Avi Schwarzschild, Tianyi Zhou, Tom Goldstein, John Langford, Anima Anandkumar, and Furong Huang. Easy2Hard-Bench: Standardized Difficulty Labels for Profiling LLM Performance and Generalization, 2024. URL https://arxiv.org/abs/2409.18433.
- Scott Emmons, Caspar Oesterheld, Vincent Conitzer, and Stuart Russell. Observation Interference in Partially Observable Assistance Games, 2024. URL https://arxiv.org/abs/2412.17797.
- Alan Fern, Sriraam Natarajan, Kshitij Judah, and Prasad Tadepalli. A Decision-Theoretic Model of Assistance. J. Artif. Int. Res., 50(1):71–104, may 2014. ISSN 1076-9757.
- Lukas Fluri, Leon Lang, Alessandro Abate, Patrick Forré, David Krueger, and Joar Skalse. The Perils of Optimizing Learned Reward Functions: Low Training Error Does Not Guarantee Low Regret, 2024. URL https://arxiv.org/abs/2406.15753.
- Google Gemini Team. Gemini: A Family of Highly Capable Multimodal Models. https://storage.googleapis.com/deepmind-media/gemini/gemini_1_report.pdf, 2023. Accessed: 2023-12-11.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere,

Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterii, Ning Zhang, Olivier Duchenne, Onur Celebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujiwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The Llama 3 Herd of Models, 2024. URL https://arxiv.org/abs/2407.21783.

Ryan Greenblatt, Carson Denison, Benjamin Wright, Fabien Roger, Monte MacDiarmid, Sam Marks, Johannes Treutlein, Tim Belonax, Jack Chen, David Duvenaud, Akbir Khan, Julian Michael, Sören Mindermann, Ethan Perez, Linda Petrini, Jonathan Uesato, Jared Kaplan, Buck Shlegeris, Samuel R. Bowman, and Evan Hubinger. Alignment faking in large language models, 2024. URL https://arxiv.org/abs/2412.14093.

Melody Y. Guan, Manas Joglekar, Eric Wallace, Saachi Jain, Boaz Barak, Alec Helyar, Rachel Dias, Andrea Vallone, Hongyu Ren, Jason Wei, Hyung Won Chung, Sam Toyer, Johannes Heidecke, Alex Beutel, and Amelia Glaese. Deliberative Alignment: Reasoning Enables Safer Language Models, 2025. URL https://arxiv.org/abs/2412.16339.

Wes Gurnee and Max Tegmark. Language Models Represent Space and Time. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=jE8xbmvFin.

Dylan Hadfield-Menell, Anca Dragan, Pieter Abbeel, and Stuart Russell. Cooperative Inverse Reinforcement Learning. arXiv e-prints, art. arXiv:1606.03137, June 2016. doi: 10.48550/arXiv.1606.03137.

Peter Hase, Mohit Bansal, Peter Clark, and Sarah Wiegreffe. The Unreasonable Effectiveness of Easy Training Data for Hard Tasks, 2024. URL https://arxiv.org/abs/2401.06751.

Stephane Hatgis-Kessell, W. Bradley Knox, Serena Booth, Scott Niekum, and Peter Stone. Influencing Humans to Conform to Preference Models for RLHF, 2025. URL https://arxiv.org/abs/2501.06416.

Evan Hubinger, Chris van Merwijk, Vladimir Mikulik, Joar Skalse, and Scott Garrabrant. Risks from Learned Optimization in Advanced Machine Learning Systems. arXiv e-prints, art. arXiv:1906.01820, June 2019. doi: 10.48550/arXiv.1906.01820.

Geoffrey Irving, Paul Christiano, and Dario Amodei. AI safety via debate, 2018. URL https://arxiv.org/abs/1805.00899.

- Hong Jun Jeon, Smitha Milli, and Anca Dragan. Reward-rational (implicit) choice: A unifying formalism for reward learning. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), Advances in Neural Information Processing Systems, volume 33, pp. 4415–4426. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/2f10c1578a0706e06b6d7db6f0b4a6af-Paper.pdf
- Zachary Kenton, Noah Y. Siegel, János Kramár, Jonah Brown-Cohen, Samuel Albanie, Jannis Bulian, Rishabh Agarwal, David Lindner, Yunhao Tang, Noah D. Goodman, and Rohin Shah. On scalable oversight with weak LLMs judging strong LLMs, 2024. URL https://arxiv.org/abs/2407.04622.
- Akbir Khan, John Hughes, Dan Valentine, Laura Ruis, Kshitij Sachan, Ansh Radhakrishnan, Edward Grefenstette, Samuel R. Bowman, Tim Rocktäschel, and Ethan Perez. Debating with More Persuasive LLMs Leads to More Truthful Answers, 2024. URL https://arxiv.org/abs/2402.06782.
- S.M. Lane. Categories for the Working Mathematician. Graduate Texts in Mathematics. Springer, 1998. ISBN 9780387984032. URL https://books.google.nl/books?id=MXboNPdTv7QC.
- Leon Lang, Davis Foote, Stuart Russell, Anca Dragan, Erik Jenner, and Scott Emmons. When Your AIs Deceive You: Challenges of Partial Observability in Reinforcement Learning from Human Feedback. In *Advances in Neural Information Processing Systems*, 2024.
- Jan Leike, David Krueger, Tom Everitt, Miljan Martic, Vishal Maini, and Shane Legg. Scalable agent alignment via reward modeling: a research direction, 2018. URL https://arxiv.org/abs/1811.07871.
- Kevin Liu, Stephen Casper, Dylan Hadfield-Menell, and Jacob Andreas. Cognitive Dissonance: Why Do Language Model Outputs Disagree with Internal Representations of Truthfulness?, 2023. URL https://arxiv.org/abs/2312.03729.
- Monte MacDiarmid, Timothy Maxwell, Nicholas Schiefer, Jesse Mu, Jared Kaplan, David Duvenaud, Sam Bowman, Alex Tamkin, Ethan Perez, Mrinank Sharma, Carson Deniagents, andEvan Hubinger. Simple probes can catch sleeper 2024. URL https://www.anthropic.com/news/probes-catch-sleeper-agents.
- Alex Mallen and Nora Belrose. Balancing Label Quantity and Quality for Scalable Elicitation, 2024. URL https://arxiv.org/abs/2410.13215.
- Samuel Marks and Max Tegmark. The Geometry of Truth: Emergent Linear Structure in Large Language Model Representations of True/False Datasets, 2024. URL https://arxiv.org/abs/2310.06824.
- Julian Michael, Salsabila Mahdi, David Rein, Jackson Petty, Julien Dirani, Vishakh Padmakumar, and Samuel R. Bowman. Debate Helps Supervise Unreliable Experts, 2023. URL https://arxiv.org/abs/2311.08702.
- Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. Linguistic Regularities in Continuous Space Word Representations. In Lucy Vanderwende, Hal Daumé III, and Katrin Kirchhoff (eds.), *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 746–751, Atlanta, Georgia, June 2013. Association for Computational Linguistics. URL https://aclanthology.org/N13-1090/.
- Neel Nanda, Andrew Lee, and Martin Wattenberg. Emergent Linear Representations in World Models of Self-Supervised Sequence Models, 2023. URL https://arxiv.org/abs/2309.00941.
- OpenAI. Introducing ChatGPT. https://openai.com/blog/chatgpt, 2022. Accessed: 2024-02-06.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che

Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. GPT-4 Technical Report, 2024. URL https://arxiv.org/abs/2303.08774.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022.

Joon Sung Park, Carolyn Q. Zou, Aaron Shaw, Benjamin Mako Hill, Carrie Cai, Meredith Ringel Morris, Robb Willer, Percy Liang, and Michael S. Bernstein. Generative Agent Simulations of 1,000 People, 2024a. URL https://arxiv.org/abs/2411.10109.

Kiho Park, Yo Joong Choe, and Victor Veitch. The Linear Representation Hypothesis and the Geometry of Large Language Models, 2024b. URL https://arxiv.org/abs/2311.03658.

Peter S Park, Simon Goldstein, Aidan O'Gara, Michael Chen, and Dan Hendrycks. Ai deception: A survey of examples, risks, and potential solutions. *Patterns*, 5(5), 2024c.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct Preference Optimization: Your Language Model is Secretly a Reward Model. arxiv e-prints, 2023.

