




A survey on interpretable reinforcement learning

Claire Glanois¹ · Paul Weng²  · Matthieu Zimmer³ · Dong Li⁴ · Tianpei Yang⁵ · Jianye Hao^{4,5} · Wulong Liu⁴

Received: 15 August 2022 / Revised: 4 March 2024 / Accepted: 24 March 2024 /

Published online: 19 April 2024

© The Author(s), under exclusive licence to Springer Science+Business Media LLC, part of Springer Nature 2024

Abstract

Although deep reinforcement learning has become a promising machine learning approach for sequential decision-making problems, it is still not mature enough for high-stake domains such as autonomous driving or medical applications. In such contexts, a learned policy needs for instance to be interpretable, so that it can be inspected before any deployment (e.g., for safety and verifiability reasons). This survey provides an overview of various approaches to achieve higher interpretability in reinforcement learning (RL). To that aim, we distinguish interpretability (as an intrinsic property of a model) and explainability (as a post-hoc operation) and discuss them in the context of RL with an emphasis on the former notion. In particular, we argue that interpretable RL may embrace different facets: interpretable inputs, interpretable (transition/reward) models, and interpretable decision-making. Based on this scheme, we summarize and analyze recent work related to interpretable RL with an emphasis on papers published in the past 10 years. We also discuss briefly some related research areas and point to some potential promising research directions, notably related to the recent development of foundation models (e.g., large language models, RL from human feedback).

Keywords Reinforcement learning · Deep reinforcement learning · Interpretability · Explainability

Editor: Bo Liu.

Claire Glanois, Paul Weng and Matthieu Zimmer have contributed equally to this work.

✉ Paul Weng
paul.weng@dukekunshan.edu.cn

¹ Real Lab, ITU, Copenhagen, Denmark

² Duke Kunshan University, Kunshan, China

³ Huawei, London, UK

⁴ Huawei, Beijing, China

⁵ Tianjin University, Tianjin, China

1 Introduction

Reinforcement learning (RL) (Sutton & Barto, 2018) is a general machine learning framework for designing systems with automatic decision-making capabilities. Research in RL has soared since its combination with deep learning, called deep RL (DRL), achieving several recent impressive successes (e.g., AlphaGo (Silver et al., 2017), video game (Vinyals et al., 2019), or robotics (OpenAI et al., 2019)). These attainments were made possible notably thanks to the introduction of the powerful approximation capability of deep learning and its adoption for sequential decision-making and adaptive control.

However, this combination has simultaneously brought all the drawbacks of deep learning to RL. Indeed, as noticed by abundant recent work in DRL, policies learned via a DRL algorithm may suffer from various weaknesses, e.g.:

- They are generally hard to understand because of the blackbox nature of deep neural network architectures (Zahavy et al., 2016).
- They are difficult to train, require a large amount of data, and DRL experiments are often difficult to replicate (Henderson et al., 2018).
- They may overfit the training environment and may not generalize well to new situations (Zhang et al., 2018b).
- Consequently, they may be unsafe and vulnerable to adversarial attacks (Huang et al., 2017).

These observations reveal why DRL is currently not ready for real-world high-stake applications such as autonomous driving or healthcare, and explain why interpretable and explainable RL has recently become a very active research direction. In this survey, we view interpretability as an intrinsic property of a model and explainability as a post-hoc operation (see Sect. 3)

Most real-world deployments of RL algorithms require that learned policies are intelligible as they provide an answer (or a basis for an answer) to various concerns encompassing ethical, legal, operational, or usability viewpoints:

Ethical concerns When designing an autonomous system, it is essential to ensure that its behavior follows some ethical and fairness principles discussed and agreed upon beforehand by the stakeholders according to the context (Crawford et al., 2016; Dwork et al., 2012; Friedler et al., 2021; Leslie, 2020; Lo Piano, 2020; Morley et al., 2020; Yu et al., 2018). The growing discussion about bias and fairness in machine learning (Mehrabi et al., 2019) suggests that mitigating measures must be taken in every aspect of an RL methodology as well. In this regard, intelligibility is essential to help assess the embedding of moral values into autonomous systems, and contextually evaluate and debate their equity and social impact.

Legal concerns As autonomous systems start to be deployed, legal issues arise regarding notably safety (Amodei et al., 2016), accountability (Commission, 2019; Doshi-Velez et al., 2019), or privacy (Horvitz & Mulligan, 2015). For instance, fully-autonomous driving cars should be permitted in the streets only once proven safe with high confidence. The question of risk management (Bonnefon et al., 2019) but also responsibility, in the case of an accident involving such systems, has become a more pressing and complex problem. Verification, accountability, but also privacy can only be ensured with more transparent systems.

Operational concerns Since transparent systems are inspectable and verifiable, they can be examined before deployment to ensure that their decision-making is based on meaningful (ideally causal) relations and not on spurious features, ensuring higher reliability and increased robustness. From the vantage point of researchers or engineers, such systems have the advantage of being more easily debugged and corrected. Moreover, one may expect that such systems are easier to train, more data efficient, and more generalizable and transferable to new domains thanks to interpretability inductive biases.

Usability concerns Interpretable and explainable models can form an essential component for building more interactive systems, where an end-user can request more information about the outcome or decision-making process. In particular, explainable systems would arguably be more trustworthy, which is a key requirement for their integration and acceptance (Mohseni et al., 2020), although the question of trust touches on many other contextual and non-epistemic factors (e.g., risk aversion or goal) beyond intelligibility.

In addition to this high-level list of concerns, we refer the interested reader to Whittlestone et al. (2021) for a more thorough discussion about the potential societal impact of the deployment of DRL-based systems. Although interpretability is a pertinent instrument to achieve more accountable AI-systems, the debate around their real-life implementation should stay active, and include diverse expertise from legal, ethical, and socio-political fields, whose coverage goes beyond the scope of this survey.

Motivated by the importance of these concerns, the number of publications in DRL specifically tackling interpretability issues has increased significantly in recent years. The surging popularity of this topic also explains the recent publication of three survey papers (Alharin et al., 2020; Heuillet et al., 2021; Puiutta & Veith, 2020) on interpretable and explainable RL. In Puiutta and Veith (2020) and Heuillet et al. (2021), a short overview is provided with a limited scope, notably in terms of surveyed papers, while Alharin et al. (2020) cover more studies, organized and categorized into explanation types. The presentation of those surveys generally leans towards explainability as opposed to interpretability (see Sect. 3 for the definitions adopted in this survey) and focuses on understanding the decision-making part of RL.

In contrast, this survey aims at providing a more comprehensive view of what may constitute interpretable RL, which we here specifically distinguish from explainable RL (see Sect. 3). In particular, while decision-making is indeed an important aspect of RL, we believe that achieving interpretability in RL should involve a more encompassing discussion of every component involved in these algorithms, and should stand on three pillars: interpretable inputs (e.g., percepts or other structural information provided to the agent), interpretable transition/reward models, and interpretable decision-making.

Based on this observation, we organize previous work that proposes methods for achieving greater interpretability in RL, along those three components, with an emphasis on DRL papers published in the last 10 years. Thus, in contrast to the previous three surveys, we cover additional work that belongs to interpretable RL such as relational RL or neuro-symbolic RL and also draw connections to other work that naturally falls in this designation, such as object-based RL, physics-based models, or logic-based task descriptions. We also briefly examine the potential impact of large language models on interpretable and explainable RL. One goal of this proposal is to discuss the work in (deep) RL that is specifically identified as belonging to interpretable RL and to draw connections to previous work in RL that is related to interpretability. Since such latter work covers a very broad research space, we can only provide a succinct account for it.

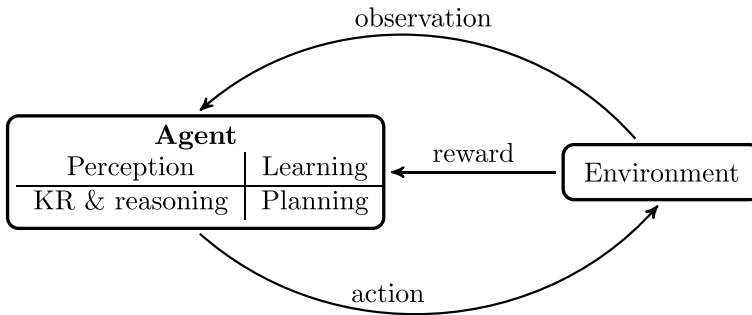


Fig. 1 Interaction loop in RL

The remaining of this survey is organized as follows. In the next section, we recall the necessary definitions and notions related to RL. Next, we discuss the definition of interpretability (and explainability) in the larger context of artificial intelligence (AI) and machine learning (Sect. 3.1), and apply it in the context of RL (Sect. 3.2). In the following sections, we present the studies related to interpretable inputs (Sect. 4) and models (Sect. 5). The work tackling interpretable decision-making, which constitutes the core part of this survey, is discussed in Sect. 6. For the sake of completeness, we also sketch a succinct review of explainable RL (Sect. 7), which helps us contrast it to interpretable RL. Based on this overview, we provide in Sect. 8 a list of open problems and future research directions, which we deem particularly relevant. Finally, we conclude in Sect. 9.

2 Background

In RL, an agent interacts with an environment through an interaction loop. The agent repeatedly receives an observation from the environment, chooses an action, and receives a new observation and usually an immediate reward. Although most RL methods solve this problem by considering the RL agent as reactive (i.e., given an observation, choose an action), Fig. 1 lists some other potential problems that an agent may tackle on top of decision-making: perception if the input is high-dimensional (e.g., image), learning from past experience, knowledge representation (KR) and reasoning, and finally planning if the agent has a model of its environment.

This RL problem is generally modeled as a Markov decision process (MDP) or one of its variants, notably partially observable MDP (POMDP).¹ An MDP is defined as a tuple $(\mathcal{S}, \mathcal{A}, T, R)$ with a set \mathcal{S} of states, a set \mathcal{A} of actions, a transition function $T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$, and a reward function $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$. The sets of states and actions, which may be finite, infinite, or even continuous, specify respectively the possible world configurations for the agent and the possible response that it can perform. In a partially observable MDP, the agent does not observe the state directly, but has access to an observation that probabilistically depends on the hidden state. The difficulty in RL is that the transition and reward functions are not known to the agent. The goal of the agent is to learn to choose actions (i.e., encoded in a policy) such that it maximizes its expected

¹ See Puterman (1994) or Bertsekas and Tsitsiklis (1996) for a more complete discussion.

(discounted) sum of rewards $\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(S_t, A_t) \right]$ where $\gamma \in [0, 1]$ is a discount factor, S_t and A_t are random variables representing the state and action at time step t . The expectation is with respect to the transition probabilities T , the action selection, and a possible distribution over initial states. This expected sum of rewards is usually encoded in a *value function* when the initial states are fixed. A policy may choose actions based on states or observations in a deterministic or randomized way. In RL, the value function often takes the form of a so-called *Q-function*, which measures the value of an action a followed by a policy from a state s , i.e., $\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(S_t, A_t) \mid S_0 = s, A_0 = a \right]$. To solve this RL problem, model-based and model-free algorithms have been proposed, depending on whether a model of the environment (i.e., transition/reward model) is explicitly learned or not.

The success of DRL is explained partly by the use of neural networks (NN) to approximate value functions or policies, but also by various algorithmic progress. Deep RL algorithms can be categorized in two main categories: value-based methods and policy gradient methods, in particular in their actor-critic version. For the first category, the model-free methods are usually variations of the DQN algorithm (Mnih et al., 2015). For the second one, the current state-of-the-art model-free methods are PPO (Schulman et al., 2017) for learning a stochastic policy, TD3 (Fujimoto et al., 2018) for learning a deterministic policy, and SAC (Haarnoja et al., 2018) for entropy-regularized learning of a stochastic policy. Model-based methods can span from simple approaches such as first learning a model and then applying a model-free algorithm using the learned model as a simulator, to more sophisticated methods that leverage the learned model to accelerate solving an RL problem (Francois-Lavet et al., 2019; Scholz et al., 2014; Veerapaneni et al., 2020).

For complex decision-making tasks, hierarchical RL (HRL) (Barto & Mahadevan, 2003) has been proposed to exploit temporal and hierarchical abstractions, which may facilitate learning and transfer, but also promote intelligibility. Although various architectures have been proposed, decisions in HRL are usually made at (at least) two levels. In the most popular framework, a higher-level controller (also called *meta-controller*) chooses temporally-extended macro-actions (also called *options*), while a lower-level controller chooses the primitive actions. Intuitively, an option can be understood as a policy with some starting and ending conditions. When it is known, it directly corresponds to the policy applied by the lower-level controller. An option can also be interpreted as a subgoal chosen by the meta-controller for the lower-level controller to reach.

3 Interpretability and explainability

In this section, we first discuss the definition of interpretability and explainability as proposed in the explainable AI (XAI) literature. Then, we focus on the instantiations of those notions in RL.

3.1 Definitions

Various terms have been used in the literature to qualify the capacity of a model to make itself understandable, such as interpretability, explainability, intelligibility, comprehensibility, transparency, or understandability. Since no consensus about the nomenclature in XAI has been reached yet, they are not always distinguished and are sometimes used interchangeably in past work or surveys on XAI. Indeed, interpretability and explainability are for instance often used as synonyms (Miller, 2019; Molnar, 2019; Ribeiro et al., 2016a).

For better clarity and specificity, in this survey, we only employ the two most common terms, *interpretability* and *explainability*, and clearly distinguish those two notions, which we define below. This distinction allows us to provide a clearer view of the different work in interpretable and explainable RL. Moreover, we use *intelligibility* as a generic term that encompasses those two notions. For a more thorough discussion of the terminology in the larger context of machine learning and AI, we refer the interested readers to surveys on XAI (Barredo Arrieta et al., 2020; Chari et al., 2020; Gilpin et al., 2019; Lipton, 2017).

Following Barredo Arrieta et al. (2020), we simply understand interpretability as a passive quality of a model, while explainability here refers to an active notion that corresponds to any external, usually post-hoc, methodology or proxy aiming at providing insights into the working and decisions of a trained model. Importantly, the two notions are not exclusive. Indeed, explainability techniques can be applied to interpretable models and explanations may potentially be more easily generated from more interpretable models. This is why, one may argue that model interpretability is more desirable than post-hoc explainability (Rudin, 2019). Moreover, both notions are epistemologically inseparable from both the observer and the context. Indeed, what is intelligible and what constitutes a good explanation may be completely different for an end-user, a system designer (e.g., AI engineer or researcher), or a legislator for instance. Except in our discussion on explainable RL, we will generally take the point of view of a system designer.