- Fabien Roger. Coup probes: Catching catastrophes with probes trained off-policy. https://www.alignmentforum.org/posts/WCj7WgFSLmyKaMwPR/coup-probes-catching-catastrophes-with-probes-2023. Accessed: 2023-04-25.
- William Saunders, Catherine Yeh, Jeff Wu, Steven Bills, Long Ouyang, Jonathan Ward, and Jan Leike. Self-critiquing models for assisting human evaluators, 2022. URL https://arxiv.org/abs/2206.05802.
- Jérémy Scheurer, Mikita Balesni, and Marius Hobbhahn. Technical Report: Large Language Models can Strategically Deceive their Users when Put Under Pressure. arxiv e-prints, 2023.
- Rohin Shah, Pedro Freire, Neel Alex, Rachel Freedman, Dmitrii Krasheninnikov, Lawrence Chan, Michael D Dennis, Pieter Abbeel, Anca Dragan, and Stuart Russell. Benefits of Assistance over Reward Learning, 2021. URL https://openreview.net/forum?id=DFIoGDZejIB.
- Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. Learning to summarize from human feedback, 2022. URL https://arxiv.org/abs/2009.01325.
- Zhiqing Sun, Longhui Yu, Yikang Shen, Weiyang Liu, Yiming Yang, Sean Welleck, and Chuang Gan. Easy-to-Hard Generalization: Scalable Alignment Beyond Human Supervision, 2024. URL https://arxiv.org/abs/2403.09472.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open Foundation and Fine-Tuned Chat Models. arxiv e-prints, 2023.
- Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J. Vazquez, Ulisse Mini, and Monte MacDiarmid. Steering Language Models With Activation Engineering, 2024. URL https://arxiv.org/abs/2308.10248.
- Elise van der Pol, Daniel E. Worrall, Herke van Hoof, Frans A. Oliehoek, and Max Welling. MDP Homomorphic Networks: Group Symmetries in Reinforcement Learning, 2021. URL https://arxiv.org/abs/2006.16908.
- Zihao Wang, Lin Gui, Jeffrey Negrea, and Victor Veitch. Concept Algebra for (Score-Text-Controlled Generative Models. In Oh, Τ. Naumann, Globerson, Based) Α. Α. K. Saenko, Hardt, and S. Levine (eds.),Advances in Neural Information Processing Systems, 36,35331–35349. Curran Associates, 2023. URL volumepp. Inc., https://proceedings.neurips.cc/paper_files/paper/2023/file/6f125214c86439d107ccb58e549e828f-Paper-Con
- Francis Rhys Ward, Francesco Belardinelli, Francesca Toni, and Tom Everitt. Honesty Is the Best Policy: Defining and Mitigating AI Deception. arxiv e-prints, 2023.
- Jiaxin Wen, Ruiqi Zhong, Akbir Khan, Ethan Perez, Jacob Steinhardt, Minlie Huang, Samuel R. Bowman, He He, and Shi Feng. Language Models Learn to Mislead Humans via RLHF, 2024. URL https://arxiv.org/abs/2409.12822.
- Marcus Williams, Micah Carroll, Adhyyan Narang, Constantin Weisser, Brendan Murphy, and Anca Dragan. On Targeted Manipulation and Deception when Optimizing LLMs for User Feedback, 2024. URL https://arxiv.org/abs/2411.02306.

- Tan Zhi-Xuan, Lance Ying, Vikash Mansinghka, and Joshua B. Tenenbaum. Pragmatic Instruction Following and Goal Assistance via Cooperative Language-Guided Inverse Planning, 2024. URL https://arxiv.org/abs/2402.17930.
- Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-Tuning Language Models from Human Preferences, 2020.
- Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson, J. Zico Kolter, and Dan Hendrycks. Representation Engineering: A Top-Down Approach to AI Transparency, 2023. URL https://arxiv.org/abs/2310.01405.
- Andy Zou, Long Phan, Justin Wang, Derek Duenas, Maxwell Lin, Maksym Andriushchenko, Rowan Wang, Zico Kolter, Matt Fredrikson, and Dan Hendrycks. Improving Alignment and Robustness with Circuit Breakers, 2024. URL https://arxiv.org/abs/2406.04313.

Appendices

In the appendices, we present mathematical details and content that goes beyond the main text. Appendix A lists preliminary results on linear algebra together with their proofs. In Appendix B, we present a theory of balanced belief models and the ambiguity for feedback that is given by binary choices between observations. This complements the core theory from Sections 2 and 3, where the feedback is given by an observation return function. In Appendix C, we complement Section 3 by showing that human belief models and their morphisms form a category, and we construct a variety of natural models composing a diagram. Appendix D contains further mathematical details for the example in Section 3.3 on symmetry-invariant belief models and reward functions. Finally, Appendix E interprets some of the related work from Section 4.2 in our framework.

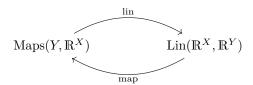
A Preliminary results on linear algebra

For general notation and conventions on linear algebra, see Section 2.1.

A.1 Different representations of linear functions

Let X,Y be two sets. Then by $\operatorname{Lin}(\mathbb{R}^X,\mathbb{R}^Y)$ we denote the set of all linear maps $F:\mathbb{R}^X\to\mathbb{R}^Y$. By $\operatorname{Maps}(Y,\mathbb{R}^X)$, we denote the set of all (simple) maps, or functions, $f:Y\to\mathbb{R}^X$. Intuitively, these encode the same information: A linear function F, when represented as a matrix, is a collection of rows indexed by $y\in Y$, which are "picked out" by a function $f:Y\to\mathbb{R}^X$. This correspondence is made precise in the following proposition, which is a "transposed" version of the classical statement that linear functions on vector spaces correspond to functions on a basis:

Proposition A.1. Define the two functions



as follows: For $f \in \text{Maps}(Y, \mathbb{R}^X)$, $v \in \mathbb{R}^X$ and $y \in Y$, we define

$$\Big[\Big[\ln(f) \Big](v) \Big](y) \coloneqq \langle f(y), v \rangle \,.$$

For $F \in \text{Lin}(\mathbb{R}^X, \mathbb{R}^Y)$, $y \in Y$ and $x \in X$, we define

$$\left[\left[\operatorname{map}(F)\right](y)\right](x) \coloneqq F_{yx}.$$

Then lin(f) has matrix elements

$$lin(f)_{yx} = [f(y)](x)$$

for $x \in X$ and $y \in Y$. Furthermore, \lim and \lim are mutually inverse bijections.

Proof. First, note that for each $f \in \text{Maps}(Y, \mathbb{R}^X)$, $\text{lin}(f) : \mathbb{R}^X \to \mathbb{R}^Y$ is indeed a linear function since the scalar product is linear in the second component. Its matrix elements are given by

$$lin(f)_{yx} = \left[\left[lin(f) \right] (e_x) \right] (y) = \langle f(y), e_x \rangle = \left[f(y) \right] (x).$$

Now we show that lin and map are mutually inverse bijections, i.e., $\lim \circ \text{map} = \text{id}_{\text{Lin}(\mathbb{R}^X,\mathbb{R}^Y)}$ and $\text{map} \circ \text{lin} = \text{id}_{\text{Maps}(Y,\mathbb{R}^X)}$. Indeed, for $F \in \text{Lin}(\mathbb{R}^X,\mathbb{R}^Y)$, $x \in X$, and $y \in Y$, we have

$$\left[(\ln \circ \operatorname{map})(F) \right]_{yx} = \left[\ln(\operatorname{map}(F)) \right]_{yx}$$

$$= \left[\left[\ln(\text{map}(F)) \right] (e_x) \right] (y)$$

$$= \left\langle \left[\text{map}(F) \right] (y), e_x \right\rangle$$

$$= \left[\left[\text{map}(F) \right] (y) \right] (x)$$

$$= F_{yx}$$

Since linear functions are fully characterized by their matrix elements, this shows $(\ln \circ \text{map})(F) = F$, and thus $\ln \circ \text{map}$ is the identity.

For the other direction, for $f \in \text{Maps}(Y, \mathbb{R}^X)$, $y \in Y$, and $x \in X$, we have

$$\left[\left[\left[(\operatorname{map} \circ \operatorname{lin})(f) \right](y) \right](x) = \left[\left[\operatorname{map}(\operatorname{lin}(f)) \right](y) \right](x) \\
= \operatorname{lin}(f)_{yx} \\
= \left[f(y) \right](x).$$

This shows that $(\text{map} \circ \text{lin})(f) = f$, and so $\text{map} \circ \text{lin}$ is also the identity.

Proposition A.2. Let X, \widehat{X}, Y be sets and $F : \mathbb{R}^X \to \mathbb{R}^Y$, $\widehat{F} : \mathbb{R}^{\widehat{X}} \to \mathbb{R}^Y$, and $\Phi : \mathbb{R}^X \to \mathbb{R}^{\widehat{X}}$ be linear functions. Let $f = \text{map}(F) : Y \to \mathbb{R}^X$ and $\widehat{f} = \text{map}(\widehat{F}) : Y \to \mathbb{R}^{\widehat{X}}$ be the functions corresponding to F and \widehat{F} by Proposition A.1. Let $\Phi^T : \mathbb{R}^{\widehat{X}} \to \mathbb{R}^X$ be the transpose of Φ , with matrix elements $\Phi^T_{\widehat{xx}} = \Phi_{\widehat{xx}}$. Then $F = \widehat{F} \circ \Phi$ if and only if $f = \Phi^T \circ \widehat{f}$:



Proof. We have

$$F = \widehat{F} \circ \Phi \iff \forall y \in Y, x \in X \colon F_{yx} = (\widehat{F} \circ \Phi)_{yx} = \sum_{\widehat{x} \in \widehat{X}} \widehat{F}_{y\widehat{x}} \Phi_{\widehat{x}x}$$

$$\iff \forall y \in Y, x \in X \colon [f(y)](x) = \sum_{\widehat{x} \in \widehat{X}} \Phi_{x\widehat{x}}^T [\widehat{f}(y)](\widehat{x}) = \langle \Phi_x^T, \widehat{f}(y) \rangle$$

$$= \left[\Phi^T (\widehat{f}(y)) \right](x) = \left[(\Phi^T \circ \widehat{f})(y) \right](x)$$

$$\iff f = \Phi^T \circ \widehat{f}.$$

That was to show.

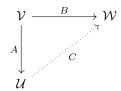
A.2 Properties of kernels and images of linear functions

The following two propositions list basic and well-known properties of kernels and images that we use in the paper:

Proposition A.3. Let $A: \mathcal{V} \to \mathcal{U}$ and $B: \mathcal{V} \to \mathcal{W}$ be linear functions. Then the following statements are equivalent:

1.
$$\ker(A) \subseteq \ker(B)$$
;

2. There exists a linear function $C: \mathcal{U} \to \mathcal{W}$ with $C \circ A = B$:



Proof. Clearly, the second claim implies the first. So now assume 1. Let $\{u^1, \ldots, u^m\}$ be a basis for $\operatorname{im}(A)$ and complement it to a basis $\{u^1, \ldots, u^n\}$ for all of \mathcal{U} , where $n \geq m$. For each $i \in \{1, \ldots, m\}$, let $v^i \in \mathcal{V}$ be any element with $A(v^i) = u^i$. Define $C(u^i) := B(v^i)$ for $i \in \{1, \ldots, m\}$ and $C(u^i) = 0$ if i > m. Linearly extend C to a linear function $C: \mathcal{U} \to \mathcal{W}$.

To show that $C \circ A = B$, let $v \in \mathcal{V}$ be arbitrary. Then $A(v) \in \text{im}(A)$, and thus there exist coefficients $\lambda_i \in \mathbb{R}$ for $i \in \{1, ..., m\}$ with

$$A(v) = \sum_{i=1}^{m} \lambda_i u^i.$$

Note that

$$A\left(v - \sum_{i=1}^{m} \lambda_i v^i\right) = A(v) - \sum_{i=1}^{m} \lambda_i A(v^i) = A(v) - \sum_{i=1}^{m} \lambda_i u^i = 0$$

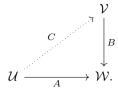
Thus, $v - \sum_{i=1}^{m} \lambda_i v^i \in \ker(A) \subseteq \ker(B)$. Consequently, we obtain

$$B(v) = B\left(\sum_{i=1}^{m} \lambda_i v^i\right) = \sum_{i=1}^{m} \lambda_i B(v^i) = \sum_{i=1}^{m} \lambda_i C(u^i)$$
$$= C\left(\sum_{i=1}^{m} \lambda_i u^i\right) = C(A(v)) = (C \circ A)(v).$$

This shows $C \circ A = B$, and thus the claim.

Proposition A.4. Let $A: \mathcal{U} \to \mathcal{W}$ and $B: \mathcal{V} \to \mathcal{W}$ be linear functions. Then the following statements are equivalent:

- 1. $\operatorname{im}(A) \subseteq \operatorname{im}(B)$.
- 2. There exists a "lift", i.e., a linear map $C: \mathcal{U} \to \mathcal{V}$ with $B \circ C = A$:



Proof. It can easily be checked that the second claim implies the first. For the other direction, let $\{u^1, \ldots, u^n\}$ be a basis for \mathcal{U} . For $i \in \{1, \ldots, n\}$, choose $v^i \in \mathcal{V}$ with $B(v^i) = A(u^i)$, which exists since $\operatorname{im}(A) \subseteq \operatorname{im}(B)$. Define $C: \mathcal{U} \to \mathcal{V}$ as the unique linear function with $C(u^i) = v^i$ for $i \in \{1, \ldots, n\}$. We obtain

$$A(u^{i}) = B(v^{i}) = B(C(u^{i})) = (B \circ C)(u^{i}).$$

Since linear functions are determined on a basis, it follows $A = B \circ C$, proving the claim.

B Balanced human belief models and choices

In this appendix, we present the core theory from Section 2 and Section 3 for the case that the feedback is in the form of *choices* instead of the observation return function $G_{\mathcal{O}}$. To still get a useful theory, we will then need to assume our human belief models to be *balanced* to ensure that constants are "propagated" appropriately through the model. In this whole appendix, we fix an MDP with a set of trajectories Ξ and observations \mathcal{O} .

B.1 Balanced belief models

Definition B.1 (Row-constant). Let X, and Y be sets. We call a linear function $A: \mathbb{R}^X \to \mathbb{R}^Y$ row-constant if for all $y, y' \in Y$ we have $0 \neq \sum_{x \in X} X_{yx} = \sum_{x \in X} X_{y'x}$.

Definition B.2 (Balanced). Let $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$ be a human belief model. We call \mathcal{M} balanced if Λ and \mathcal{E} are row-constant, and if $\mathcal{V} \subseteq \mathbb{R}^{\Omega}$ contains all constant functions.

 Λ and \mathcal{E} being row-constant means that the corresponding functions $\lambda:\Xi\to\mathbb{R}^\Omega$ and $\epsilon:\mathcal{O}\to\mathbb{R}^\Omega$ (cf. Proposition A.1) map to vectors of feature strengths with a constant *total* weighting, as is for example the case for probability distributions. That \mathcal{V} contains all constant functions is often naturally the case. To demonstrate this in a simple example we first prove the following lemma:

Lemma B.3. Let X, Y, Z be sets and $A : \mathbb{R}^X \to \mathbb{R}^Y$, $B : \mathbb{R}^Y \to \mathbb{R}^Z$ be row-constant linear functions. Then the composition $B \circ A : \mathbb{R}^X \to \mathbb{R}^Z$ is also row-constant.

Proof. Let a, b be the row-sums of A and B, respectively. Then for all $z \in Z$, we obtain

$$\sum_{x \in X} (B \circ A)_{zx} = \sum_{x \in X} \sum_{y \in Y} B_{zy} A_{yx} = \sum_{y \in Y} B_{zy} \sum_{x \in X} A_{yx} = b \cdot a \neq 0,$$

which shows the claim.

Example B.4. We continue Example 2.4 and show that the model is balanced. Note that for all $\xi \in \Xi$, we have

$$\sum_{(s,a,s')\in\mathcal{S}\times\mathcal{A}\times\mathcal{S}} \mathbf{\Gamma}_{\xi,(s,a,s')} = \sum_{(s,a,s')\in\mathcal{S}\times\mathcal{A}\times\mathcal{S}} \sum_{t=0}^{T-1} \gamma^t \delta_{(s,a,s')}(s_t,a_t,s_{t+1})$$

$$= \sum_{t=0}^{T-1} \gamma^t \sum_{(s,a,s')\in\mathcal{S}\times\mathcal{A}\times\mathcal{S}} \delta_{(s,a,s')}(s_t,a_t,s_{t+1})$$

$$= \sum_{t=0}^{T-1} \gamma^t$$

$$\neq 0.8$$

Thus, Γ is row-constant. \mathbf{B} is also row-constant since all rows are probability distributions, and so Lemma B.3 implies that also $\mathcal{E} = \mathbf{B} \circ \Gamma$ is row-constant. $\mathcal{V} = \mathbb{R}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}}$ also clearly contains all constant functions. Overall, this means the model $(\mathcal{S} \times \mathcal{A} \times \mathcal{S}, \Gamma, \mathbf{B} \circ \Gamma, \mathbb{R}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}})$ is balanced.

For more examples of balanced human belief models, see Appendix C.2.

B.2 The ambiguity for balanced belief models and choices

Let $\sigma: \mathbb{R} \to (0,1)$ be any known bijective function with $\sigma(r) + \sigma(-r) = 1$ for all $r \in \mathbb{R}$ (e.g., the sigmoid function). For $o, o' \in \mathcal{O}$, we define the probability that a human with feature belief function $\mathcal{E}: \mathbb{R}^{\Omega} \to \mathbb{R}^{\mathcal{O}}$ and reward object $\tilde{R}_{\Omega} \in \mathbb{R}^{\Omega}$ prefers o over o' by

$$P_{\mathcal{E}}^{\tilde{R}_{\Omega}}(o \succ o') := \sigma\Big(\big[\mathcal{E}(\tilde{R}_{\Omega})\big](o) - \big[\mathcal{E}(\tilde{R}_{\Omega})\big](o')\Big). \tag{18}$$

For the rest of the section, we fix a human belief model $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$. Furthermore, we fix the *implicit* true reward object $R_{\Omega} \in \mathcal{V}$ together with the return function $G = \Lambda(R_{\Omega})$, the observation return function $G_{\mathcal{O}} = \mathcal{E}(R_{\Omega})$ and the choice probability function $P_{\mathcal{O}} := P_{\mathcal{E}}^{R_{\Omega}}$ that serves as our operationalization of "feedback".