While interpretability is achieved by resorting to intrinsically more transparent models, explainability requires carrying out additional processing steps to explicitly provide a kind of explanation aiming to clarify, justify, or rationalize the decisions of a trained black-box model. At first sight, it seems that interpretability is involving an *objectual and mechanistic* understanding of the model, whereas explainability mostly restricts itself to a more *functional*.²—and often model-agnostic—understanding of the outcomes of a model. Yet, as advocated by Páez (2019), post-hoc intelligibility in AI should require some degree of objectual understanding³ of the model, since a thorough understanding of a model’s decisions, also encompasses the ability to think counterfactually (“What if...”) and contrastively (“How could I alter the data to get outcome X?”).

Since the main focus of this review is interpretability, we further clarify this notion by recalling three potential definitions as proposed by Lipton (2017): simulatability, decomposability, and algorithmic transparency.

A model is *simulatable* if its inner working can be simulated by a human. Examples of simulatable models are small linear models or decision trees. The concept of simplicity, and quantitative aspects, consequently underlie any definition of simulatability. In that sense, a hypothesis class is not inherently interpretable with respect to simulatability. Indeed, a decision tree may not be simulatable if its depth is huge, whereas a NN may be simulatable if it has only a few hidden nodes. A model is *decomposable* if each of its parts (input, parameter, and calculation) can be understood intuitively. Since a decomposable model assumes its inputs to be intelligible, any simple model based on complex highly-engineered features is not decomposable. Examples of decomposable models are linear models or decision trees using interpretable features. While the other two definitions

² A functional understanding “relies on an appreciation for functions, goals, and purpose” while a mechanistic understanding “relies on an appreciation of parts, processes, and proximate causal mechanisms” (Páez, 2019).

³ Some objectual understanding is particularly beneficial when considering legal accountability and public responsibility.

focuses on the model, the third one shifts the attention to the learning process and requires it to be intelligible. Thus, an algorithm is *transparent* if its properties are well-understood (e.g., convergence). In that sense, standard learning methods for linear regression or support vector machine may be considered transparent. However, since the training of deep learning models is currently still not very well-understood, it results in a regrettable lack of transparency of these algorithms.

Although algorithmic transparency is important, the most relevant notions for this survey are the first two since our main focus is the intelligibility of trained models. More generally, it would be useful and interesting to try to provide finer definitions of those notions, however this is out of the scope of this paper, whose goal is to provide an overview of work aiming at enhancing interpretability in (deep) RL.

3.2 Interpretability in RL

Based on the previous discussion of interpretability in the larger context of AI, we now turn to the RL setting. To solve an RL problem, the agent may need to solve different AI tasks (notably perception, knowledge representation, reasoning, learning, planning) depending on the assumptions made about the environment and the capability of the agent (see Fig. 1).

For this whole process to be interpretable, all its components have arguably to be intelligible, such as: (1) the inputs (e.g., observations or any other information the agent may receive) and its processing, (2) the transition and preference models, and (3) the decision-making model (e.g., policy and value functions). The preference model describes which actions or policies are preferred. It can simply be based on the usual reward function, but can also take more abstract forms such as logic programs. Note that making those components more intelligible supposes a certain disentanglement of the underlying factors, and entails a certain representation structure. With this consideration, this survey can be understood as discussing methods to achieve *structure* in RL, which consequently enhances interpretability.

The three definitions of interpretability (i.e., simulatability, decomposability, and algorithmic transparency) discussed previously can be applied in the RL setting. For instance, for an RL model to be simulatable, it has to involve simple inputs, simple preference (possibly also transition) models, and simple decision-making procedures, which may be hard to achieve in practical RL problems. In this regard, applying those definitions of interpretability at the global level in RL does not lead to any interesting insights in our opinion. However, because RL is based on different components, it may be judicious to apply the different definitions of interpretability to them, possibly in a differing way. A more modular view provides a more revealing analysis framework to understand previous and current work related to interpretable RL. Thus, an RL approach can be categorized for instance, as based on a non-interpretable input model, but simulatable reward and decision-making models (Penkov & Ramamoorthy, 2019) or as based on simple inputs, a simulatable reward model, and decomposable transition and decision-making models (Degrís et al., 2006).

In the end, interpretability is a difficult notion to delineate, whose definition may depend on both the observer and the context. Moreover, it is generally not a Boolean property, but there is a continuum from black-box (e.g., deep NN) to undoubtedly interpretable (e.g., structured program) models. For these reasons, we discuss interpretable methods, but also DRL approaches that are not necessarily considered interpretable, but bring some intelligibility in the RL framework (e.g., via NN architectural bias or

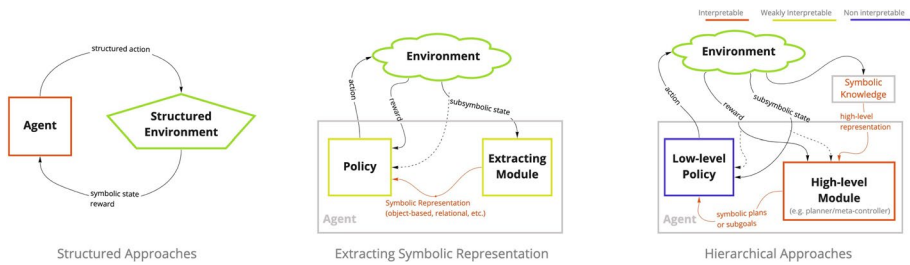


Fig. 2 Illustrations of the different approaches for interpretable inputs: (Left) Structured Approach, (Center) Extracting Symbolic Representation, (Right) Hierarchical Approach. Dashed lines represent optional links, depending on the methods

regularization). We organize these studies along three categories: interpretable inputs, interpretable transition and preference models, and interpretable decision-making. We argue that they are essential aspects of interpretable RL since they each pave the way towards higher interpretability in RL. As in any classification attempt, the boundary between the different clusters is not completely sharp. Indeed, some propositions could arguably belong to several categories. However, to avoid repetition, we generally discuss them only once with respect to their most salient contributions.

4 Interpretable inputs

A first step towards interpretable RL regards the inputs that an RL agent uses to learn and make its decisions. Arguably, these inputs must be intelligible if one wants to understand the decision-making process later on. Note that we define inputs in a very general sense. They can be **any interpretable information** specified by the system designer. Thus, they include the typical (interpretable) RL observations, but also any structural information, which enforces higher intelligibility to the data provided to the agent, such as the *relational or hierarchical* structure of the problem.

Interpretable numeric observations (e.g., kinematic information) can be directly provided by sensors or estimated from high-dimensional observations (e.g., depth from images in Michels et al., 2005). To keep this survey concise, we do not cover such methods, which are more related to the state estimation problem. Instead, we mainly focus on RL approaches exploiting (pre-specified or learned) symbolic and structural information, since they specifically enforce intelligibility in the agents' inputs. Moreover, they can help with faster learning, better generalizability and transferability but also be smoothly integrated with reasoning and planning.

Diverse approaches have been investigated to provide interpretable inputs to the agent (see Fig. 2). They may be pre-given as seen in the literature of structured RL (Sect. 4.1), or may need to be extracted from high-dimensional observations (Sect. 4.2). Tangentially, additional interpretable knowledge can be provided to help the RL agent, in addition to the observations (Sect. 4.3).

Table 1 Overview of *structured approaches*

Relational representations	Approach	Interpretability	References
Relational, Q linear approx.	Linear Programming	Simulatable ⁵	Guestrin et al. (2003)
Relational, FOL Q–decision tree	Q-learning	Simulatable	Džeroski et al. (2001)
Relational, hierarchical, FOL Q–decision tree	Q-learning	Simulatable	Driessens and Blokkeel (2001)
Relational, HOL Q–decision tree	Q-learning	Decomposable	Cole et al. (2003)
Relational, feature-based, Q linear approx.	Q-learning	Partial decomp	Walker et al. (2004)
Relational, feature-based, Q w/ rel. Naive Bayes Net	Q-learning	Partial decomp	Sanner (2005)
Relational, Q w/ graph kernels	Q-learning w/ Gaussian processes	Partial decomp	Driessens et al. (2006)
Relational, Graph NN State rep., Neural policy/value	Actor-Critic	Partial decomp	Garg et al. (2020), Janisch et al. (2021)

Simulatability holds assuming small domains and also implies decomposability here

4.1 Structured approaches

The literature explored in this subsection assumes a pre-given structured representation of the environment which may be modelled through a collection of objects, and their relations as in *object-oriented* or *relational* RL (RRL, Dzeroski et al., 1998). Thus, the MDP is assumed to be structured and the problem is solved within that structure: task, reward, state transition, policies and value function are defined over objects and their interactions—e.g., using first-order logic⁴ (FOL, Barwise, 1977) as in RRL. A non-exhaustive overview of these approaches is provided in Table 1.

A first question, tied to knowledge representation (Swain, 2013), is the choice of the specific structured representations for the different elements (i.e., state, transition, rewards, value function, policy). For instance, a first step in this literature was to depart from *propositional representation*, and turn to relational representations, which not only seems to better fit the way we reason about the environment—in terms of objects and relations—but may bring other benefits, such as the easy incorporation of logical background knowledge. Indeed, in propositional representations, the number of objects is fixed, all relations have to be grounded—a computationally heavy operation—but most importantly it is not suitable to generalize over objects and relations, and is unable to capture the structural aspect of the domain (e.g., Blocks World, Slaney & Thiébaux, 2001). Variations of this structure are presented below.

Relational MDP

⁴ Recall, in contrast to propositional logic (i.e., Boolean vector representation), FOL describes the world in terms of objects, predicates (i.e., relations between objects), and functions (i.e., objects defined from other objects).

Relational MDPs (RMDPs) (Guestrin et al., 2003) are first-order representation of factored MDPs (Boutilier et al., 2000) and are based on probabilistic relational model⁵ (PRM, Koller, 1999). The representation involves different classes of objects (over which binary relations are defined), each having attributes attached to a specific domain. Transition and reward models are assumed given, e.g., as a dynamic Bayesian network (DBN, Dean & Kanazawa, 1990), although the specific representational language may vary.⁶ Closely related, Object Oriented-MDPs (see Diuk et al., 2008, presented in Sect. 5.1)—later extended to deictic representations (Marom & Rosman, 2018)—use similar state-representation yet differ in the way their transition dynamics are described: transitions are assumed deterministic and learned within a specific propositional form in the first step of their algorithm.

Relational RL

Following the initial work on Relational RL (Dzeroski et al., 1998), a consequent line of work summarized below extends previous work dealing with MDPs modelled in a relational language to the learning setting, at the crossroad of RL and logical machine learning—such as inductive logic programming (ILP, Cropper et al., 2020) and probabilistic logic learning (De Raedt & Kimmig, 2015). We also refer the interested readers to the surveys by van Otterlo (2009, 2012).

In Dzeroski et al. (2001), the Q-function is learned with a relational regression tree using Q-learning extended to situations where states, actions, and policies are represented using first-order logic. However, explicitly representing value functions in relational learning is difficult, partly due to *concept drift* (van Otterlo, 2005), which occurs since the policy providing examples for the Q-function is being constantly updated. It may motivate to turn towards policy learning (as Dzeroski et al., 1998 relying on P-trees), and employ approximate policy iteration methods which would keep explicit representation of the policy but not the value function, for larger probabilistic domains.

Diverse extensions of relational MDPs and of Relational RL (RRL) have been proposed, either exact or approximate methods, in model-free and in model-based, with more or less expressive representations and within a plain or more hierarchical approach (Driessens & Blockeel, 2001). Regarding the representations, previous model-free RRL work is based on explicit logical representation such as logical (FOL or more rarely Higher Order Logic (HOL, e.g., Cole et al., 2003)) regression trees, which, in a top-down way, recursively partition the state space; in contrast, other bottom-up and feature-based approaches (Sanner, 2005; Walker et al., 2004) aim to learn useful relational features which they would combine to estimate the value function, either by feeding them to a regression algorithm (Walker et al., 2004), or into a relational naive Bayes network (Sanner, 2005). Other alternatives to regression trees have been implemented such as through Gaussian processes—incrementally learnable Bayesian regression—with graph kernels, defined over a set of state and action (Driessens et al., 2006). Finally, other work in quest of more expressivity turns towards neural representations. For instance, after extracting a graph instance expressed in RDDDL⁷ (Sanner, 2011), Garg et al. (2020) compute nodes embedding via

⁵ PRMs may be understood as “relational” extensions of “propositional” Bayesian networks.

⁶ Relational Dynamic Influence Diagram Language (RDDL, Sanner, 2011), extending DBN using state-dependent rewards aggregated over objects, is able to model parallel effects. In contrast, Probabilistic Planning Domain Definition Language (PPDDL, Younes & Littman, 2004) employs action-transition-based rewards and models correlated effect. Note that Guestrin et al. (2003) assume static representations, which are unfit for real-world dynamics or relational environments such as Blocks World.

Table 2 Overview of approaches for *Learning Symbolic Representations*

Representations	Approach	Interpretability	References
Probabilistic symbols for planning	Unsupervised clustering (e.g., DBSCAN algorithm)	HL modular	Konidaris et al. (2015, 2018)
Probabilistic symbols for planning	Unsupervised clustering (Bayesian hierarchical)	HL modular	Andersen and Konidaris (2017)
Symbols as classifiers	Human-teaching	HL modular	Kulick et al. (2013)
Symbols as classifiers	Program-guided	HL modular	Penkov and Ramamoorthy (2019)
Symbols as classifiers	Program-guided AE	HL modular	Sun et al. (2020)
Objects	Object recognition w/ template matching	Partial decomp.	Li et al. (2017b)
Objects, relational	Unsupervised object extraction w/ activation spectrum	Partial decomp	Garnelo et al. (2016)
Objects	Unsupervised video segmentation w/ optical flow	Weak	Goel et al. (2018)
Relational	Relational MLP-modules	Weak	Adjodah et al. (2018)
Objects	CNN w/ attention	Weak	Zambaldi et al. (2019)

graph propagation steps, which are then fed to value and policy decoders (multi-layer perceptrons, MLP) attached to each action symbol. In Janisch et al. (2021), graph NNs are similarly used to build a relational state representation in relational problems. The authors resort to auto-regressive policy decomposition (Vinyals et al., 2017) to tackle multi-parameter actions (attached to unary or binary predicates).