We use the following adaptation of Definition 2.6:

Definition B.5. We define the set of return functions that are **feedback-compatible** with $P_{\mathcal{O}}$ as

$$\mathrm{FC}^{\mathcal{M}}(P_{\mathcal{O}}) \coloneqq \Big\{ \tilde{G} \in \mathbb{R}^{\Xi} \; \big| \; \exists \tilde{R}_{\Omega} \in \mathcal{V} \colon P_{\mathcal{E}}^{\tilde{R}_{\Omega}} = P_{\mathcal{O}} \; and \; \mathbf{\Lambda}(\tilde{R}_{\Omega}) = \tilde{G} \Big\}.$$

We define the ambiguity left in the return function G after the choice probability function $P_{\mathcal{O}}$ is known by

$$\mathrm{Amb}^{\mathcal{M}}(G, P_{\mathcal{O}}) := \Big\{ G' \in \mathbb{R}^{\Xi} \ \big| \ G' = \tilde{G} - G \ \text{for} \ \tilde{G} \in \mathrm{FC}^{\mathcal{M}}(P_{\mathcal{O}}) \Big\}.$$

Clearly, we have

$$FC^{\mathcal{M}}(G_{\mathcal{O}}) = G + Amb^{\mathcal{M}}(G, P_{\mathcal{O}}).$$

Recall the ambiguity $\mathrm{Amb}^{\mathcal{M}}(G,G_{\mathcal{O}})$ defined in Definition 2.6. We obtain:

Proposition B.6. Assume $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$ is balanced. Let $\mathbf{1} \in \mathbb{R}^{\Xi}$ denote the function that is constant 1. Then:

$$\mathrm{Amb}^{\mathcal{M}}(G, P_{\mathcal{O}}) = \mathrm{Amb}^{\mathcal{M}}(G, G_{\mathcal{O}}) + \mathbb{R}\langle \mathbf{1} \rangle = \mathbf{\Lambda} \big(\ker(\mathcal{E}) \cap \mathcal{V} \big) + \mathbb{R}\langle \mathbf{1} \rangle.$$

Proof. The second equality follows from Proposition 2.7, so we are left with proving the first. By abuse of notation, we will write 1 for the three functions that are constant 1 on Ξ , \mathcal{O} , or Ω .

Let $G' \in \mathrm{Amb}^{\mathcal{M}}(G, P_{\mathcal{O}})$. Then $G' = \mathbf{\Lambda}(\tilde{R}_{\Omega}) - G$ for $\tilde{R}_{\Omega} \in \mathcal{V}$ with $P_{\mathcal{E}}^{\tilde{R}_{\Omega}} = P_{\mathcal{O}}$. The latter means the following for all $o, o' \in \mathcal{O}$:

$$\sigma\Big(\big[\mathcal{E}(\tilde{R}_{\Omega})\big](o) - \big[\mathcal{E}(\tilde{R}_{\Omega})\big](o')\Big) = \sigma\Big(\big[\mathcal{E}(R_{\Omega})\big](o) - \big[\mathcal{E}(R_{\Omega})\big](o')\Big).$$

Since σ is invertible, we obtain

$$\left[\mathcal{E}(\tilde{R}_{\Omega})\right](o) - \left[\mathcal{E}(\tilde{R}_{\Omega})\right](o') = \left[\mathcal{E}(R_{\Omega})\right](o) - \left[\mathcal{E}(R_{\Omega})\right](o')$$

Fix any $o' \in \mathcal{O}$ and set $c_{\mathcal{O}} := \left[\mathcal{E}(\tilde{R}_{\Omega})\right](o') - \left[\mathcal{E}(R_{\Omega})\right](o')$. Then for all $o \in \mathcal{O}$, we have

$$[\mathcal{E}(\tilde{R}_{\Omega})](o) = [\mathcal{E}(R_{\Omega})](o) + c_{\mathcal{O}}.$$

Or, equivalently:

$$\mathcal{E}(\tilde{R}_{\Omega}) = \mathcal{E}(R_{\Omega}) + c_{\mathcal{O}} \cdot \mathbf{1} = G_{\mathcal{O}} + c_{\mathcal{O}} \cdot \mathbf{1}.$$

Since \mathcal{E} is row-constant, there exists a $c_{\Omega} \in \mathbb{R}$ with $\mathcal{E}(c_{\Omega} \cdot \mathbf{1}) = c_{\mathcal{O}} \cdot \mathbf{1}$. This implies

$$\mathcal{E}(\tilde{R}_{\Omega} - c_{\Omega} \cdot \mathbf{1}) = G_{\mathcal{O}}.$$

We also have $\tilde{R}_{\Omega} - c_{\Omega} \cdot \mathbf{1} \in \mathcal{V}$ since \mathcal{V} contains all constant functions. Furthermore, $\mathbf{\Lambda}(c_{\Omega} \cdot \mathbf{1}) = c_{\Xi} \cdot \mathbf{1}$ for another constant $c_{\Xi} \in \mathbb{R}$ since $\mathbf{\Lambda}$ is row-constant. Overall, we thus obtain

$$G' = \mathbf{\Lambda}(\tilde{R}_{\Omega}) - G = \mathbf{\Lambda}(\tilde{R}_{\Omega} - c_{\Omega} \cdot \mathbf{1}) - G + c_{\Xi} \cdot \mathbf{1} \in Amb^{\mathcal{M}}(G, G_{\mathcal{O}}) + \mathbb{R}\langle \mathbf{1} \rangle.$$

For the other direction, let $G' \in \text{Amb}^{\mathcal{M}}(G, G_{\mathcal{O}}) + \mathbb{R}\langle \mathbf{1} \rangle$. Then $G' = \tilde{G} - G + c_{\Xi} \cdot \mathbf{1}$ for $\tilde{G} = \mathbf{\Lambda}(\tilde{R}_{\Omega})$ with $\tilde{R}_{\Omega} \in \mathcal{V}$ with $\mathcal{E}(\tilde{R}_{\Omega}) = G_{\mathcal{O}}$. We have

$$\tilde{G} + c_{\Xi} \cdot \mathbf{1} = \mathbf{\Lambda} (\tilde{R}_{\Omega} + c_{\Omega} \cdot \mathbf{1})$$

for a constant $c_{\Omega} \in \mathbb{R}$ with $\Lambda(c_{\Omega} \cdot \mathbf{1}) = c_{\Xi} \cdot \mathbf{1}$. Since \mathcal{V} contains all constant functions, we have $\tilde{R}_{\Omega} + c_{\Omega} \cdot \mathbf{1} \in \mathcal{V}$. We also have

$$P_{\mathcal{E}}^{\tilde{R}_{\Omega}+c_{\Omega}\cdot\mathbf{1}}=P_{\mathcal{E}}^{\tilde{R}_{\Omega}}=P_{\mathcal{O}}$$

since the constant gets cancelled out in the definition of the choice probabilities, and since $\mathcal{E}(\tilde{R}_{\Omega}) = G_{\mathcal{O}}$. All of this implies

$$G' = (\tilde{G} + c_{\Xi} \cdot \mathbf{1}) - G \in Amb^{\mathcal{M}}(G, P_{\mathcal{O}}).$$

That proves the claim.

Note that for the purpose of policy optimization it is not an issue that the ambiguity has an "irreducible" constant term since this does not change the ordering of policies under the policy evaluation function $J(\pi) = \mathbf{E}_{\xi \in P^{\pi}(\cdot)}[G(\xi)]$.

Remark B.7. In light of the previous proposition, it turns out that the ambiguity does not depend on the true return function and choice probabilities, and we can thus write it as $Amb_P^{\mathcal{M}} = Amb^{\mathcal{M}}(G, P_{\mathcal{O}})$. The "P" is added to distinguish from the ambiguity $Amb^{\mathcal{M}} = Amb^{\mathcal{M}}(G, G_{\mathcal{O}})$ that we study in the main paper.

Using this result, we also obtain a version of Theorem 3.2, with the ambiguity replaced by the one we use in this appendix:

Theorem B.8. Let $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$ and $\widehat{\mathcal{M}} = (\widehat{\Omega}, \widehat{\Lambda}, \widehat{\mathcal{E}}, \widehat{\mathcal{V}})$ be two balanced human belief models and assume that $\widehat{\mathcal{M}}$ covers \mathcal{M} . We think of \mathcal{M} as the "true" human belief model with reward object $R_{\Omega} \in \mathcal{V}$ and corresponding return function $G = \Lambda(R_{\Omega})$ and choice probability function $P_{\mathcal{O}} = P_{\mathcal{E}}^{R_{\Omega}}$. Then we have:

- 1. $\operatorname{Amb}_{P}^{\mathcal{M}} \subseteq \operatorname{Amb}_{P}^{\widehat{\mathcal{M}}}$.
- 2. If \mathcal{M} also covers $\widehat{\mathcal{M}}$, then $\operatorname{Amb}_{P}^{\mathcal{M}} = \operatorname{Amb}_{P}^{\widehat{\mathcal{M}}}$.
- 3. There is an $R_{\widehat{\Omega}} \in \widehat{\mathcal{V}}$ with $P_{\widehat{\mathcal{E}}}^{R_{\widehat{\Omega}}} = P_{\mathcal{O}}$ and $\widehat{\Lambda}(R_{\widehat{\Omega}}) = G$.
- 4. Assume $\widehat{\mathcal{M}}$ is complete. Then every reward object $\widetilde{R}_{\widehat{\Omega}} \in \widehat{\mathcal{V}}$ with $P_{\widehat{\mathcal{E}}}^{\widetilde{R}_{\widehat{\Omega}}} = P_{\mathcal{O}}$ also satisfies $\widehat{\Lambda}(\widetilde{R}_{\widehat{\Omega}}) = G + c_{\Xi} \cdot \mathbf{1}$ for a constant $c_{\Xi} \in \mathbb{R}$.

Proof. Statements 1 and 2 follow from the ambiguity characterization in Proposition B.6 and the analogous statements in Theorem 3.2. Statements 3 and 4 can be proved with similar arguments as the corresponding statements in Theorem 3.2.

C A diagram in the category of human belief models

Let us consider an MDP together with a fixed set of trajectories Ξ and observations \mathcal{O} . Then in Definition 2.2, we defined the notion of a human belief model $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$. In Definition 3.3, we then introduced the notion of a morphism $\Phi: \mathcal{M} \to \widehat{\mathcal{M}}$ between human belief models, which is defined as a linear function $\Phi: \mathbb{R}^{\Omega} \to \mathbb{R}^{\widehat{\Omega}}$ such that $\Phi(\mathcal{V}) \subseteq \mathcal{V}'$, $\widehat{\Lambda} \circ \Phi|_{\mathcal{V}} = \Lambda$ and $\widehat{\mathcal{E}} \circ \Phi|_{\mathcal{V}} = \mathcal{E}$. This notion turned out important since it is equivalent to model covering (Definition 3.1), which implies that the covering model can be used for the return function inference from human feedback (Theorem 3.2), especially if its ambiguity disappears.

Belief models for fixed sets of trajectories Ξ and observations \mathcal{O} , together with their morphisms, form a category (Lane, 1998), meaning that they satisfy the following simple properties:

- Composition: Assume $\mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3$ are three human belief models and $\Phi : \mathcal{M}_1 \to \mathcal{M}_2, \Phi' : \mathcal{M}_2 \to \mathcal{M}_3$ morphisms between them. Then also the composition $\Phi' \circ \Phi : \mathcal{M}_1 \to \mathcal{M}_3$ is a morphism.
- Identities: For any human belief model $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$, the identity $id_{\mathbb{R}^{\Omega}} : \mathcal{M} \to \mathcal{M}$ is a morphism.
- Associativity: $(\Phi'' \circ \Phi') \circ \Phi = \Phi'' \circ (\Phi' \circ \Phi)$ for any three morphisms that can be composed in the specified order.

All of these properties can be trivially checked, and so human belief models and their morphisms indeed form a category.

In this appendix, we want to write down a simple commutative diagram of morphisms in this category. Here, a diagram means a graph of human belief models and morphisms between them. For this to be *commutative* means that any pathway from one human belief model to another is the same morphism. We prepare

this in Appendix C.1 by writing down all linear functions from which the functions Λ , \mathcal{E} , and Φ will be constructed. In Appendix C.2 we then specify the resulting human belief models and briefly consider their properties. In Appendix C.3 we interpret the matrix elements that appear in the feature belief functions of all models. Finally, in Appendix C.4, we write down the resulting commutative diagram and the resulting relations for the ambiguities.

C.1 Preparing the models

We build on Example 2.4. The idea is that we consider reward objects at four different levels: Return functions, classical reward functions, and return- and reward functions of abstractions of trajectories and transitions that the human might care about. By modeling the human as having features at all four of these different levels, we can create a multitude of human belief models.

Let $\mathbf{B}: \mathbb{R}^{\Xi} \to \mathbb{R}^{\mathcal{O}}$ be the matrix corresponding to a trajectory-belief function $b: \mathcal{O} \to \Delta(\Xi) \subseteq \mathbb{R}^{\Xi}$ via Proposition A.1. Let $\Gamma: \mathbb{R}^{S \times \mathcal{A} \times S} \to \mathbb{R}^{\Xi}$ be the linear function mapping reward functions to their corresponding return functions.

Let \mathcal{F} be a set of "abstractions of transitions" and $h: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathcal{F}$ a function mapping each transition to its abstraction. Write reward objects over abstractions as $R_{\mathcal{F}} \in \mathbb{R}^{\mathcal{F}}$. Then we obtain the induced map

$$h^*: \mathbb{R}^{\mathcal{F}} \to \mathbb{R}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}}, \quad R_{\mathcal{F}} \mapsto R_{\mathcal{F}} \circ h.$$

 $h^*(R_{\mathcal{F}})$ measures a transition (s, a, s') by evaluating $R_{\mathcal{F}}$ at the transition's abstraction: $R_{\mathcal{F}}(h(s, a, s'))$. Thus, $h^*(R_{\mathcal{F}})$ is guaranteed to give the same reward to transitions with the same abstraction.

We can then also consider the space of abstraction sequences \mathcal{F}^T together with the function $h^T:\Xi\to\mathcal{F}^T$ given by

$$h^{T}(s_{0}, a_{0}, \dots, s_{T-1}, a_{T-1}, s_{T}) := (h(s_{0}, a_{0}, s_{1}), \dots, h(s_{T-1}, a_{T-1}, s_{T})).$$

Write return functions over abstraction sequences as $G_{\mathcal{F}^T} \in \mathbb{R}^{\mathcal{F}^T}$. h^T then gives rise to the dual function

$$h^{T^*}: \mathbb{R}^{\mathcal{F}^T} \to \mathbb{R}^{\Xi}, \quad G_{\mathcal{F}^T} \mapsto G_{\mathcal{F}^T} \circ h^T.$$

Thus, $h^{T^*}(G_{\mathcal{F}^T})$ evaluates a trajectory by evaluating the sequence of abstractions using $G_{\mathcal{F}^T}$. As before, two trajectories with the same sequences of abstractions then obtain the same return.

Recall the function $\Gamma: \mathbb{R}^{S \times A \times S} \to \mathbb{R}^{\Xi}$ mapping a reward function to the corresponding return function. Then we obtain an analogous function for reward objects on abstractions:

$$\mathbf{\Gamma}_{\mathcal{F}}: \mathbb{R}^{\mathcal{F}} \to \mathbb{R}^{\mathcal{F}^T}, \quad \left[\mathbf{\Gamma}_{\mathcal{F}}(R_{\mathcal{F}})\right](f_1, \dots, f_T) \coloneqq \sum_{t=0}^{T-1} \gamma^t R_{\mathcal{F}}(f_t).$$

Proposition C.1. The diagram

of linear functions commutes, meaning that all pathways with the same start and end are the same function.

Proof. We have

$$[(\Gamma \circ h^*)(R_{\mathcal{F}})](s_0, a_0, \dots, a_{T-1}, s_T) = [\Gamma (h^*(R_{\mathcal{F}}))](s_0, a_0, \dots, a_{T-1}, s_T)$$

$$= \sum_{t=0}^{T-1} \gamma^{t} [h^{*}(R_{\mathcal{F}})] (s_{t}, a_{t}, s_{t+1})$$

$$= \sum_{t=0}^{T-1} \gamma^{t} R_{\mathcal{F}} (h(s_{t}, a_{t}, s_{t+1}))$$

$$= [\Gamma_{\mathcal{F}}(R_{\mathcal{F}})] (h(s_{0}, a_{0}, s_{1}), \dots, h(s_{T-1}, a_{T-1}, s_{T}))$$

$$= [\Gamma_{\mathcal{F}}(R_{\mathcal{F}})] (h^{T}(s_{0}, a_{0}, \dots, s_{T-1}, a_{T-1}, s_{T}))$$

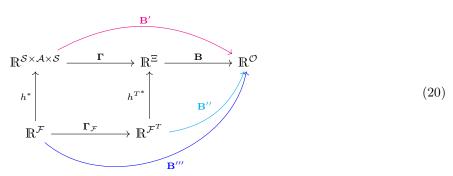
$$= [h^{T^{*}} (\Gamma_{\mathcal{F}}(R_{\mathcal{F}}))] (s_{0}, a_{0}, \dots, a_{T-1}, s_{T})$$

$$= [(h^{T^{*}} \circ \Gamma_{\mathcal{F}})(R_{\mathcal{F}})] (s_{0}, a_{0}, \dots, a_{T-1}, s_{T}).$$

This shows $\Gamma \circ h^* = h^{T^*} \circ \Gamma_{\mathcal{F}}$. Consequently, the diagram commutes.

The idea will be that the rows of \mathbf{B} , $\mathbf{B} \circ \mathbf{\Gamma}$, $\mathbf{B} \circ h^{T^*}$ and $\mathbf{B} \circ \mathbf{\Gamma} \circ h^* = \mathbf{B} \circ h^{T^*} \circ \mathbf{\Gamma}_{\mathcal{F}}$ all correspond (via Proposition A.1) to feature beliefs over trajectories, transitions, trajectory abstractions, and transition abstractions, respectively. We explain this interpretation in detail in Appendix C.3. All of these functions "factorize" over trajectories, but of course this need not be the case in reality: A realistic human could have an intrinsic belief over state transitions, sequences of abstractions, or single abstractions, without this belief "factorizing" in a rational way over state sequences.

Thus, let the following be an extended version of the diagram from Proposition C.1, with new linear functions $\mathbf{B}', \mathbf{B}'', \mathbf{B}'''$. This extension is now *not* necessarily commutative anymore:



To interpret $\mathbf{B}', \mathbf{B}'', \mathbf{B}'''$ on similar grounds as \mathbf{B} , it makes sense to assume that they are row-constant (Definition $\mathbf{B}.1$), but otherwise they can be arbitrary.