Let us point out that despite the “reinforcement” appellation, a significant proportion of work in RRL assumes that environment models (transitions and reward structures) are known to the agent, which may be unrealistic. RRL has also been applied to diverse domains, such as for efficient exploration within robotics (Martínez et al., 2017b).

Discussion

In structured MDP and relational RL, by borrowing from symbolic reasoning, most work leads to agents that can learn and reason about objects. Such an explicit and logical representation of learned structures may help both generalize or transfer efficiently and robustly to similar representational frameworks by e.g., reusing learned representations or policies. However, some major drawbacks are that these approaches necessitate the symbolic representation to be hand-designed, and often rely on non-differentiable operations. They are therefore not very flexible over framework variations (e.g., task or input) and not well suited for more complex tasks, or noisy real-life environments.

4.2 Learning symbolic representations

When inputs are given as high-dimensional raw data, it seems judicious—although challenging—to extract explicit symbolic representations on which we can arguably reason and

plan in a more efficient and intelligible way. This process of abstraction, which is very familiar to human cognition,⁷ reduces the complexity of an environment to low dimensional, discrete, abstract features. By abstracting away lower-level details and irrelevant variations, this paradigm brings undeniable advantage and could greatly leverage the learning and generalization abilities of the agent. Moreover, it provides the possibility of reusing high-level features through environments, space and time.

Some promising new research directions—tackling the key problem of symbol grounding (Harnad, 1990)—are adopting an end-to-end training, therefore tying the semiotic emergence not only to control but also to efficient high-level planning (Andersen & Konidaris, 2017; Konidaris et al., 2014, 2015), or model-based learning (Francois-Lavet et al., 2019), to encourage more meaningful abstractions. Work presented below (see Table 2) ranges from extracting symbols to relational representations, which are in turn used for control or planning. In the next two sections, we distinguish actual *high-level (HL) decomposability*—meaning the HL module is decomposable—from *HL modularity*, which denotes the gain in interpretability brought by the task decomposition, which may be seen as a partial high-level decomposability.

Symbol grounding

Some previous approaches (Andersen & Konidaris, 2017; Konidaris et al., 2014, 2015, 2018) have tackled the problem of learning symbolic representations adapted for high-level planning from raw data. As they are concerned about evaluating the feasibility and success probability of a high-level plan, they only need to construct symbols both for the initiation set and the termination set of each option. There, the state-variables are gathered into factors, which can be seen as sub-goals, and are tied to a set of symbols; through unsupervised clustering, each option is attached to a partition of the symbolic state; it leads to a probabilistic distribution over symbolic options (Sutton et al., 1999) which guides the higher-level policy to evaluate the plan.

In a different direction, some researchers have involved human teaching (Kulick et al., 2013) or programs (Penkov & Ramamoorthy, 2019; Sun et al., 2020) to guide the learning of symbols. For instance, in the work by Penkov and Ramamoorthy (2019), the pre-given program mapping the perceived symbols (from an auto-encoder network) to actions, imposes semantic priors over the learned representations, and may be seen as a regularization which structures the latent space. Meanwhile, Sun et al. (2020) present a perception module which aims to answer the conditional queries (“if”) within the program, which accordingly provides a symbolic goal to the low-level controller. However, these studies, by relying on a human-designed program, or human-teaching, partly bypass the problem of autonomous and enacted symbol extraction.

Object Recognition

When the raw input is given as an image, diverse techniques within computer vision and within Object Recognition (OR) or Instance Segmentation (combining semantic segmentation and object localization) are beneficial to extract a symbolic representation. Such extracted information, fed as input to the policy or Q-network, should arguably lead to more interpretable networks. OR aims specifically to find and identify objects in an image or video sequence, despite possible changes in sizes, scales or obstruction when objects are being moved; it ranges from classical techniques—such as template matching (Brunelli, 2009), or Viola-Jones algorithm (Viola & Jones, 2001)—to more advanced ones.

⁷ Physical theories are a typical example of this practice, where laws—such as laws of motions—are reused across instantiations and scenes with various primitive entities.

As an example of recent RL work learning symbolic representation, O-DRL (Iyer et al., 2018; Li et al., 2017b) can exploit and incorporate object characteristics such as the presence and positions of game objects, which are extracted with template matching before being fed into the Q-network. In the case of moving rigid objects, techniques have been developed to exploit information from object movement to further improve object recognition. For instance, Goel et al. (2018) first autonomously detect moving objects by exploiting structure from motion, and then use this information for action selection.

Aiming to delineate a general end-to-end RL framework, Deep Symbolic RL (DSRL, Garnelo et al., 2016) combines a symbolic front end with a neural back end learning to map high-dimensional raw sensor data into a symbolic representation in a lower-dimensional conceptual space. Such proposition may be understood within emblematic neuro-symbolic approaches' scheme (as presented in Bader and Hitzler 2005, Fig.4), where a symbolic system and a connectionist system share information back and forth. The authors also reflect on a few key notions for an ideal implementation such as *conceptual abstraction* (e.g., how to detect high level similarity), *compositional structure*, *common sense priors* or *causal reasoning*. However, their first prototype proposal is relatively limited, with a symbolic front end carrying out very little high-level reasoning, and a simple neural back end for unsupervised symbol extraction. Further work has questioned the generalization abilities of DSRL (Dutra & d'Avila Garcez, 2017) or aimed to incorporate common sense within DSRL (d'Avila Garcez et al., 2018) to improve learning efficiency and accuracy, albeit still within quite restricted settings.

Relational Representations

To obtain more interpretable inputs for the policy or Q-value network, some work deals with specifically relation-centric representations, e.g., graph-based representation. Such relational representation can be leveraged for decision-making, e.g., once fed to the value or policy network or even in a hierarchical setting. Within this line of research, graph networks (Battaglia et al., 2018) stand out as an effective way to compute interactions between entities and can support combinatorial generalization to some extent (Battaglia et al., 2018; Cranmer et al., 2020; Gilmer et al., 2017; Li et al., 2017c; Sanchez-Gonzalez et al., 2018; Scarselli et al., 2009). Roughly, the inference procedure is a form of propagation process similar to a message passing system. Having high capacity, graph networks have been thoroughly exploited in a diverse range of problem domains, either for supervised, unsupervised or in model-free or model-based RL, for tasks ranging from visual scene understanding, to physical systems dynamics via chemical molecule properties, image segmentation, point clouds data, combinatorial optimization, or dynamic of multi-agent systems.

Aiming to represent relations between objects, relational modules have been designed to inform the Q-network (Adjodah et al., 2018), and/or the policy network (Zambaldi et al., 2019). In a work by Zambaldi et al. (2019), the pairwise interactions are computed via a self-attention mechanism (Vaswani et al., 2017), and used to update each entity representation which—according to the authors' claim—is led to reflect important structure about the problem and the agent's intention. Unlike most prior work in relational inductive bias (e.g., Wang et al., 2018), it does not rely on a priori knowledge of the problem and the relations, yet is hard to scale to large input space, suffering from quadratic complexity. Other relational NN modules could be easily incorporated into any RL framework, e.g., (Chang et al., 2017; Santoro et al., 2017) which aim to factorize dynamics of physical systems into pairwise interactions.

Discussion

Extracting symbolic and relational representations from high dimensional raw data is crucial as it would avoid the need of hand-designing the symbolic domain, and could

therefore unlock enhanced adaptability when facing new environments. Indeed, the symbolic way of representing the environment inherently benefits from its compositional and modular perspective, and enables the complexity of the state-space to be reduced thanks to abstraction. Moreover, most of the modules presented in this section may be incorporated as a preprocessing step for any object-oriented or relational RL, either being trained beforehand or more smoothly end-to-end with the subsequent model.

Nevertheless, despite a certain history of work in this area, in the absence of pre-given hand-crafted schemes, it remains a consequent challenge for an agent to autonomously extract relevant abstractions model from a high-dimensional continuous complex and noisy environment. Advanced computer vision techniques could be leveraged in RL, e.g., for object detection (e.g., YOLO, Redmon et al., 2016), tracking (e.g., Deep Sort, Wojke et al., 2017), or for extracting structured representation (e.g., scene graphs, Bear et al., 2020).

For scene interpretation, going beyond the traditional spatial or subsumption (“part-of”) relations, some logic-based approaches have emerged, using decidable fragments of FOL, such as Description Logic (DL), in order to infer new facts in a scene, given basic components (object type or spatial relation), e.g., redefining labels. For instance, Donadello et al. (2017) extending the work of Serafini and d’Avila Garcez (2016) employ fuzzy FOL. Another idea would be to first learn—ideally causally—*disentangled* subsymbolic representations from low-level data (e.g., Higgins et al., 2018), to bootstrap the subsequent learning of higher-level symbolic representations, on which more logical and reasoning-based frameworks can then be deployed. In partially observable domains, an alternative approach is to enforce interpretability of the memory of the agent (Paischer et al., 2023).

As a side note, this line of work touches sensitive questions on how symbols acquire their meanings; e.g., with the well-known *symbol grounding problem*, inquiring on how to ground the representations and symbolic entities from raw observations.⁸ On top of the challenges of “when” and “how” to invent a new symbol, enters also the question of how to assess of its quality.

4.3 Hierarchical approaches

Instead of working entirely within a symbolic structure as in Sect. 4.1, many researchers have tried to incorporate elements of symbolic knowledge with sub-symbolic components, with notable examples within hierarchical RL (HRL, e.g., Hengst., 2010). Their aim was notably to leverage both symbolic and neural worlds, with a more structured high-level and a more flexible low-level. In this type of work, the agent can be understood as taking as inputs this interpretable structural information in addition to its usual observations. One may argue that this hierarchical structure helps make more sense of the low-level high-dimensional observations.

Table 3 introduces these approaches, focusing notably on their high-level interpretability, as their lower-level components—especially when based on neural networks—rarely claim to be interpretable. Indeed, different levels of temporal and hierarchical abstractions within human decision-making arguably participate to make it more intelligible. For instance, Beyret et al. (2019) demonstrate how a meta-controller providing subgoals to a controller achieves both performance and interpretability for robotic tasks.

⁸ A common assumption in contemporary cognitive science is that these representations have to emerge in strong dependency to the actions and goals of the agent (enacted) and the environment (situated).

Table 3 Overview of *Hierarchical Approaches*

Symbolic knowledge	Learned components	HL interpretability	References
Subgoals	Controllers (LL, HL) ^a	Modular	Beyret et al. (2019)
HL Domain	Controller, HL RL Agent	Partial Decomp	Sridharan et al. (2019)
Symbolic Plans	Full & Subpolicies	Modular	Andreas et al. (2017)
STG ^c	Selector, Subpolicies	Modular	Shu et al. (2018)
Model Primitive	Gate, Subpolicies	Modular	Wu et al. (2019a)
HL domain, SP ^c	RL Agent	Simulatable	Leonetti et al. (2016)
HL domain, SP	RL Agent, SP	Simulatable	Yang et al. (2018a)
HL domain, SP	Controllers (LL, HL) ^a Task & Motion Planner	Partial. Decomp	Jiang et al. (2018)
HL domain, SP	Controllers (LL, HL) ^a , SP	Partial. Decomp	Lyu et al. (2019)
HL domain	Subgoal FSA ^b , RL Agent	Simulatable	Furelos-Blanco et al. (2021)
HL domain, FSA ^b	FSA-guided RL Agent	Simulatable	Li et al. (2019)
HL filtering rules	RL Agent	Partial decomp	Zhang et al. (2019)
HL rules	Model-based agent	Partial decomp	Lu et al. (2018)
HL domain, FSA	RL Agent, RS ^d	Partial decomp	Camacho et al. (2019)
HL domain, SP	Controller, RS ^d	Partial decomp	Grzes and Kudenko (2008)

^aLow-level High Level

^bFinite State Automaton

^cSymbolic Planner

^dReward Shaping

^eStochastic Temporal grammar may be given as a prior, or trained

Working at a higher level of abstraction than the controller, it seems reasonable that the high-level module (e.g., meta-controller, planner) manipulates symbolic representations and handles the reasoning part, while the low-level controller could still benefit from the flexibility of neural approaches. Yet, we could also imagine less dichotomic architectures: for instance, another neural module may coexist at the high-level along with the logic-based module to better inform the decisions under uncertainty (as in Sridharan et al., 2019).

Various symbolic domain knowledge may be incorporated to HRL frameworks, in order to leverage both learning and high-level reasoning: high-level domain-knowledge (e.g., with high-level transitions and mappings from low-level to high level), high-level plans, task-decomposition, or specific decisions rules (e.g., for safety filtering). High-level domain knowledge may be described through symbolic logic-based or *action language* such as PDDL or RDDDL (Santer, 2011) or via temporal logic (Camacho et al., 2019; Li et al., 2019).

Modularity

Modular approaches in HRL have been relevantly applied to decompose a possibly complex task domain into different regions of specialization, but also to multitask DRL, in order to reuse previously learned skills across tasks (Andreas et al., 2017; Shu et al., 2018; Wu et al., 2019a). Tasks may be annotated by handcrafted instructions (as “policy sketches” in Andreas et al., 2017), and symbolic subtasks may be associated with subpolicies which a full task-specific policy aims to successfully combine. Distinctively, Shu et al.

(2018) propose to train (or provide as a prior) a stochastic temporal grammar (STG), in order to capture temporal transitions between tasks; an STG encodes priorities of sub-tasks over others and enables to better learn how to switch between base or augmented subpolicies. The work in Wu et al. (2019a) seeks to learn a mixture of subpolicies by modeling the environment assuming given a set of (imperfect) models specialized by regions, referred to as *model primitives*. Each policy is specialized on those regions and the weights in the mixture correspond to the posterior probability of a model given the current state. Echoing the hierarchical abstractions involved in human decision-making, such modularity and task-decomposability—referred to as *high-level modularity* as previously—would arguably participate to make decision-making more intelligible.