C.2 Various human belief models

We take the previous diagrams as the starting point to construct human belief models Remember that a belief model is of the form $\mathcal{M} = (\Omega, \Lambda, \mathcal{E}, \mathcal{V})$. The sets $\Xi, \mathcal{S} \times \mathcal{A} \times \mathcal{S}$, \mathcal{F} , and \mathcal{F}^T are four different possible feature sets Ω . Λ is given by a composition of linear functions that maps to \mathbb{R}^{Ξ} . \mathcal{E} is given by a composition mapping to $\mathbb{R}^{\mathcal{O}}$. The space \mathcal{V} is either given by the full vector space \mathbb{R}^{Ω} , or by images of functions mapping to \mathbb{R}^{Ω} . Overall, using the diagram from Proposition C.1, this leads to the following 9 models, with the superscript denoting the feature space, and the subscript indicating where the valid reward objects "originate from":

$$\mathcal{M}_{\mathcal{F}}^{\mathcal{F}} \coloneqq \left(\mathcal{F}, \ \mathbf{\Gamma} \circ h^*, \ \mathbf{B} \circ \mathbf{\Gamma} \circ h^*, \ \mathbb{R}^{\mathcal{F}} \right)$$

$$\mathcal{M}_{\mathcal{S} \times \mathcal{A} \times \mathcal{S}}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}} \coloneqq \left(\mathcal{S} \times \mathcal{A} \times \mathcal{S}, \ \mathbf{\Gamma}, \ \mathbf{B} \circ \mathbf{\Gamma}, \ \mathbb{R}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}} \right)$$

$$\mathcal{M}_{\mathcal{F}}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}} \coloneqq \left(\mathcal{S} \times \mathcal{A} \times \mathcal{S}, \ \mathbf{\Gamma}, \ \mathbf{B} \circ \mathbf{\Gamma}, \ \operatorname{im}(h^*) \right)$$

$$\mathcal{M}_{\mathcal{F}}^{\mathcal{F}^T} \coloneqq \left(\mathcal{F}^T, \ h^{T^*}, \ \mathbf{B} \circ h^{T^*}, \ \mathbb{R}^{\mathcal{F}^T} \right)$$

$$\mathcal{M}_{\mathcal{F}}^{\mathcal{F}^T} \coloneqq \left(\mathcal{F}^T, \ h^{T^*}, \ \mathbf{B} \circ h^{T^*}, \ \operatorname{im}(\mathbf{\Gamma}_{\mathcal{F}}) \right)$$

$$\mathcal{M}_{\Xi}^{\Xi} \coloneqq \left(\Xi, \ \mathrm{id}_{\mathbb{R}^{\Xi}}, \ \mathbf{B}, \ \mathbb{R}^{\Xi}\right)$$

$$\mathcal{M}_{\mathcal{S} \times \mathcal{A} \times \mathcal{S}}^{\Xi} \coloneqq \left(\Xi, \ \mathrm{id}_{\mathbb{R}^{\Xi}}, \ \mathbf{B}, \ \mathrm{im}(\Gamma)\right)$$

$$\mathcal{M}_{\mathcal{F}^{T}}^{\Xi} \coloneqq \left(\Xi, \ \mathrm{id}_{\mathbb{R}^{\Xi}}, \ \mathbf{B}, \ \mathrm{im}(h^{T^{*}})\right)$$

$$\mathcal{M}_{\mathcal{F}}^{\Xi} \coloneqq \left(\Xi, \ \mathrm{id}_{\mathbb{R}^{\Xi}}, \ \mathbf{B}, \ \mathrm{im}(\Gamma \circ h^{*})\right)$$

For example, $\mathcal{M}_{S \times \mathcal{A} \times S}^{S \times \mathcal{A} \times S}$ is the model from Example 2.4; $\mathcal{M}_{\mathcal{F}}^{S \times \mathcal{A} \times S}$ is the same model, but with valid reward functions restricted to those that only "care about" abstractions; \mathcal{M}_{Ξ}^{Ξ} is a model in which the features are given by full trajectories, and there are no restrictions on the valid return functions; etc.

Now, \mathbf{B}' naturally gives rise to the following three models, which we color differently to distinguish them more easily:

$$\mathcal{M}_{\mathcal{F}}^{\prime\mathcal{F}} = \left(\mathcal{F}, \ \mathbf{\Gamma} \circ h^*, \ \mathbf{B}^\prime \circ h^*, \ \mathbb{R}^{\mathcal{F}}\right)$$
$$\mathcal{M}_{\mathcal{F}}^{\prime\mathcal{S} \times \mathcal{A} \times \mathcal{S}} = \left(\mathcal{S} \times \mathcal{A} \times \mathcal{S}, \ \mathbf{\Gamma}, \ \mathbf{B}^\prime, \ \operatorname{im}(h^*)\right)$$
$$\mathcal{M}_{\mathcal{S} \times \mathcal{A} \times \mathcal{S}}^{\prime\mathcal{S} \times \mathcal{A} \times \mathcal{S}} = \left(\mathcal{S} \times \mathcal{A} \times \mathcal{S}, \ \mathbf{\Gamma}, \ \mathbf{B}^\prime, \ \mathbb{R}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}}\right).$$

Similarly, \mathbf{B}'' gives rise to the following three models:

$$\mathcal{M}_{\mathcal{F}}^{"\mathcal{F}} = \left(\mathcal{F}, \ h^{T^*} \circ \Gamma_{\mathcal{F}}, \ \mathbf{B}'' \circ \Gamma_{\mathcal{F}}, \ \mathbb{R}^{\mathcal{F}}\right)$$
$$\mathcal{M}_{\mathcal{F}}^{"\mathcal{F}^T} = \left(\mathcal{F}^T, \ h^{T^*}, \ \mathbf{B}'', \ \operatorname{im}(\Gamma_{\mathcal{F}})\right)$$
$$\mathcal{M}_{\mathcal{F}^T}^{"\mathcal{F}^T} = \left(\mathcal{F}^T, \ h^{T^*}, \ \mathbf{B}'', \ \mathbb{R}^{\mathcal{F}^T}\right)$$

Finally, \mathbf{B}''' gives rise to a single model:

$$\mathcal{M}_{\mathcal{F}}^{\prime\prime\prime\mathcal{F}} = \left(\mathcal{F}, \ \mathbf{\Gamma} \circ h^*, \ \mathbf{B}^{\prime\prime\prime}, \ \mathbb{R}^{\mathcal{F}}\right)$$

Note that all component linear functions appearing in any of these models (identities, $h^*, \Gamma, \Gamma_{\mathcal{F}}, h^{T^*}, \mathbf{B}, \dots, \mathbf{B}'''$) are row-constant. By Lemma B.3 then, also all compositions are row-constant, which then implies that all 16 models are balanced, as defined in Definition B.2. The first 9 models are also faithful (Definition 2.13) since all feature belief functions factorize as in Proposition 2.14, with Y given by **B** in all cases. The other 7 models will typically not be faithful.

C.3 Interpreting the matrix elements

We now interpret the different feature belief functions that appeared in the nine first models of the previous subsection. Recall that the linear function $\mathbf{B}: \mathbb{R}^\Xi \to \mathbb{R}^\mathcal{O}$ "comes from" a function $b: \mathcal{O} \to \Delta(\Xi) \subseteq \mathbb{R}^\Xi$. Thus, all matrix elements $\mathbf{B}_{o\xi}$ can be interpreted as a probability $[b(o)](\xi)$ for the trajectory ξ when viewing observation o. We now explain similar interpretations for the matrix elements of all the other feature belief functions:

 $\mathbf{B} \circ \mathbf{\Gamma}$: It contains matrix elements

$$(\mathbf{B} \circ \mathbf{\Gamma})_{o,(s,a,s')} = \sum_{\xi} [b(o)](\xi) \sum_{t=0}^{T-1} \gamma^t \delta_{(s,a,s')}(s_t, a_t, s_{t+1}),$$

the expected discounted number of times the transition (s, a, s') is present in the trajectory.

 $\mathbf{B} \circ h^{T^*}$: Write **f** for (f_1, \dots, f_T) . Then this contains matrix elements

$$\begin{split} (\mathbf{B} \circ \boldsymbol{h}^{T^*})_{o\mathbf{f}} &= \sum_{\xi \in \Xi} \left[b(o) \right] (\xi) \boldsymbol{h}_{\xi\mathbf{f}}^{T^*} \\ &= \sum_{\xi \in \Xi} \left[b(o) \right] (\xi) \delta_{\mathbf{f}} (\boldsymbol{h}^T(\xi)) \end{split}$$

$$\begin{split} &= \sum_{\xi \colon h^T(\xi) = \mathbf{f}} \left[b(o) \right] (\xi) \\ &= \left[b(o) \right] \left((h^T)^{-1} (\mathbf{f}) \right) \\ &= \left[b(o)_{h^T} \right] (\mathbf{f}). \end{split}$$

In the second to last step, we view b(o) as a probability distribution that, when evaluated on a set, evaluates to the sum of the probabilities of the set's elements. In the last step, we use the definition of the distributional law of a random variable X with respect to a probability distribution P on the sample space: $P_X(x) = P(X^{-1}(x))$. The result is the believed probability, after observing o, of a trajectory with sequence of abstractions \mathbf{f} .