Symbolic Planning

A specific line of work within HRL has emerged trying to fuse symbolic planning (SP, Cimatti et al., 2008) with RL (SP+RL), to guide the agent's task execution and learning (Leonetti et al., 2016). In classical SP, an agent uses a symbolic planner to generate a sequence of symbolic actions (*plan*) based on its symbolic knowledge. Yet, this pre-defined notion of planning seems unfit to most RL or real-world domains which present both domain uncertainties and execution failures. Recent work then usually interleaves RL and SP, aiming to send feedback signals to the planner in order to handle such scenarios. Planning agents often carry prior knowledge of the high-level dynamics, typically hand-designed, and assumingly consistent with the low-level environment.⁹

Framing SP in the context of automatic option discovery, in PEORL (Yang et al., 2018a), a constraint answer set solver generates a symbolic plan, which is then turned into a sequence of options to guide the reward-based learning. Unlike earlier work in SP-RL (e.g., Leonetti et al., 2016), RL is intertwined with SP, such that more suitable options could be selected. Some work has extended PEORL, with two planning-RL loops for more robust and adaptive task-motion planning (Jiang et al., 2018), or with an additional meta-controller in charge of subtask evaluation to propose new intrinsic goals to the planner (Lyu et al., 2019). Recent work (Jin et al., 2022) starts to deal with learning action models in the RL loop.

Declarative Domain Knowledge

Declarative and common sense knowledge has been incorporated in RL frameworks in diverse ways to guide exploration, such as to filter out unreasonable or risky actions with finite state automaton (Li et al., 2019) or high-level rules (Zhang et al., 2019). In contrast, Furelos-Blanco et al. (2021) propose to learn a finite-state automaton for the higher-level with an ILP method and solve the lower-level with an RL method. There, the automaton is used to generate subgoals for the lower level. A deeper integration of knowledge representation and reasoning with model-based RL has been advocated in Lu et al. (2018), where the learned dynamics are fed into the logical-probabilistic reasoning module to help it select a task for the controller. In a different direction, exploiting declarative knowledge to construct actions sequences can also help reward shaping to find the optimal policy, as a few studies (Camacho et al., 2019; Grzes & Kudenko, 2008) have demonstrated. We refer to the survey by Zhang and Sridharan (2020) for further examples of studies both in probabilistic planning and RL aiming to reason with declarative domain knowledge.

Discussion

⁹ In contrast, the work in RL+SP mentioned in Sect. 4.2 does not assume similar HL domain knowledge, and aims to learn the mapping from the low-level domain to high-level symbols.

Symbolic knowledge and reasoning elements have been consistently incorporated to (symbolic) planning or hierarchical decision-making, in models which may reach a certain high-level interpretability, albeit neglecting action-level interpretability. Moreover, most work still typically relies on manually-crafted symbolic knowledge—or has to rely on a (pretrained or jointly-trained) perception model for symbol grounding—assuming a pre-given symbolic structure hand-engineered by a human expert. Due to the similarities shared by the majority of discrete dynamic domains, some researchers (Lyu et al., 2019) have argued that a laborious crafting of symbolic model is not always necessary, as the symbolic formulation could adapt to different problems, by instantiating new types of objects and a few additional rules for each new task; such claim would still need to be backed by further work in order to demonstrate such flexibility. Moreover, when adopting a symbolic or logical high-level framework, one also needs to face a new trade-off between expressivity and complexity of the symbolic representation.

5 Interpretable transition/preference models

In this section, we overview the work that focuses on exploiting an interpretable model of the environment or task. This model can take the form of a transition model (Sect. 5.1) or preference model (Sect. 5.2). Such interpretable models can help an RL agent reason about its decision-making, but also help humans understand and explain its decision-making. As such, they can be used in an interpretable RL algorithm, but also in a post-hoc procedure to explain the agent's decision-making. Note that those models may be learned, or not, and, when not learned, may possibly be fully provided to the RL agent or not. For instance, the reward function, which is one typical way of defining the preference model, is generally specified by the system designer in order to guide the RL agent to learn and perform a specific task. This function is usually not learned directly by the RL agent (except in inverse RL, Ng & Russell, 2000), but still has to be intelligible in some sense, otherwise the agent's decision-making may be based on spurious reasons and the agent may not accomplish the desired task since a non-intelligible reward function is hard to verify and may be incorrectly specified.

5.1 Interpretable transition models

Interpretable transition models can help discover the structure and potential decomposition in a problem that are useful for more data-efficient RL (via e.g., more effective exploration), but also allow for larger generalizability and better transfer learning. The work discussed here either solely focuses on model learning or belongs to model-based RL. Various interpretable representations have been considered for learning transition models, such as decision trees or graphical models for probabilistic models, physics-based or graph for deterministic models, or NNs with architectural inductive bias. The NN-based approaches, which are more recent, are presented separately to emphasize them. We provide an overview of all the methods for interpretable transition models in Table 4.

Probabilistic Models

While the setting of factored MDPs generally assumes that the structure is given, Degris et al. (2006) propose a general model-based RL approach that can both learn the structure of the environment and its dynamics. This method is instantiated with decision

Table 4 Overview of approaches for *Interpretable Transition Models*

Type	Model	Interpretability	References
Probabilistic	Decision tree	Simulatable	Degrís et al. (2006)
Probabilistic	Noisy deictic rule	Decomposable	Pasula et al. (2007)
Probabilistic	First-order rule	Decomposable	Walker et al. (2008)
Probabilistic	Relational action schema	Decomposable	Walsh (2010)
Probabilistic	Relational planning operator	Decomposable	Martínez et al. (2016, 2017a)
Probabilistic	Weighted labeled multigraph	Decomposable	Metzen (2013)
Probabilistic	Graphical model	Decomposable	Kansky et al. (2017)
Probabilistic	Gaussian process	Decomposable	Kaiser et al. (2019)
Deterministic	Object-oriented representation	Decomposable	Diuk et al. (2008)
Deterministic	Deictic object-oriented rep	Decomposable	Miarom and Rosman (2018)
Deterministic	Physics engine	Decomposable	Scholz et al. (2014)
Deterministic	Graph of transitions between user-defined Boolean attributes	Decomposable	Zhang et al. (2018a)
Deterministic	State-space graph	Decomposable	Eysenbach et al. (2019)
Neural	Graph NNs	Partial decomp	Battaglia et al. (2016)
Neural	Object dynamic predictor	Partial decomp	Sanchez-Gonzalez et al. (2018)
Neural	from pixels	Partial decomp	Finn et al. (2016)
Neural	Object dynamic predictor with relations	Partial decomp	Finn and Levine (2017)
Neural	Object-level dynamics model with unsupervised learning	Partial decomp	Zhu et al. (2018), Zhu et al. (2020)
Neural	Models for entity grounding dynamics & observation distrib	Partial decomp	Agnew and Domingos (2018)
			Veerapaneni et al. (2020)

trees to represent the environment model. Statistical χ^2 tests are used to decide for the structural decomposition.

Various approaches are based on generative models under the form of graphical models. In Metzén (2013), a graph-based representation of the transition model is learned in continuous domains. In Kansky et al. (2017), schema networks are proposed as generative models to represent the transition and reward models in a problem described in terms of entities (e.g., objects or pixels) and their attributes. As a graphical model, the authors explain how to learn its structure and how to use it for planning using inference. The work in Kaiser et al. (2019) learns via variational inference an interpretable transition model by encoding high-level knowledge in the structure of a graphical model.

In the relational setting, a certain number of studies explored the idea of learning a relational probabilistic model for representing the effects of actions (see Walsh, 2010 for discussions of older work). In summary, those approaches are either based on batch learning (e.g., (Pasula et al., 2007; Walker et al., 2008)) or online methods (e.g., Walsh, 2010) with or without guarantees using more or less expressive relational languages. Some recent work (Martínez et al., 2016, 2017a) proposes to learn a relational probabilistic model via ILP and uses optimization to select the best planning operators.

Deterministic Models

In Diuk et al. (2008), an efficient model-based approach is proposed for an object-oriented representation of the world. This approach is extended by Marom and Rosman (2018) to deictic object-oriented representations, which use partially grounded predicates, in the KWIK framework (Walsh, 2010).

Alternatively, Scholz et al. (2014) explore the use of a physics engine as a parametric model for representing the deterministic dynamics of the environment. The parameters of this engine are learned by a Bayesian learning approach. Finally, the control problem is solved using the A* algorithm.

In a hierarchical setting, several recent studies have proposed to rely on search algorithms on state-space graphs or planning algorithms for the higher-level policy. Given a mapping from the state space to a set of binary high-level attributes, Zhang et al. (2018a) learn a model of the environment predicting if a low-level policy would successfully transition from an initial set of binary attributes to another set. The low-level policy observes the current state and the desired set of attributes to reach. Once the transition model and policy are learned, a planning module can be applied to reach the specified high-level goals. Thus, the high-level plan is interpretable, but the low-level policy is not.

Eysenbach et al. (2019) extract a state-space graph from a replay buffer and apply the Dijkstra algorithm to find a shortest path to reach a goal. This graph represents the state space for the high level, while the low level is dealt with a goal-conditioned policy. The approach is validated in navigation problems with high-dimensional inputs.

NN-based Model

Various recent propositions have tried to learn dynamics model using NNs with specific architectural inductive bias taking graphs as inputs (Battaglia et al., 2016; Sanchez-Gonzalez et al., 2018). While this line of work can provide high-fidelity simulators, the learned model may suffer from a lack of interpretability.

The final set of work we would like to mention aims at learning object-based dynamics models from low-level inputs (e.g., frames) using NNs. In that sense, they can be understood as an extension of the work in relational domain where the input is now generally high-level. One early work (Finn et al., 2016; Finn & Levine, 2017) tries to take

Table 5 Overview of approaches for *Interpretable Preference Models*

Model	Approach	Interpretability	References
Relational	Inverse RL	Simulatable	Munzer et al. (2015)
Relational	Active learning	Decomposable	Martínez et al. (2017a)
Deep decision trees	Batch adversarial inverse RL	Partial decomp	Srinivasan and Doshi-Velez (2020)
Tree structure	Active preference learning inverse RL	Decomp	Bewley and Lecue (2022)
LTL	Pre-specified	Simulatable	Aksaray et al. (2016)
Geometric LTL	Pre-specified	Simulatable	Littman et al. (2017)
Truncated LTL	Pre-specified	Simulatable	Li et al. (2017a), Li et al. (2019)
LTL	Pre-specified	Simulatable	Toro Icarte et al. (2018b)
Finite-state machine	Pre-specified	Simulatable	Toro Icarte et al. (2018a)
Formal languages	Pre-specified	Simulatable	Camacho et al. (2019)
LTL	Pre-specified	Simulatable	Hasanbeig et al. (2020)
Finite-state machine	Local search	Simulatable	Toro Icarte et al. (2019)
Finite-state machine	Automata learning	Simulatable	Xu et al. (2020), Gaon and Brafman (2020), Corazza et al. (2022)
Symbolic plan	Pre-specified	Simulatable	Illanes et al. (2020)
Boolean task algebra	Pre-specified	Simulatable	Tasse et al. (2020)

into account moving objects. However, the model predicts pixels and does not take into account relations between objects, which limits its generalizability. Zhu et al. (2018) propose a novel NN, which can be trained in an unsupervised way, for object detection and object dynamic prediction conditioned on actions and object relations. This work has been extended to deal with multiple dynamic objects (Zhu et al., 2020).

Another work (Agnew & Domingos, 2018) proposes an unsupervised method called Object-Level Reinforcement Learner (OLRL), which detects objects from pixels and learns a compact object-level dynamics model. The method works according to the following steps. Frames are first segmented into blobs of pixels, which are then tracked over time. An object is defined as blobs having similar dynamics. Dynamics of those objects are then predicted with a gradient boosting decision tree.

In Veerapaneni et al. (2020) an end-to-end object-centric perception, prediction, and planning (OP3) framework is developed. The model has different components jointly-trained: for (dynamic) entity grounding, for modeling the dynamics and for modeling the observation distribution. The variable binding problem is treated as an inference problem: being able to infer the posterior distribution of the entity variables given a sequence of observations and actions. One further specificity of this work is to model a scene not globally but locally, i.e., for each entity and its local interactions (locally-scoped entity-centric functions), avoiding the complexity to work with the full combinatorial space, and enabling generalization to various configurations and number of objects.

Discussion

Diverse approaches have been considered for learning an interpretable transition model. They are designed to represent either deterministic or stochastic environments. Recent work is based on NNs in order to process high-dimensional inputs (e.g., images) and has adopted an object-centric approach. While NNs hinder the intelligibility of the method, the

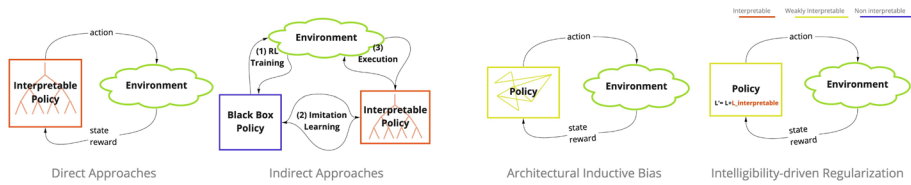


Fig. 3 Illustrations of the different approaches for interpretable decision-making. From left to right: Direct Approach, Indirect Approach, Architectural Inductive Bias, Intelligibility-driven Regularization

decomposition into object dynamics can help scale and add transparency to the transition model.

5.2 Interpretable preference models

In RL, the usual approach to describe the task to be learned or performed by an agent consists in defining suitable rewards, which is often a hard problem to solve for the system designer. The difficulty of reward specification has been recognized early (Russell, 1998). Indeed, careless reward engineering can lead to undesired learned behavior (Randlov & Alstrom, 1998) and to value misalignment (Arnold et al., 2017). Different approaches have been proposed to circumvent or tackle this difficulty: imitation learning (or behavior cloning) (Osa et al., 2018; Pomerleau, 1989), inverse RL (Arora & Doshi, 2018; Ng & Russell, 2000), learning from human advice (Kunapuli et al., 2013; Maclin & Shavlik, 1996), or preference elicitation (Rothkopf & Dimitrakakis, 2011; Weng et al., 2013). Table 5 summarizes the related work for interpretable preference models. Note that some of the approaches in Sect. 5.1 apply here as well (e.g., Walker et al., 2008).

Closer to interpretable RL, Munzer et al. (2015) extend inverse RL to relational RL, while Martínez et al. (2017a) learn from demonstrations in relational domains by active learning. In Srinivasan and Doshi-Velez (2020), more interpretable rewards are learned using tree-structured representations (Bewley & Lecue, 2022; Yang et al., 2018b).

Recent work has started to investigate the use of temporal logic (or variants) to specify an RL task (Aksaray et al., 2016; Littman et al., 2017; Li et al., 2017a, 2019). Related to this direction, Kasenberg and Scheutz (2017) investigate the problem of learning from demonstration an interpretable description of an RL task under the form of linear temporal logic (LTL) specifications.

Another related work Toro Icarte et al. (2018a) propose to specify and represent a reward function as a finite-state machine, called reward machine, which clarifies the reward function structure. Reward machines can be specified using inputs in a formal language, such as LTL (Camacho et al., 2019; Hasanbeig et al., 2020). The model has also been extended to stochastic reward machines (Corazza et al., 2022). Reward machines can also be learned by local search (Toro Icarte et al., 2019) or with various automata learning techniques (Gaon & Brafman, 2020; Xu et al., 2020). Inspired by reward machine, Illanes et al. (2020) propose the notion of taskable RL, where RL tasks can be described as symbolic plans.

Using a different approach, Tasse et al. (2020) show how a Boolean task algebra, if such structure holds for a problem, can be exploited to generate solutions for new tasks by task composition (Todorov, 2009). Such approach, which was extended to the lifelong

RL setting (Tasse et al., 2022), can arguably provide interpretability of the solutions thus obtained.

Discussion

As can be seen, research in interpretable preference representation has been less developed than for transition models. However, we believe interpretability of preference models is as important if not more than interpretability for describing the environment dynamics, when trying to understand the action selection of an RL agent. Thus, more work is needed in this direction to obtain more transparent systems based on RL.

6 Interpretable decision-making

We now turn to the main part of this survey paper, which deals with the question of interpretable decision-making. Among the various approaches that have been explored to obtain interpretable policies (or value functions), we distinguish four main families (see Fig. 3). Interpretable policies can be learned directly (Sect. 6.1) or indirectly (Sect. 6.2), and in addition, in DRL, interpretability can also be enforced or favored at the architectural level (Sect. 6.3) or via regularization (Sect. 6.4).

6.1 Direct approaches

Work in the direct approach aims at directly searching for a policy in a policy space chosen and accepted as interpretable by the system designer. They can be categorized according to their search space and the method to search in this space (Table 6). Several models have been considered in the literature:

Decision Trees

A decision tree is a directed acyclic graph where the nodes can be categorized into decision nodes and leaf nodes. It is interpretable by nature, but learning it can be computationally expensive. The decision nodes will determine the path to follow in the tree until a leaf node is reached, this selection is mostly done according to the state features. Decision trees can represent value functions or policies. For instance, some older work (Ernst et al., 2005) uses a decision tree to represent the Q-value function where each leaf node represents the Q-value of an action in a state. Their optimization method is based on decision tree-based supervised learning methods which do not rely on differentiability.

In contrast, Likmeta et al. (2020) propose to learn parameterized decision nodes. In this approach, the policy instead of the value function is represented by the decision tree where a leaf represents the action to take. The structure of the tree is assumed to be given by experts. To update the tree parameters, policy gradient with parameter-based exploration is employed. Similarly, Silva et al. (2020) design a method to discretize differentiable decision trees such that policy gradient can be used during learning. Therefore, the whole structure of the tree can be learned. The analysis in Silva et al. (2020) also suggests that representing the policy instead of the value function with a decision tree is more beneficial. Expert or other prior knowledge may also bootstrap the learning process, as demonstrated by Silva and Gombolay (2020), where the policy tree is initialized from human-provided knowledge, before being dynamically learned. In addition, Topin et al. (2021) introduce a method defining a meta-MDP from a base MDP with additional actions where any policy in the meta-MDP can be transformed in a decision tree policy in the base MDP. In this way,

Table 6 Overview of *Direct Approaches*

Type	Approach	Interpretability	References
Decision tree	Tree-based algorithms	Simulatable	Ernst et al. (2005)
Decision tree	Gradient descent	Simulatable	Gupta et al. (2015), Likmeta et al. (2020), Pace et al. (2022), Silva et al. (2020), Silva and Gombolay (2020), Topin et al. (2021)
Formulas	Multi-armed bandit with depth search	Decomposable	Maes et al. (2012a, 2012b)
Formulas	Genetic algorithms	Decomposable	Hein et al. (2018, 2019)
Formulas	Gradient descent	Decomposable	Ault et al. (2020)
Fuzzy controllers	Gradient descent	Decomposable	Akrour et al. (2019)
Fuzzy controllers	Particle swarm	Decomposable	Hein et al. (2017)
Logic Rules	Gradient Descent	Decomposable	Delfosse et al. (2023), Evans and Grefenstette (2018), Glanois et al. (2022), Jiang and Luo (2019), Payani and Fekri (2019a), Payani and Fekri (2019b), Payani and Fekri (2020), Zimmer et al. (2021)
Chained Logic Rules	Gradient descent	Decomposable	Ma et al. (2020), Yang and Song (2019)
Programs	Gradient descent	Decomposable	Anderson et al. (2020), Qiu and Zhu (2022), Verma et al. (2019)
Logic Programs	Program synthesis	Decomposable	Cao et al. (2022)
Graphical Models	Gradient descent	Decomposable	Chen et al. (2020), Levine (2018)

the meta-MDP can be solved by classic DRL algorithms. Recently, Pace et al. (2022) propose to learn a tree-based policy in the offline and partially-observable setting.

In a different approach, Gupta et al. (2015) propose to learn a binary decision tree where each leaf is itself a parametric policy. Linear Gibbs softmax policies are learned in discrete action spaces, while in continuous action spaces Gaussian distributions are learned. These parametric policies remain interpretable since their parameters are directly interpretable (probabilities for Gibbs softmax policies, mean and standard deviation for Gaussian distribution) and do not depend on states. Hence, the composition of the decision tree with the parametric policies is interpretable. Policy gradient is employed to update the parametric policies and to choose how the tree should grow.

Formulas

Maes et al. (2012a) represent the Q-value function (used to define a greedy policy) with a simple closed-form formula constructed from a pre-specified set of allowed components: binary operations (addition, subtraction, multiplication, division, minimum, and maximum), unary operations (square root, logarithm, absolute value, negation, and inverse), variables (components of states or actions, both described as vectors), and a fixed set of constants. Because of the combinatorial explosion, the total number of operators, constants, and variables occurring in a formula was limited to 6 in their experiments. To search among this space, the authors formulate a multi-armed bandit problem and used a depth-limited search approach (Maes et al., 2012b).

In Hein et al. (2018, 2019), a formula is used to directly represent a policy. In this work, expressivity is improved by adding more operators (tanh, if, and, or) and deeper formulas (a maximum depth of 5 and around 30 possible variables). Genetic programming is used to search for the formula when a batch of RL transitions is available.

A different approach is proposed in the context of traffic light control by Ault et al. (2020) who design a dedicated interpretable polynomial function where the parameters are learned by a variant of DQN. This function is then used similarly to a Q-value function to derive a policy.

Fuzzy controllers

Fuzzy controllers define a policy as a set of fuzzy rules of the form: “IF *fuzzy_condition(state)* DO *action*.” Akrouf et al. (2019) assume that a state is categorized in a discrete number of clusters with a fuzzy membership function. Then *fuzzy_condition(state)* is defined as a distance to a centroid. A policy is defined as a Gaussian distribution such that the closer a state is to a centroid, the more the mean associated to the centroid is taken into account for the global mean. The mean associated to each cluster is learned via policy gradient given a non-interpretable critic. Similarly, Hein et al. (2017) learn fuzzy rules for deterministic policies with particle swarm optimization in continuous action domains. In both approaches, the number of rules (and clusters) are adapted automatically.

Logic Rules

Neural Logic Reinforcement Learning (NLRL Jiang & Luo, 2019) aims at representing policies by first-order logic. NLRL combines policy gradient methods with a new differentiable ILP architecture adapted from Evans and Grefenstette (2018). All the possible rules are generated given expert-designed rule templates. To represent the importance of rules in the deduction, a weight is associated to each rule. As all rules are applied with a softmax over their weights, the resulting predicate takes its value over the continuous interval [0, 1] during learning. Such approach is able to generalize to domains with more objects than it was trained on, but computing all the applications of all possible rules during training is costly.

In the previous approach, to limit the number of possible rules, the templates are generally formulated such that the number of atoms in the body of a rule is restricted to two. To overcome this limitation, Payani and Fekri (2019a, 2019b) design an alternative model, enforcing formulas to be in disjunctive normal form, where weights are associated to atoms (instead of rules) in a clause and extend it to RL (Payani & Fekri, 2020). Similarly, Zimmer et al. (2021) also define weights associated to atoms. However, the architecture proposed by Dong et al. (2019) is adapted to enforce interpretability and relies on a Gumbel-Softmax distribution to select the arguments in a predicate. This approach can be more interpretable than previous similar work, since it can learn a logic program instead of a weighted combination of logic formulas. Recently, alternative approaches (Delfosse et al., 2023; Glanois et al., 2022) have been further explored.

Alternatively, Yang and Song (2019) propose a differentiable ILP method extending multi-hop reasoning (Lao & Cohen, 2010; Yang et al., 2017). Instead of performing forward-chaining on predefined templates, weights are associated to every possible relational paths where a path corresponds to a multi-step chain-like logic formula. Compared to previous work, it is less expressive since it is not able to represent full Horn clauses, but has a better scalability. This approach is extended by Ma et al. (2020) to the RL setting.

Programs

Verma et al. (2019) propose a novel approach to learn directly a policy written as a program. Their approach can be seen as inspired by (constrained functional) mirror descent. Thus, their algorithm iteratively updates the current policy using a gradient step in the continuous policy space that mixes neural and programmatic representations, then projects the resulting policy in the space of programmatic policies via imitation learning. This approach is extended by Anderson et al. (2020) to safe RL in order to avoid unsafe states during exploration with formal verification. Recently, Qiu and Zhu (2022) propose a differentiable method to learn a programmatic policy, while Cao et al. (2022) propose a framework to synthesize hierarchical and cause-effect logic programs.

Graphical Models

Most previously-discussed work uses a deterministic interpretable representation. However, graphical models can also be considered interpretable. Thus, for instance, in the context of autonomous driving, Chen et al. (2020) solve the corresponding DRL problem as a probabilistic inference problem (Levine, 2018): for both the RL model and policy, they learn graphical models with hidden states, which are trained to be interpretable by enforcing semantic meanings available at training time. The drawback of this approach is that it can only provide interpretability to the learned latent space.

Discussion

Using the direct approach to find interpretable policies is hard since we must be able to solve two potentially conflicting problems at the same time: (1) finding a good policy for the given (PO)MDP and (2) keeping that policy interpretable. These two objectives may become contradictory when the RL problems are large, resulting in a scalability issue with the direct approach. Most work learning a fully interpretable policy focuses only on small toy problems.

The direct approach is related to discrete optimization where the objective function is not differentiable and looking for a policy in such a space is very difficult. Another limitation of these approaches is their poor robustness to noise. To overcome those issues, several approaches use a continuous relaxation to make the objective function differentiable (i.e., search in a smoother space) and more robust to noise, but the scalability issue remains open.

6.2 Indirect approach

In contrast to the direct approach, the indirect approach follows two steps: first train a non-interpretable policy with any efficient RL algorithm, then transfer this trained policy to an interpretable one. Thus, this approach is related to imitation learning (Hussein et al., 2017) and policy distillation (Rusu et al., 2016). Note a similar two-step approach can be found in post-hoc explainability for RL. However, a key difference concerns how the obtained interpretable policy is used, either as a final controller or as a policy that explains a black-box controller, which leads to different considerations about how to learn and evaluate such interpretable policy (see Sect. 7). Similarly to RL algorithms that can be subdivided into value-based methods and policy-based methods, the focus in the indirect approach may be to obtain an interpretable representation of either a learned Q-value function (which provides an implicit representation of a policy) or a learned policy (often called *oracle*), although most work focuses on the latter case. Regarding the types of interpretable policies, decision trees (or variants) are often chosen due to their interpretability, however other representations like programs have also been considered.

Decision Trees and Variants

Liu et al. (2018)'s work is representative among the value-based methods using decision trees. The authors introduce Linear Model U-trees (LMUTs) to approximate Q-functions estimated by NNs in DRL. LMUTs is based on U-tree (Maes et al., 1996), which is a tree-structured representation specifically designed to approximate a value function. A U-Tree, whose structure and parameters are learned online, can be viewed as a compact decision tree where each arc corresponds to the selection of the feature of a current or past observation, and each path from the root to a leaf represents a cluster of observation histories having the same Q-values. LMUTs extend U-Trees by having in each leaf a linear model, which is trained by stochastic gradient descent. Although LMUT is undoubtedly a more interpretable model than a NN, it shows its limit when dealing with high-dimensional features spaces (e.g., images). In Liu et al. (2018), rules extraction and super-pixels (Ribeiro et al., 2016b) are used to explain the decision-making of the resulting LMUT-based agent.

Many policy-based methods propose to learn a decision tree policy. The difficulty of this approach is that a high-fidelity policy may require a large-sized decision tree. To overcome this difficulty, Bastani et al. (2018) present a method called VIPER that builds on DAGGER (Ross et al., 2011), a state-of-the-art imitation learning algorithm, but exploits the available learned Q-function. The authors show that their proposition can achieve comparable performance to the original non-interpretable policy, and is amenable to verification. As an alternative approach to control the decision tree size, Roth et al. (2019) propose to increase its size only if the novel decision tree increases sufficiently the performance. As an improvement to work like VIPER using only one decision tree, Vasic et al. (2019) use a mixture of Expert Trees (MOET). The approach is based on a gating function that partitions the state space and then within each partition, a decision tree expert (via VIPER) approximates the policy.

For completeness, we mention a few other relevant studies, mostly based on imitation learning: Natarajan et al. (2011) learn a set of relational regression trees in relational domains by functional gradient boosting; Cichosz and Pawełczak (2014) learn decision tree policies for car driving; Nagesh Rao et al. (2019) extract a set of fuzzy rules from a neural oracle.

Table 7 Overview of approaches with *Architectural Inductive Bias*

Inductive bias	Model	Intelligibility	References
Relational	Graph NN	Partial decomp	Wang et al. (2018)
Logical	Modular architecture: MLPs wired w/ tensor operators	Weak partial decomp	Dong et al. (2019)
Attention	Self-attention bottleneck LSTM controller	Partial decomp	Tang et al. (2020)
Attention	ConvLSTM, attention module, LSTM controller	Partial decomp	Mott et al. (2019)
Attention	DQN architecture w/ attention	Partial decomp	Annasamy and Sycara (2019)

Programs

As an alternative to decision trees, Verma et al. (2018) introduce Programmatically Interpretable RL (PIRL) to generate policies represented in a high-level, domain-specific programming language. To find a program that can reproduce the performance of a neural oracle, they propose a new method, Neurally Directed Program Search (NDPS). NPDS performs a local search over the non-smooth space of programmatic policies in order to minimize a distance from this neural oracle computed over a set of adaptively chosen inputs. To restrict the search space, a policy sketch is assumed to be given. Unlike the imitation learning setting where the goal is to match the expert demonstrations perfectly, a key feature of NPDS is that the expert trajectories only guide the local program search in the program space to find a good policy.