Finally, we look at the matrix $\mathbf{B} \circ \mathbf{\Gamma} \circ h^*$ (for which we give two slightly different formulas): The matrix elements are given as

$$\begin{split} (\mathbf{B} \circ \mathbf{\Gamma} \circ h^*)_{of} &= \sum_{(s,a,s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}} (\mathbf{B} \circ \mathbf{\Gamma})_{o,(s,a,s')} \cdot h^*_{(s,a,s'),f} \\ &= \sum_{(s,a,s') \colon h(s,a,s') = f} \sum_{\xi \in \Xi} \left[b(o) \right] (\xi) \sum_{t=0}^{T-1} \gamma^t \delta_{(s,a,s')} (s_t, a_t, s_{t+1}) \\ &= \sum_{\xi \in \Xi} \left[b(o) \right] (\xi) \sum_{t=0}^{T-1} \gamma^t \sum_{(s,a,s') \colon h(s,a,s') = f} \delta_{(s,a,s')} (s_t, a_t, s_{t+1}) \\ &= \sum_{\xi \in \Xi} \left[b(o) \right] (\xi) \sum_{t=0}^{T-1} \gamma^t \delta_f (h(s_t, a_t, s_{t+1})). \end{split}$$

This is the expected discounted number of times that one encounters the abstraction f. Using that $\Gamma \circ h^* = h^{T^*} \circ \Gamma_{\mathcal{F}}$ by Proposition C.1, we can also write this as

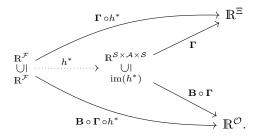
$$\begin{split} \left(\mathbf{B} \circ \boldsymbol{h}^{T^*} \circ \boldsymbol{\Gamma}_{\mathcal{F}} \right)_{of} &= \sum_{\mathbf{f} \in \mathcal{F}^T} (\mathbf{B} \circ \boldsymbol{h}^{T^*})_{o\mathbf{f}} \cdot (\boldsymbol{\Gamma}_{\mathcal{F}})_{\mathbf{f}f} \\ &= \sum_{\mathbf{f} \in \mathcal{F}^T} \left[b(o)_{\boldsymbol{h}^T} \right] (\mathbf{f}) \sum_{t=0}^T \gamma^t \delta_f(f_t). \end{split}$$

This can also be described as the expected discounted number of times that one encounters the abstraction f.

C.4 The resulting commutative diagram

Building on Appendix C.2, the following is a commutative diagram of belief models and model morphisms, with four differently colored "connected components"; one can sometimes use Proposition C.1 in the process of showing that every linear function in the diagram is a morphism of belief models, and that the final diagram commutes:

For example, the following diagram visualizes the fact that $h^*: \mathcal{M}_{\mathcal{F}}^{\mathcal{F}} \to \mathcal{M}_{\mathcal{F}}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}}$ is a morphism:



This gives rise to the following diagram of ambiguities:

The ambiguities are computed using Proposition 2.7, and the inclusions and equalities of ambiguities follow from Theorem 3.2 and Theorem 3.5. Here, the ambiguity $\Gamma\left(\ker(\mathbf{B}\circ\Gamma)\right) = \ker(\mathbf{B}) \cap \operatorname{im}(\Gamma)$ is the special case discussed in depth in Lang et al. (2024). Note that the models $\mathcal{M}_{\mathcal{F}}^{\mathcal{F}}, \mathcal{M}_{\mathcal{F}}^{\mathcal{S}\times\mathcal{A}\times\mathcal{S}}$ and $\mathcal{M}_{\mathcal{S}\times\mathcal{A}\times\mathcal{S}}^{\mathcal{S}\times\mathcal{A}\times\mathcal{S}}$ are closely related to the models $\mathcal{M}_1, \mathcal{M}_2$ and \mathcal{M}_3 from Section 3.3.

Assume we would use one of these models in practice. The further right or up it is in the diagram, the more ambiguity there is, but it is then also more likely that the model covers the true belief model (should it appear in the diagram in the first place). Thus, there is a trade-off between covering the true belief model, and keeping the ambiguity small.

D Details on the example with invariant features

Here, we present more mathematical details for Section 3.3. This appendix is not self-contained and we recommend reading it alongside the section in the main paper.

D.1 Details on the MDP and observations

Formally, the states are given by $S = (\{L, R\} \times \{U, D\})^2$, with the first component being the hand-position, and the second component being the button position. For example, the state in Equation (10) is given by ((L, U), (R, D)).

Furthermore, we define functions $\operatorname{Pos}_H : \mathcal{S} \to \{L, R\} \times \{U, D\}$ and $\operatorname{Pos}_B : \mathcal{S} \to \{L, R\} \times \{U, D\}$ as the first and second projection. These are the position of a "hand" H and a "button" B, in a 2x2 gridworld. Then the state s from Equation (10) satisfies $\operatorname{Pos}_H(s) = (L, U)$ and $\operatorname{Pos}_B(s) = (R, D)$.

The set of trajectories is formally given by $\Xi = (S \times A)^3 \times S$. The set of observations is formally given by

$$\mathcal{O} = \left[\{L, R\}^2 \times \{P, \overline{P}\} \right]^3 \times \{L, R\}^2.$$

Here, \overline{P} means that it was *not* observed that a button was pressed.

D.2 Details on the belief models and symmetries

Let $G = D_4$ be the dihedral group of order 8, i.e., the symmetry group of the square. It is given by

$$G = D_4 = \{e, r, r^2, r^3, f, rf, r^2f, r^3f\},$$

where r is a clockwise rotations by 90° and f is a flip over the horizontal axis. In compositions, we apply f first. G acts on $S \times A$ by individually acting on states and actions:

$$g.(s,a) \coloneqq (g.s,g.a),$$

where $g.s = g.(\operatorname{Pos}_H(s), \operatorname{Pos}_B(s)) := (g.\operatorname{Pos}_H(s), g.\operatorname{Pos}_B(s))$, where on the generators g = r and g = f we have

$$r.(L,U) = (R,U), \quad r.(R,U) = (R,D), \quad r.(R,D) = (L,D), \quad r.(L,D) = (L,U)$$

 $f.(L,U) = (L,D), \quad f.(R,U) = (R,D), \quad f.(R,D) = (R,U), \quad f.(L,D) = (L,U).$

This specifies the action on states. On actions, we specify

$$r.L=U, \quad r.U=R, \quad r.R=D, \quad r.D=L, \quad r.P=P.$$

$$f.L=L, \quad f.U=D, \quad f.R=R, \quad f.D=U, \quad f.P=P.$$

Thus, the "pressing" action remains invariant.

With this group action, we obtain a set of equivalence classes of state-action pairs, given by $\overline{S \times A}$. A set of representatives for the equivalence classes is given by

$$(\lbrace s_0 \rbrace \times \mathcal{A}^{s_0}) \cup (\lbrace s_1 \rbrace \times \mathcal{A}^{s_1}) \cup (\lbrace s_2 \rbrace \times \mathcal{A}^{s_2}), \tag{21}$$

where

$$s_0 = ((R, D), (R, D)), \quad s_1 = ((L, D), (R, D)), \quad s_2 = ((L, U), (R, D)),$$

and where the (state-dependent) set of actions are given by

$$\mathcal{A}^{s_0} = \mathcal{A}^{s_2} = \{L, D, P\}, \quad \mathcal{A}^{s_1} = \{L, R, U, D, P\}.$$

We then have a function

$$h: \mathcal{S} \times \mathcal{A} \to \bigcup_{i \in \{0,1,2\}} \{s_i\} \times \mathcal{A}^{s_i}$$

that maps each state-action pair to a representative, given by h(s,a) = g.(s,a) for the unique $g \in D_4$ for which g.(s,a) is in the set of representatives from Equation (21). Via h, we now identify $\overline{S} \times \overline{A}$ with $\bigcup_{i \in \{0.1.2\}} \{s_i\} \times A^{s_i}$.

Equation (13) can be showed by

$$\Lambda_{\xi,(s,a)} = [\lambda(\xi)](s,a)$$

$$= \sum_{t=0}^{2} \delta_{(s,a)}(h(s_{t},a_{t}))$$

$$= \sum_{(s',a')\in\mathcal{S}\times\mathcal{A}} \delta_{(s,a)}(h(s',a')) \sum_{t=0}^{2} \delta_{(s',a')}(s_{t},a_{t})$$

$$= \sum_{(s',a')\in\mathcal{S}\times\mathcal{A}} h^{*}_{(s',a'),(s,a)} \Gamma_{\xi,(s',a')}$$

$$= \sum_{(s',a')\in\mathcal{S}\times\mathcal{A}} \Gamma_{\xi,(s',a')} h^{*}_{(s',a'),(s,a)}$$

$$= (\Gamma \circ h^{*})_{\xi,(s,a)}.$$
(22)

For the matrix elements of h^* , we used Equation (1).

The human's belief $\mathcal{E}(s, a \mid o)$ for $(s, a) \in \overline{\mathcal{S} \times \mathcal{A}}$ and $o \in \mathcal{O}$ is then given as follows: The human has a uniform prior $B(s) = P_0$ over possible start-states sampled from P_0 , and a uniform prior over possibly next actions given the current state, leading to a prior distribution over $B(\xi) \in \Delta(\Xi)$. Then, upon seeing o, the human implicitly computes a posterior belief over trajectories compatible with the observation, simply given by

$$B(\xi \mid o) \propto \delta_o(O(\xi)) \cdot B(\xi). \tag{23}$$

Equation (14) can be showed by

$$\mathcal{E}_{o,(s,a)} = [\epsilon(o)](s,a)$$

$$= \sum_{\xi \in \Xi} [b(o)](\xi) \cdot [\lambda(\xi)](s,a)$$

$$= \sum_{\xi \in \Xi} \mathbf{B}_{o\xi} \cdot \mathbf{\Lambda}_{\xi,(s,a)}$$

$$= (\mathbf{B} \circ \mathbf{\Lambda})_{o,(s,a)}$$

$$= (\mathbf{B} \circ \mathbf{\Gamma} \circ h^*)_{o,(s,a)}$$
(24)

D.3 Details on the ambiguity analysis for \mathcal{M}_2

Recall the observation o_2 :

Since we assumed that the starting state is one of the states in Equation (11), the human has the belief $[\epsilon(o_2)](s_2, P) = 3$, i.e., the human is certain that, up to symmetry, the hand performed a pressing action three times in s_2 . Thus,

$$0 = \left[\mathcal{E}(\overline{R})\right](o_2) = \left[\epsilon(o_2)\right](s_2, P) \cdot \overline{R}(s_2, P) = 3 \cdot \overline{R}(s_2, P)$$

This implies $\overline{R}(s_2, P) = 0$.