Zhu et al. (2019) also propose a search technique to find a program to mimic a trained NN policy for verification and shielding (Alshiekh et al., 2018). The novelty in their approach is to exploit the information of safe states, assumed to be given. If a generated program is found to be unsafe from an initial state, this information is used to guide the generation of subsequent programs.

In Burke et al. (2019), a method is proposed to learn a program from demonstration for robotics tasks that are solvable by applying a sequence of low-level proportional controllers. In a first step, the method fits a sequence of such controllers to a demonstration using a generative switching controller task model. This sequence is then clustered to generate a symbolic trace, which is then used to generate a programmatic representation by a program induction method.

Finally, although strictly speaking not a program, Koul et al. (2019) propose to extract from a trained recurrent NN policy a finite-state representation (i.e., Moore machine) that can approximate the trained policy and possibly match its performance by fine-tuning if needed. This representation is arguably more interpretable than the original NN.

Discussion

As mentioned previously, the direct approach requires tackling simultaneously two difficulties: (1) solve the RL problem and (2) obtain an interpretable policy. In contrast to the direct approach, the indirect approach circumvents the first above-mentioned difficulty at the cost of solving two consecutive (hopefully easier) problems: (1) solve the RL problem with any efficient RL algorithm, (2) mimic the good learned policy with an interpretable one by solving a supervised learning problem. Therefore, any imitation learning (Hussein et al., 2017) and policy distillation (Rusu et al., 2016) methods could be applied to

obtain an interpretable policy in the indirect approach. However, the indirect approach is more flexible than standard imitation learning because of the unrestricted access to (1) an already-trained expert policy using the same observation/action spaces, and (2) its value function as well. The teacher-student framework (Torrey & Taylor, 2013) fits particularly well this setting. As such, it would be worthwhile to investigate the applications of techniques proposed for this framework (e.g., Zimmer et al., 2014) to the indirect approach.

In addition, the work by Verma et al. (2019) seems to be a promising approach to combine the direct and indirect approaches. While the authors show that the performance of their proposition outperforms NDPS, it is currently still not completely clear which of a direct method or an indirect one should be preferred to learn good interpretable policies.

6.3 Architectural inductive bias

To favor interpretable decision-making, specific architectural choices may be adopted for the policy network or value function, may it be through relational, logical, or attention-based bias; some examples are presented in Table 7.

Relational Inductive Bias

Such bias refers to inductive bias imposing constraints on relationships and interactions among entities in a learning process. It can take various forms ranging from convolutional NNs to Graph NNs (GNN) as mentioned in Section 4.2 for representation learning, here designed for the policy network. While still difficult to apprehend since the approach is based on NNs, enforcing specific structures in the NN architecture arguably makes the model (and thus the decision-making computed by it) more interpretable. A representative example is NerveNet (Wang et al., 2018) which aims at learning a structured policy—parametrized as a GNN, and executing some graph propagation steps.

Logical Inductive Bias

For instance, Neural Logic Machine (NLM) (Dong et al., 2019) is an end-to-end differentiable neural-symbolic architecture for inductive learning and logic reasoning. Predicates are represented by probabilistic tensors, i.e., grounded on any possible combination of objects. From a set of premises (base predicates), the forward pass in NLM, mimicking a sequence of forward chaining steps, outputs some conclusive tensors. Some logical architectural bias is embedded, as through the explicit wiring among the neural modules to realize the logical existential quantifiers as tensorial operations. Such approach can be seen as learning on a continuous relaxation of logic programs. Some undeniable advantages of NLM compared with the neuro-symbolic literature is the improved inference time, and that it does not rely on hand-engineered rule templates. However, what it gains in scalability, it loses in interpretability.

Attention-based Inductive Bias

Another intelligibility incentive is the use of selective attention mechanisms for the policy network (Mott et al., 2019; Tang et al., 2020), or the Q-network (Annasamy & Sycara, 2019). The work of Tang et al. (2020) evolves RL agents which are encouraged to attend to a small fraction of its visual input, by selecting which spatial patches of the input representation they feed to the LSTM controller. Similarly, Mott et al. (2019) present a soft, top-down, spatial attention mechanism applied to the visual input, while allegedly uncovering part of the underlying decision process, in terms of space (“where”) and content (“what”). Although the authors argue that these attentions mechanisms yield more informative and reliable explanations than other methods for analyzing saliency, the correlation between attention and explainability has been both supported (Wiegrefe & Pinter, 2019) and

Table 8 Overview of approaches with *Intelligibility-Driven Regularization*

Regularizer	Model	RL Training	Interpretability	References
Smoothness regularizer	NN	Yes	Weak	Jia et al. (2019)
Alignment regularizer	Model-based model-free, w/ double DQN	Yes	Weak	Francois-Lavet et al. (2019)
Model compression	NN	No	Weak	Buciluă et al. (2006)
L1-Regularizer	NN	No	Weak	Zhang et al. (2016)
Legibility/Predictability Regularizer	–	No	Weak	Dragan et al. (2013)
Tree-Regularizer	NN	No	Weak	Wu et al. (2019b)

disputed (Brunner et al., 2020; Jain & Wallace, 2019) in further work along different scenarios. For related work focusing on explainability (see Sect. 7).

Discussion

As recent work suggests, relational or logical inductive bias can foster reasoning and generalization over structured data, may it be a graph or predicates, and can improve learning efficiency and robustness, while still benefiting from the flexibility of statistical learning, in contrast to pure symbolic approaches. Although, as soon as the environment and dynamics are complex, these learned relational, logical or attentive representations would not be sufficient to non-ambiguously make sense of the decision-making process.

It is worth mentioning some related deep learning work, which has used logical background knowledge as a way to shape the neural architecture itself, such as Franca et al. (2014), whose neural ILP-solver builds recursive NNs, made with AND-OR type of networks. However, strong architectural bias may drastically decrease the model's expressivity.

6.4 Intelligibility-driven regularization

An alternative to structural bias is to encompass a soft bias on the hypothesis space, through some additional cost function favoring a certain notion of interpretability. This additional term, as ultimately aiming at improving generalization error over training error, can be interpreted as a regularization technique. Although this approach is natural and has received some attention in the broader scope of machine learning, it is relatively less explored in DRL. For this reason, we also discuss some non-RL studies in that direction, which could potentially be fruitfully adapted to RL. These approaches are summarized in Table 8.

Classical regularization methods in deep learning which foster lower complexity should be beneficial for interpretability, although far from being sufficient, e.g., L1-regularization (Zhang et al., 2016) encouraging sparsity or *model compression* (Buciluă et al., 2006).

Other interpretability-oriented penalty formulations have been proposed such as erratic-behavior penalties to improve smoothness (Jia et al., 2019), or objectives targeting legible or predictable motions (Dragan et al., 2013); another example is given by Francois-Lavet et al. (2019) who introduce an additional loss term based on cosine similarity to encourage the predicted abstract state change to align with a chosen embedding vector. This

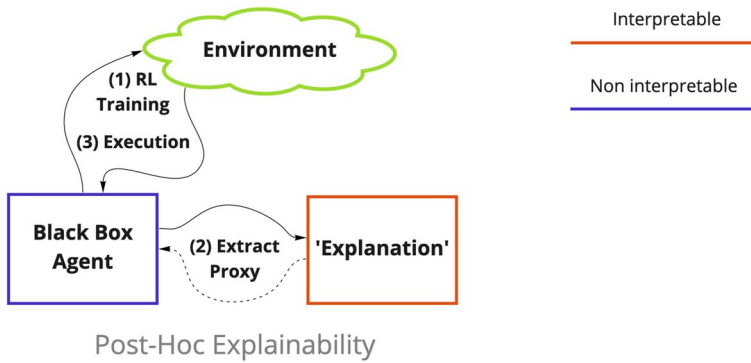


Fig. 4 Illustration for explainable RL approaches. Dashed lines represent optional links, depending on the methods

regularization arguably drives the abstract state to be more meaningful and generalizable, and thereupon may enable more efficient planning.

More typical interpretability-oriented regularizers have been proposed, with FOL (in DL, with Serafini & d'Avila Garcez, 2016), or tree-regularizers (Wu et al., 2019b). The (regional) tree-regularization proposed by Wu et al. (2019b) aims to specifically learn deep policy networks whose decision boundaries are well approximated by small decision tree(s), hence targets human simulatability. By considering interpretability from the very start—in contrast to indirect approaches aiming at approximating a black-box policy network with a decision tree a posteriori—it should be more accessible to reach both good performance and simulatability, due to the *multiple optima* property of deep NN. Indeed, indirect approaches may be unreliable as the original unregularized black box NN has no incentive to be simulatable or decomposable.

Discussion

Embedding interpretability bias through regularizers has the advantage to be easily integrated with any optimization algorithm, such as stochastic gradient descent, if the hypothesis class is made differentiable. As deep models have—infamously—multiple optima of similar predictive accuracy (Goodfellow et al., 2016), we can hope that using interpretability-oriented regularizers may not impact much the performance. However, since such approaches do not restrict the search space per se, they do not provide interpretability guarantee.

There are a few noticeable studies in deep learning aiming to distil logical knowledge through loss functions and regularizers during the NN training, such as (Demeester et al., 2016; Donadello et al., 2017; Diligenti et al., 2017; Rocktäschel et al., 2015; Serafini & d'Avila Garcez, 2016; Wang & Pan, 2019; Xu et al., 2018)¹⁰ or (Minervini et al., 2017) with adversarial training. Bridging the gap between XRL and interpretable literature, Plumb et al. (2020) propose some explainability-regularizers, differentiable, and model agnostic, which would encourage the learned models, trained end-to-end, to be well explainable. Although intelligibility-enhancers seem numerous, the question of defining

¹⁰ For instance, Serafini and d'Avila Garcez (2016) use FOL-based loss-function to constrain the learned semantic representations to be logically consistent.

Table 9 Overview of approaches in *Explainable RL*

Type	Approach	References
Visual	t-SNE, saliency map from Jacobian	Zahavy et al. (2016)
Visual	Saliency map from perturbation	Greydanus et al. (2018)
Visual	Saliency map by balancing specificity and relevance	Gupta et al. (2020)
Visual	SHAP	Wang et al. (2020)
Visual	Attention mask	Shi et al. (2020)
Visual	Attention mask with information bottleneck	Kim and Bansal (2020)
Visual	Summary from history	Sequeira and Gervasio (2020)
Textual	State predicates	Hayes and Shah (2017)
Textual	State and outcome predicates	van der Waa et al. (2018)
Textual	Reuse of provided instructions	Fukuchi et al. (2017)
Causal	Causal model	Madumal et al. (2020b)
Causal	Opportunity chain	Madumal et al. (2020a)
Policy	Soft decision tree	Coppens et al. (2019)
Policy	Decision tree	Bewley and Lawry (2021)
Policy	Decision tree	Kenny et al. (2023)
Other	Reward decomposition	Juozapaitis et al. (2019)
Other	Markov chain on abstract state space	Topin and Veloso (2019)
Other	Probability of success, # steps to reach goal	Cruz et al. (2019)

specific regularizers leading to a reasonably-interpretable decision-making in complex environments is far from being obvious.

7 Explainable RL

Although the focus of this survey is on interpretable RL, we also provide a succinct overview of explainable RL (XRL) for completeness and in order to contrast it with the work in interpretable RL. Figure 4 illustrates the high-level procedure in XRL, which can be contrasted with the approaches for interpretable decision-making described in Fig. 3. A more thorough discussion on XRL can be found in the recent surveys by Alharin et al. (2020) or Heuillet et al. (2021). The methods discussed in this section are summarized in Table 9.

The goal in XRL is to provide some explanations regarding an RL agent's decisions, e.g., highlighting the main features that influenced a decision and their importance. This is commonly done via a post-hoc and often model-agnostic procedure after a black-box model is already trained, which usually only aims to offer a functional understanding. Many contextual parameters should be taken into consideration when defining what constitutes a "good" explanation for a scenario, e.g., background knowledge and levels of expertise of the explanation recipients, their needs and expectations, but also (often neglected) the time available to them. Explanations can take various forms:

Visual explanation

Using the DQN algorithm, Zahavy et al. (2016) build two graphical representations in order to analyze the decisions made by the DQN network: (1) t-SNE maps (van der Maaten & Hinton, 2008) from the activations of the last hidden layer of the network and

(2) saliency maps from the Jacobian of the network. Motivated by the limitations of Jacobian saliency maps, Greydanus et al. (2018) propose to build saliency maps using a perturbation-based approach, which provides information about the importance of a perturbed region. Continuing this line of research, Gupta et al. (2020) introduce the idea of balancing specificity and relevance in order to build saliency maps to highlight more relevant regions. In order to take into account non-visual inputs as well, Wang et al. (2020) extend a generic explanation technique called SHAP (SHapley Additive exPlanation) (Lundberg & Lee, 2017) to select important features for RL. Another approach is based on attention mechanisms. Shi et al. (2020) propose to learn attention masks in a self-supervised way to highlight information important for a decision. In Kim and Bansal (2020), attention is further combined with an information bottleneck mechanism in order to generate sparser attention maps. Using a different kind of explanation, Sequeira and Gervasio (2020) investigate the use of visual summaries extracted from histories to explain an agent's behavior.

Textual explanation

The work of Hayes and Shah (2017) generates explanations for choosing an action by finding state predicates that co-occur with that action. Inspired by that approach, van der Waa et al. (2018) extend it by introducing outcome predicates and provide contrastive explanations using both state and outcome predicates. In a setting where the agent learns from instructions given by a human tutor, Fukuchi et al. (2017) propose to explain the agent's decisions by reusing the provided instructions.

Causal explanation

In the proposition of Madumal et al. (2020b), a causal model is learned from a given graph of causal relation in order to generate contrastive explanations of action choices. Building on this work, Madumal et al. (2020a) instead generate explanations based on potential future actions using the concept of opportunity chains, which include information of what is enabled or caused by an action.

Interpretable policy

Some work tries to obtain a more intelligible policy in order to explain a trained RL agent using, e.g., soft decision trees (Coppens et al., 2019), decision trees (Bewley & Lawry, 2021) or prototypes (Kenny et al., 2023). Note that the indirect approach for interpretability (as presented in Sect. 6.2) should not be confused with this approach for post-hoc explainability. In the latter case, a more intelligible policy is learned to explain a black-box policy that is used as the proper controller. In contrast, in the former case, a more interpretable policy is learned to be used as the final controller that replaces the intermediate black-box policy, which therefore does not need to be explained anymore. Therefore, in the latter case, it is important that the intelligible policy mimics the black-box policy well, while in the former, the performance of the interpretable policy is more important than its ability to mimic the black-box policy. When learning such an explanatory policy, a compromise between its intelligibility and the fidelity of its approximation needs to be found. One common drawback is that such an approximation may be valid only on a restricted domain.

Other

Furthermore, Juozapaitis et al. (2019) propose to learn a vector Q-function, where each component corresponds to a given attribute called reward type. This decomposition of the Q-function is then used to explain preferences between actions. In contrast, in Topin and Veloso (2019), a policy is explained with a Markov chain built on an abstract state space. In addition, in goal-oriented RL, Cruz et al. (2019) justify an action choice based on its probability of success and the number of time steps to reach the goal.

Discussion

Most work we discussed takes the target audience of the explanations to be the end-user. Even in this case, explanations can take multiple forms. Thus, it can be presented to the user in different modes (e.g., visual, textual, tabular,...) and it can be either local or global. Beyond their forms, explanations may also answer intelligibility queries of different nature and granularity: certainty, contextual, case-based or analogies, contrastive, counterfactual (“what if”), simulation-based (consequences), trace/steps, why not, etc (e.g., Chari et al. (2020), Lim et al. (2019), Mittelstadt et al. (2019)). Hence, an explanation can be used to clarify, justify, or rationalize an action choice. Future work on XRL should make those aspects clear, since this information would impact how an explanation technique should be evaluated and taken into consideration.

One issue with post-hoc explanation approaches is that while the generated explanation may seem to make sense, it may in fact be specious (e.g., Atrey et al., 2020 for saliency maps) and may not reflect the true inner working of the model. While this may not be an issue if the explanation is used as a tool to justify an action choice to a user, this is problematic for understanding the decision-making process. Note this issue does not occur if an interpretable policy is used for decision-making.

While the explainable and interpretable literature refers to usually divergent approaches, some recent work aimed at bridging this gap, (e.g., Plumb et al., 2020 previously mentioned in Sect. 6.4). Through regularizers, it gracefully integrates explainability considerations during the training of the model. It stands at odds with traditional XRL literature, assuming they could extract a posteriori explanations, without any incentive for the model to be intelligible.

8 Open problems and research directions

Before concluding this survey, we discuss a selection of open problems, which we regard as essential within the quest for interpretable RL.

Full Interpretability in RL

The work we have reviewed falls in various intermediate levels on the interpretability scale, some being more interpretable than others for different RL components. Moreover, few DRL approaches accepting high-dimensional inputs, if any, can achieve full interpretability, i.e., interpretable inputs, interpretable models, and interpretable decision-making. Designing a fully-interpretable RL method with a high-degree of interpretability seems not to be achievable with the current methods, especially for complex tasks like autonomous driving, although such tasks calls for such methods. Thus, for DRL to be considered as a practical method to solve those difficult tasks, fully interpretable RL methods must be developed for all the RL components. Given the complexity of those tasks, this may only be achievable by abstraction and composition in the programming language sense, where interpretable methods can be composed to solve more difficult problems.

Interpretability vs Performance

A commonly-held opinion is that using a more transparent model or approach impacts negatively the final performance (Ribeiro et al., 2016a). In the light of impressive results achieved by deep learning methods, this opinion seems hard to be challenged. However, some different voices (Rudin & Carlson, 2019) suggest that black-box models like those based on deep learning may not always be needed and that in some domains simple models should be favored and can obtain excellent performance without the drawbacks of deep learning methods. Similar remarks have also been made in

deep RL by Mania et al. (2018), who showed that simple linear models with stochastic search can fare well against more advanced DRL methods. Extrapolating those observations, one may wonder if this could be achieved with all aspects of RL (inputs, models, decision-making) and if interpretability can be considered as a regularization technique, which would bring more transparency obviously, but also larger generalizability.

Interpretability vs Scalability

In addition to the challenge of designing a fully-interpretable RL method, running such a method in order to learn a fully-interpretable solution would probably be also more costly in terms of computation than a standard DRL algorithm. Indeed, for decision-making for instance, learning a fully-interpretable policy corresponds to a task similar to program synthesis (Gulwani et al., 2017), which requires a search over a discrete solution space whose size increases exponentially fast with the solution size. Therefore, there may be a trade-off between the degree of interpretability one may want to achieve and the scalability of the interpretable algorithm. This question is crucial to investigate as the research moves to more and more interpretable methods, critically needed for high-stake tasks.

Evaluation of Interpretability and Explainability

We finish this discussion by a more classic question that has been frequently raised within XAI, but that we mention here due to its importance. Given the various meanings of interpretability and explainability and more precisely the various purposes they can serve, there is no common ground for the definition of good evaluation metrics for XAI in general, but also for interpretable and explainable RL. For interpretability, is there a good metric for deciding if one model is more interpretable than another? For explainability, is it possible to evaluate what constitutes a good explanation in a specific context? This state of affairs prevents a comparative evaluation of the different methods that have been proposed, which also impedes the rapid progress in this research direction. While achieving more precise definitions for interpretability and explainability can help, evaluation metrics and protocols could be proposed depending on precise goals regarding ethical, legal, operational, or usability concerns, which may help them to be adopted by the research community.

Impact of Foundation Models on Interpretable and Explainable RL

Recently, foundation models (i.e., large models pretrained on a massive quantity of generally-unlabelled data) (Bommasani et al., 2022) have shown impressive capabilities, notably for generation under user prompt of texts (OpenAI et al., 2023; Glaese et al., 2022) or images (Ramesh et al., 2021; Rombach et al., 2022). Although such models may arguably not be very transparent, various research effort (e.g., (Friedman et al., 2023; Singh et al., 2023; Wu et al., 2023)) is currently spent to make their use more interpretable. Moreover, foundation models, in particular large language models (Liu et al., 2023), which are pretrained on text, have potentially the capability of generating textual explanation, although one should be careful about various well-known risks such as bias or hallucinations. Interestingly, they can also implicitly learn human preferences during pretraining and can therefore serve directly as a source of rewards (Kwon et al., 2023; Rafailov et al., 2023).

In addition, these large models can be fine-tuned using RL with human feedback (RLHF) (Casper et al., 2023). RLHF is a framework that allows humans to specify their preferences, usually via pairwise comparisons of trajectories, circumventing the need of specifying a reward function, which is usually a difficult task for the system designer. In RLHF, a reward function is usually learned from human preferences. In such approach, the preference input is naturally understandable to humans and learning a reward function with a model such as a tree-structured model (Bewley & Lecue, 2022) can make the whole preference model interpretable.

9 Conclusion

We surveyed recent work in RL related to the important concern of interpretability (and its related notion of explainability). We argued for a definition of interpretability in RL, which contrary to the general setting of explainable AI, leads to different levels of transparency in the components that play a role in RL. In particular, we first discussed studies that focus on interpretable inputs (e.g., observations, but possibly other structural information). Moreover, we provided an overview of approaches that deal with learning an interpretable transition model (e.g., important for interpretable model-based RL), but also those that deal with learning an interpretable preference model (e.g., fundamental to justify action selection). Then, we surveyed methods learning interpretable policies, which constitute arguably the most critical part of interpretable RL. For completeness, we also provided a short review of work related to post-hoc explainability. Finally, we highlighted a few open problems and future research directions that we deemed as particularly relevant.

Although concerns around the ethical implications of algorithmic and automation deployment are nothing new (Wiener, 1954), the field of AI ethics still seems at its infancy as we begin to witness the extent of the influence and impact that these systems may have on our societal fabric when deployed. In this regard, a responsible practice for the design, implementation, use, and monitoring/auditing of AI-driven systems is greatly impeded by the non-intelligibility of current algorithms. As RL-based systems become more widespread, questions related to interpretability become consequently increasingly pressing. One could even contend that interpretable RL is one of the key deadlocks to overcome to make RL a more functional method for being deployed in real-life. While it may be hard to achieve a fully intelligible RL model, one may envision hierarchical RL approaches where some parts may not be completely transparent—e.g., at the low-level—but a maximum of other parts—e.g., at the subgoal level—are thoroughly interpretable.

While the act of opening up the blackbox do not suffice to instantly disclose a thorough understanding of its social implications—since we “*need to look across the system, rather than merely inside*” (as noted by Ananny & Crawford, 2018)—algorithmic intelligibility appears as a promising step towards further *algorithmic accountability* and more trustworthy AI. We encourage the curious reader to look further at the generous work of other researchers investigating these tangent questions (such as Crawford et al. (2016), Raji et al. (2020), Daly et al. (2019), Yu et al. (2018), Commission (2019) to mention only a few).

Author contributions All the authors participated in the initial discussions to define the scope of the survey and find the relevant papers. Once defined, the first three authors wrote the major part of the survey. DL and TY helped with improving earlier versions of this manuscript.

Funding This research work was funded by Huawei Technology Ltd.

Declarations

Conflict of interest The authors declare that they have no conflict of interest.

References

- Adjudah, D., Klinger, T., & Joseph, J. (2018). Symbolic relation networks for reinforcement learning. In *NeurIPS workshop on representation learning*.
- Agnew, W., & Domingos, P. (2018). Unsupervised object-level deep reinforcement learning. In *NeurIPS workshop on deep RL*.
- Akrou, R., Tateo, D., & Peters, J. (2019). Towards reinforcement learning of human readable policies. In *Workshop on deep continuous-discrete machine learning*.
- Aksaray, D., Jones, A., Kong, Z., et al. (2016). Q-Learning for robust satisfaction of signal temporal logic specifications. In *CDC*.
- Alharin, A., Doan, T. N., & Sartipi, M. (2020). Reinforcement learning interpretation methods: A survey. *IEEE Access*, 8, 171058–171077.
- Alshiekh, M., Bloem, R., Ehlers, R., et al. (2018). Safe reinforcement learning via shielding. In *AAAI*.
- Amodei, D., Olah, C., Steinhardt, J., et al. (2016). Concrete Problems in AI Safety. [arXiv: 1606.06565](https://arxiv.org/abs/1606.06565)
- Ananny, M., & Crawford, K. (2018). Seeing without knowing: Limitations of the transparency ideal and its application to algorithmic accountability. *New Media and Society*, 20(3), 973–89.
- Andersen, G., & Konidaris, G. (2017). Active exploration for learning symbolic representations. In *NeurIPS*.
- Anderson, G., Verma, A., Dillig, I., et al. (2020). Neurosymbolic reinforcement learning with formally verified exploration. In *NeurIPS*.
- Andreas, J., Klein, D., & Levine, S. (2017). Modular multitask reinforcement learning with policy sketches. In *ICML*.
- Annasamy, R.M., & Sycara, K. (2019). Towards better interpretability in deep Q-networks. In *AAAI*.
- Arnold, T., Kasenberg, D., & Scheutz, M. (2017). Value alignment or misalignment: What will keep systems accountable? In *AAAI workshop*.
- Arora, S., & Doshi, P. (2018). A survey of inverse reinforcement learning: Challenges, methods and progress. [arXiv:1806.06877](https://arxiv.org/abs/1806.06877)
- Atrey, A., Clary, K., & Jensen, D. (2020). Exploratory not explanatory: Counterfactual analysis of saliency maps for deep reinforcement learning. In *ICLR*.
- Ault, J., Hanna, J. P., & Sharon, G. (2020). Learning an interpretable traffic signal control policy. In *AAMAS*.
- Bader, S., & Hitzler, P. (2005). Dimensions of neural-symbolic integration: A structured survey. In *We Will Show Them: Essays in Honour of Dov Gabbay*.
- Barredo Arrieta, A., Díaz-Rodríguez, N., Ser, J. D., et al. (2020). Explainable artificial intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58, 82–115.
- Barto, A. G., & Mahadevan, S. (2003). Recent advances in hierarchical reinforcement learning. *Discrete Event Dynamic Systems*
- Barwise, J. (1977). An introduction to first-order logic. *Studies in Logic and the Foundations of Mathematics*, 90, 5–46.
- Bastani, O., Pu, Y., & Solar-Lezama, A. (2018). Verifiable reinforcement learning via policy extraction. In *NeurIPS*.
- Battaglia, P., Pascanu, R., Lai, M., et al. (2016). Interaction networks for learning about objects, relations and physics. In *NeurIPS*.
- Battaglia, P. W., Hamrick, J. B., Bapst, V., et al. (2018). Relational inductive biases, deep learning, and graph networks. [arXiv:1806.01261](https://arxiv.org/abs/1806.01261)
- Bear, D., Fan, C., Mrowca, D., et al. (2020). Learning physical graph representations from visual scenes. In *NeurIPS*.
- Bertsekas, D., & Tsitsiklis, J. (1996). *Neuro-dynamic programming*. Athena Scientific.
- Bewley, T., & Lawry, J. (2021). TripleTree: A versatile interpretable representation of black box agents and their environments. In *AAAI*.
- Bewley, T., & Lécué, F. (2022). Interpretable preference-based reinforcement learning with tree-structured reward functions. In *AAMAS*.
- Beyret, B., Shafti, A., & Faisal, A. A. (2019). Dot-to-dot: Explainable hierarchical reinforcement learning for robotic manipulation. In *IROS*.
- Bommasani, R., Hudson, D. A., Adeli, E., et al. (2022). On the opportunities and risks of foundation models. [arXiv:2108.07258](https://arxiv.org/abs/2108.07258)
- Bonnefon, J., Shariff, A., & Rahwan, I. (2019). The trolley, the bull bar, and why engineers should care about the ethics of autonomous cars [point of view]. *Proceedings of the IEEE*, 107(3), 502–4.
- Boutillier, C., Dearden, R., & Goldszmidt, M. (2000). Stochastic dynamic programming with factored representations. *Artificial Intelligence*, 121(1–2), 49–107.