Recall observation o_1 :

Now, the first action could either not change anything, or horizontally align H and B. An action that does not change anything is more likely (chance 2/3 since there are two actions, in the direction of two different adjacent walls, that achieve this, which both correspond to action L up to symmetry), and so we obtain

$$[\epsilon(o_1)](s_2, L) = 2/3, \quad [\epsilon(o_1)](s_2, P) = 4/3,$$

 $[\epsilon(o_1)](s_2, D) = 1/3, \quad [\epsilon(o_1)](s_1, P) = 2/3.$

Compare also with (24). Thus, we obtain

$$0 = \left[\mathcal{E}(\overline{R}) \right] (o_1)$$

= $2/3 \cdot \overline{R}(s_2, L) + 4/3 \cdot \overline{R}(s_2, P) + 1/3 \cdot \overline{R}(s_2, D) + 2/3 \cdot \overline{R}(s_1, P)$
= $2/3 \cdot \overline{R}(s_1, P)$.

Here, we used that $\overline{R}(s_2, P) = 0$ by what we showed before, and $\overline{R}(s_2, L) = \overline{R}(s_2, D) = 0$ since $\overline{R}(s, a) = 0$ whenever $a \neq P$ (i.e., since $\overline{R} \in \mathcal{V}$). Thus, we have $\overline{R}(s_1, P) = 0$ as well.

Finally, we look at the observation sequence o_0 given as follows:

Again, there is a chance of 2/3 that the first action does not change anything. Given the first step, everything which follows is deterministic, leading to these feature beliefs:

$$[\epsilon(o_0)](s_2, L) = 2/3, \quad [\epsilon(o_0)](s_2, R) = 2/3, \quad [\epsilon(o_0)](s_1, P) = 2/3$$
$$[\epsilon(o_0)](s_2, D) = 1/3, \quad [\epsilon(o_0)](s_1, R) = 1/3, \quad [\epsilon(o_0)](s_0, P) = 1/3.$$

Compare again with (24). This means that

$$0 = \left[\mathcal{E}(\overline{R}) \right] (o_0)$$

$$= 2/3 \cdot \overline{R}(s_2, L) + 2/3 \cdot \overline{R}(s_2, R) + 2/3 \cdot \overline{R}(s_1, P) + 1/3 \cdot \overline{R}(s_2, D) + 1/3 \cdot \overline{R}(s_1, R) + 1/3 \cdot \overline{R}(s_0, P)$$

$$= 1/3 \cdot \overline{R}(s_0, P).$$

Here, we used that $\overline{R}(s_1, P)$ by what we showed before, together, again, with the fact that $\overline{R}(s, a) = 0$ for all $a \neq P$. That shows $\overline{R}(s_0, P) = 0$.

D.4 Details on the ambiguity analysis for \mathcal{M}_3

We have

$$[(\mathbf{B} \circ \mathbf{\Gamma})(R')](o) = (\mathbf{B} \circ \mathbf{\Gamma})_{o,(s_1',P)} \cdot R'(s_1',P) + (\mathbf{B} \circ \mathbf{\Gamma})_{o,(s_1'',P)} \cdot R'(s_1'',P)$$

$$= (\mathbf{B} \circ \mathbf{\Gamma})_{o,(s_1',P)} - (\mathbf{B} \circ \mathbf{\Gamma})_{o,(s_1'',P)}$$

$$= 0$$
(25)

In the computation, the second step follows from the definition of R'. The last step follows from the symmetry remarked on before.

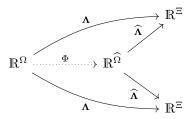
E Mathematical interpretations of related work in our framework

In this appendix, we briefly interpret some of the related work from Section 4.2 in our framework for the special case that they learn linear reward probes. Note that these interpretations are not meant to capture everything there is to say about that work — the summaries we provide are quite coarse. In all examples below, we assume access to a very capable foundation model $\hat{\lambda}:\Xi\to\mathbb{R}^{\widehat{\Omega}}$ that allows for a linear ontology translation $\Psi:\mathbb{R}^{\widehat{\Omega}}\to\mathbb{R}^{\widehat{\Omega}}$ to the human's ontology λ , as in Section 3.4.1: $\Psi\circ\hat{\lambda}=\lambda$. Define $\Phi:=\Psi^T:\mathbb{R}^\Omega\to\mathbb{R}^{\widehat{\Omega}}$, which then satisfies $\hat{\Lambda}\circ\Phi=\Lambda$ by Proposition A.2. In all approaches below, we define $\hat{\epsilon}$ and assume that the return function is learned with the same method as in Section 3.4.2. Notably, in all of the approaches one essentially just defines $\hat{\epsilon}:=\hat{\lambda}$, i.e., no explicit modeling of humans is performed.

E.1 Amplified oversight and eliciting latent knowledge

In amplified oversight, one *amplifies* the human to give accurate feedback, which means we can assume $\mathcal{O} = \Xi$ and $\epsilon = \lambda$. One approach to achieve this would be to essentially *define* $\epsilon := \Psi \circ \widehat{\lambda}$ by giving the human access to the linear ontology translation Ψ for understanding the foundation model's thoughts. This would roughly be in the spirit of eliciting latent knowledge (Christiano et al., 2021), where the human can query a reporter to give information about arbitrary latent knowledge of an AI.

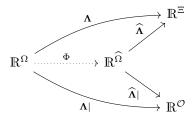
Accordingly, one can also choose $\hat{\epsilon} = \hat{\lambda}$, leading to the following coverage diagram:



The ontologies and feature belief functions are then the same, which automatically means that the ambiguity disappears: $\Lambda(\ker(\Lambda)) = 0$.

E.2 Easy-to-hard generalization

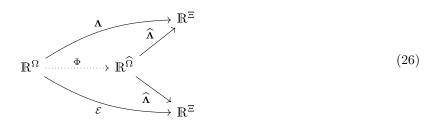
In this setting, $\mathcal{O} \subseteq \Xi$ is a *subset* of trajectories that the human correctly understands. Thus, for $\xi \in \mathcal{O}$, one has $\epsilon(\xi) = \lambda(\xi)$, and so $\epsilon = \lambda|_{\mathcal{O}}$ is simply a restriction. In this setting, one can also set $\hat{\epsilon} = \widehat{\lambda}|_{\mathcal{O}}$. Now, let Λ and $\widehat{\Lambda}$ be the linear functions corresponding to $\lambda|_{\mathcal{O}}$ and $\widehat{\lambda}|_{\mathcal{O}}$, respectively, via Proposition A.1. One obtains the following diagram:



For Φ in this diagram to be a morphism, we need that the lower diagram commutes. With Proposition A.2, this follows from the assumption that $\Phi^T = \Psi$ is an ontology translation: $\Psi \circ \widehat{\lambda} = \lambda$ implies $\Psi \circ \widehat{\lambda}|_{\mathcal{O}} = \lambda|_{\mathcal{O}}$. The ambiguity is now given by $\widehat{\mathbf{\Lambda}}\big(\ker(\widehat{\mathbf{\Lambda}}|)\big)$. By using Theorem 2.11, this ambiguity disappears if and only if for all $\xi \in \Xi$, we have $\widehat{\lambda}(\xi) \in \mathbb{R}\langle \widehat{\lambda}(\xi) \mid \xi \in \mathcal{O} \rangle$. Thus, the ambiguity vanishes if the trajectories that the human understands have enough variety in the vector space of feature strengths.

E.3 Classical RLHF and weak-to-strong generalization

In classical RLHF, without any safeguards, one just uses the model $\hat{\epsilon} = \hat{\lambda}$ as the feature belief function even though $\epsilon \neq \lambda$ and hopes for the best:



In this case, the lower triangle does not commute since $\widehat{\Lambda} \circ \Phi \neq \mathcal{E}$, which, using Proposition A.2, is due to $\Psi \circ \widehat{\lambda} = \lambda \neq \epsilon$. This means that the second model does not cover the true belief model, and so the guarantee from Theorem 3.2 breaks. In fact, the return function that would be inferred using this model is $G_{\mathcal{O}}$, the observation return function itself (cf. the definition of feedback-compatible return functions, Definition 2.6, applied to this faulty model). Lang et al. (2024) extensively discuss failure modes in this case, called deceptive inflation and overjustification. We note that weak-to-strong generalization (Burns et al., 2023a),

when used without additional techniques, also considers this setting, but tries to ensure that the learning process or $\hat{\lambda}$ contains inductive biases that steer the learning process to learn G anyway.