- Brunelli, R. (2009). *Template matching techniques in computer vision: Theory and practice*. Wiley Publishing.
- Brunner, G., Liu, Y., Pascual, D., et al. (2020). On identifiability in transformers. In *ICLR*
- Buciluă, C., Caruana, R., & Niculescu-Mizil, A. (2006). Model compression. In *KDD*.
- Burke, M., Penkov, S., & Ramamoorthy, S. (2019). From explanation to synthesis: Compositional program induction for learning from demonstration. In *RSS*.
- Camacho, A., Toro Icarte, R., Klassen, T. Q., et al. (2019). LTL and beyond: Formal languages for reward function specification in reinforcement learning. In *IJCAI*.
- Cao, Y., Li, Z., Yang, T., et al. (2022). GALOIS: Boosting deep reinforcement learning via generalizable logic synthesis. In *NeurIPS*.
- Casper, S., Davies, X., Shi, C., et al. (2023). Open problems and fundamental limitations of reinforcement learning from human feedback. [arXiv:2307.15217](https://arxiv.org/abs/2307.15217)
- Chang, M. B., Ullman, T., Torralba, A., et al. (2017). A compositional object-based approach to learning physical dynamics. In *ICLR*.
- Chari, S., Gruen, D. M., Seneviratne, O., et al. (2020). Directions for explainable knowledge-enabled systems. [arXiv:2003.07523](https://arxiv.org/abs/2003.07523)
- Chen, J., Li, S. E., & Tomizuka, M. (2020). Interpretable end-to-end urban autonomous driving with latent deep reinforcement learning. In *ICML workshop on AI for autonomous driving*.
- Cichosz, P., & Pawelczak, L. (2014). Imitation learning of car driving skills with decision trees and random forests. *International Journal of Applied Mathematics and Computer Science*, 24, 579–97.
- Cimatti, A., Pistore, M., & Traverso, P. (2008). Automated planning. In *Handbook of knowledge representation*.
- Cole, J., Lloyd, J., & Ng, K. S. (2003). Symbolic learning for adaptive agents. In *Annual partner conference*.
- Commission, E. (2019). Ethics guidelines for trustworthy AI. <https://ec.europa.eu/digital-single-market/en/news/ethics-guidelines-trustworthy-ai>
- Coppens, Y., Efthymiadis, K., Lenaerts, T., et al. (2019). Distilling deep reinforcement learning policies in soft decision trees. In *IJCAI workshop on XAI*.
- Corazza, J., Gavran, I., & Neider, D. (2022). Reinforcement learning with stochastic reward machines. In *AAAI*.
- Cranmer, M., Sanchez Gonzalez, A., Battaglia, P., et al. (2020). Discovering symbolic models from deep learning with inductive biases. In *NeurIPS*.
- Crawford, K., Dobbe, R., Dryer, T., et al. (2016). *AI Now Report*. AI Now Institute: Tech. rep.
- Cropper, A., Dumančić, S., & Muggleton, S.H. (2020). Turning 30: New ideas in inductive logic programming. In *IJCAI*.
- Cruz, F., Dazeley, R., & Vamplew, P. (2019). Memory-based explainable reinforcement learning. In *Advances in artificial intelligence*.
- Daly, A., Hagendorff, T., Li, H., et al. (2019). *Artificial Intelligence, Governance and Ethics: Global Perspectives*. SSRN Scholarly Paper: Chinese University of Hong Kong.
- d'Avila Garcez, A., Dutra, A. R. R., & Alonso, E. (2018). Towards Symbolic Reinforcement Learning with Common Sense. [arXiv:1804.08597](https://arxiv.org/abs/1804.08597)
- De Raedt, L., & Kimmig, A. (2015). Probabilistic (logic) programming concepts. *Machine Learning*, 100(1), 5–47.
- Dean, T., & Kanazawa, K. (1990). A model for reasoning about persistence and causation. *Computational Intelligence*, 5(3), 142–150.
- Degrís, T., Sigaud, O., & Willemin, P. H. (2006). Learning the structure of factored Markov decision processes in reinforcement learning problems. In *ICML*.
- Delfosse, Q., Shindo, H., Dhami, D., et al. (2023). Interpretable and explainable logical policies via neurally guided symbolic abstraction. In *NeurIPS*.
- Demeester, T., Rocktäschel, T., & Riedel, S. (2016). Lifted rule injection for relation embeddings. In *EMNLP*.
- Diligenti, M., Gori, M., & Saccà, C. (2017). Semantic-based regularization for learning and inference. *Artificial Intelligence*, 244, 143–65.
- Diuk, C., Cohen, A., & Littman, M. L. (2008). An object-oriented representation for efficient reinforcement learning. In *ICML*.
- Donadello, I., Serafini, L., & D'Avila Garcez, A. (2017). Logic tensor networks for semantic image interpretation. In *IJCAI*.
- Dong, H., Mao, J., Lin, T., et al. (2019). Neural logic machines. In *ICLR*.
- Doshi-Velez, F., Kortz, M., Budish, R., et al. (2019). Accountability of AI under the law: The role of explanation. [arXiv:1711.01134](https://arxiv.org/abs/1711.01134)
- Dragan, A. D., Lee, K. C., & Srinivasa, S. S. (2013). Legibility and predictability of robot motion. In *HRI*.

- Driessens, & Blockeel, H. (2001). Learning digger using hierarchical reinforcement learning for concurrent goals. In *EWRL*.
- Driessens, K., Ramon, J., & Gartner, T. (2006). Graph kernels and Gaussian processes for relational reinforcement learning. *Machine Learning*
- Dutra, A. R., & d'Avila Garcez, A. S. (2017). A Comparison between deep Q-networks and deep symbolic reinforcement learning. In *CEUR workshop proceedings*.
- Dwork, C., Hardt, M., Pitassi, T., et al. (2012). Fairness through awareness. In *ICTS*.
- Dzeroski, S., Raedt, L. D., & Blockeel, H. (1998). Relational reinforcement learning. In *ICML*.
- Džeroski, S., De Raedt, L., & Driessens, K. (2001). Relational reinforcement learning. *Machine Learning*, 43(1), 7–52.
- Ernst, D., Geurts, P., & Wehenkel, L. (2005). Tree-based batch mode reinforcement learning. *JMLR*, 6, 503–556.
- Evans, R., & Grefenstette, E. (2018). Learning explanatory rules from noisy data. *Journal of Artificial Intelligence Research*, 61, 1–64.
- Eysenbach, B., Salakhutdinov, R. R., & Levine, S. (2019). Search on the replay buffer: Bridging planning and reinforcement learning. In *NeurIPS*.
- Finn, C., Goodfellow, I., & Levine, S. (2016). Unsupervised learning for physical interaction through video prediction. In *NeurIPS*.
- Finn, C., & Levine, S. (2017). Deep visual foresight for planning robot motion. In *ICRA*.
- Franca, M. V. M., Zaverucha, G., & Garcez, A. (2014). Fast relational learning using bottom clause propositionalization with artificial neural networks. *Machine Learning*, 94(1), 81–104.
- Francois-Lavet, V., Bengio, Y., Precup, D., et al. (2019). Combined reinforcement learning via abstract representations. In *AAAI*.
- Friedler, S. A., Scheidegger, C., & Venkatasubramanian, S. (2021). The (Im)possibility of fairness: Different value systems require different mechanisms for fair decision making. *Communications of the ACM*, 64(4), 136–143.
- Friedman, D., Wettig, A., & Chen, D. (2023). Learning transformer programs. In *NeurIPS*.
- Fujimoto, S., Hoof, H., & Meger, D. (2018). Addressing function approximation error in actor-critic methods. In *ICML*.
- Fukuchi, Y., Osawa, M., Yamakawa, H., et al. (2017). Autonomous self-explanation of behavior for interactive reinforcement learning agents. In *International conference on human agent interaction*.
- Furelos-Blanco, D., Law, M., Jonsson, A., et al. (2021). Induction and exploitation of subgoal automata for reinforcement learning. *JAIR*, 70, 1031–1116.
- Gaon, M., & Brafman, R. I. (2020). Reinforcement learning with non-Markovian rewards. In *AAAI*.
- Garg, S., Bajpai, A., Mausam. (2020). Symbolic network: Generalized neural policies for relational MDPs. [arXiv:2002.07375](https://arxiv.org/abs/2002.07375)
- Garnelo, M., Arulkumaran, K., & Shanahan, M. (2016). Towards deep symbolic reinforcement learning. In *NeurIPS workshop on DRL*.
- Gilmer, J., Schoenholz, S. S., Riley, P. F., et al. (2017). Neural message passing for quantum chemistry. In *ICML*.
- Gilpin, L. H., Bau, D., Yuan, B. Z., et al. (2019). Explaining explanations: An overview of interpretability of machine learning. In *DSAA*.
- Glaese, A., McAleese, N., Trebacz, M., et al. (2022). Improving alignment of dialogue agents via targeted human judgements. [arXiv:2209.14375](https://arxiv.org/abs/2209.14375)
- Glanois, C., Jiang, Z., Feng, X., et al. (2022). Neuro-symbolic hierarchical rule induction. In *ICML*.
- Goel, V., Weng, J., & Poupart, P. (2018). Unsupervised video object segmentation for deep reinforcement learning. In *NeurIPS*.
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT Press.
- Greydanus, S., Koul, A., Dodge, J., et al. (2018). Visualizing and understanding atari agents. In *ICML*.
- Grzes, M., & Kudenko, D. (2008). Plan-based reward shaping for reinforcement learning. In *International conference intelligent systems*.
- Guestrin, C., Koller, D., Gearhart, C., et al. (2003). Generalizing plans to new environments in relational MDPs. In *IJCAI*.
- Gulwani, S., Polozov, O., & Singh, R. (2017). Program synthesis. *Foundations and Trends in Programming Languages*, 4(1–2), 1–119.
- Gupta, P., Puri, N., Verma, S., et al. (2020). Explain your move: Understanding agent actions using focused feature saliency. In *ICLR*.
- Gupta, U. D., Talvitie, E., & Bowling, M. (2015). Policy tree: Adaptive representation for policy gradient. In *AAAI*.

- Haarnoja, T., Zhou, A., Abbeel, P., et al. (2018). Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *ICML*.
- Harnad, S. (1990). The symbol grounding problem. *Physica D-Nonlinear Phenomena*, 42, 335–346.
- Hasanbeig, M., Kroening, D., & Abate, A. (2020). Deep reinforcement learning with temporal logics. In *Formal modeling and analysis of timed systems*.
- Hayes, B., & Shah, J. A. (2017). Improving robot controller transparency through autonomous policy explanation. In *International conference on HRI*.
- Hein, D., Hentschel, A., Runkler, T., et al. (2017). Particle swarm optimization for generating interpretable fuzzy reinforcement learning policies. *Engineering Applications of AI*, 65, 87–98.
- Hein, D., Udluft, S., & Runkler, T. A. (2018). Interpretable policies for reinforcement learning by genetic programming. *Engineering Applications of AI*, 76, 158–169.
- Hein, D., Udluft, S., & Runkler, T. A. (2019). Generating interpretable reinforcement learning policies using genetic programming. In *GECCO*.
- Henderson, P., Islam, R., Bachman, P., et al. (2018). Deep reinforcement learning that matters. In *AAAI*.
- Hengst, B. (2010). Hierarchical reinforcement learning. *Encyclopedia of machine learning* (pp. 495–502). Springer.
- Heuillet, A., Couthouis, F., & Díaz-Rodríguez, N. (2021). Explainability in deep reinforcement learning. *Knowledge-Based Systems*, 214, 106685.
- Higgins, I., Amos, D., Pfau, D., et al. (2018). Towards a definition of disentangled representations. [arXiv:1812.02230](https://arxiv.org/abs/1812.02230)
- Horvitz, E., & Mulligan, D. (2015). Data, privacy, and the greater good. *Science*, 349(6245), 253–255.
- Huang, S., Papernot, N., Goodfellow, I., et al. (2017). Adversarial attacks on neural network policies. In *ICLR workshop*.
- Hussein, A., Gaber, M. M., Elyan, E., et al. (2017). Imitation learning: A survey of learning methods. *ACM Computing Surveys*, 50(2), 211–2135.
- Illanes, L., Yan, X., Icarte, R. T., et al. (2020). Symbolic plans as high-level instructions for reinforcement learning. In *ICAPS*.
- Iyer, R., Li, Y., Li, H., et al. (2018). Transparency and explanation in deep reinforcement learning neural networks. In *AIES*.
- Jain, S., & Wallace, B. C. (2019). Attention is not explanation. In *NAACL*.
- Janisch, J., Pevný, T., & Lisý, V. (2021). Symbolic relational deep reinforcement learning based on graph neural networks. [arXiv:2009.12462](https://arxiv.org/abs/2009.12462)
- Jia, R., Jin, M., Sun, K., et al. (2019). Advanced building control via deep reinforcement learning. In *Energy Procedia*.
- Jiang, Y., Yang, F., Zhang, S., et al. (2018). Integrating task-motion planning with reinforcement learning for robust decision making in mobile robots. In *ICAPS*.
- Jiang, Z., & Luo, S. (2019). Neural logic reinforcement learning. In *ICML*.
- Jin, M., Ma, Z., Jin, K., et al. (2022). Creativity of ai: Automatic symbolic option discovery for facilitating deep reinforcement learning. In *AAAI*.
- Juozapaitis, Z., Koul, A., Fern, A., et al. (2019). Explainable reinforcement learning via reward decomposition. In *IJCAI/ECAI workshop on explainable artificial intelligence*.
- Kaiser, M., Otte, C., Runkler, T., et al. (2019). Interpretable dynamics models for data-efficient reinforcement learning. In *ESANN*.
- Kansky, K., Silver, T., Mély, D. A., et al. (2017). Schema networks: Zero-shot transfer with a generative causal model of intuitive physics. In *ICML*.
- Kasenberg, D., & Scheutz, M. (2017). Interpretable apprenticeship learning with temporal logic specifications. In *CDC*.
- Kenny, E. M., Tucker, M., Shah, J. (2023). Towards interpretable deep reinforcement learning with human-friendly prototypes. In *ICLR*.
- Kim, J., & Bansal, M. (2020). Attentional bottleneck: Towards an interpretable deep driving network. In *CVPR workshop*.
- Koller, D. (1999). Probabilistic relational models. In *Inductive logic programming* (pp. 3–13).
- Konidaris, G., Kaelbling, L. P., & Lozano-Perez, T. (2014). Constructing symbolic representations for high-level planning. In *AAAI*.
- Konidaris, G., Kaelbling, L. P., & Lozano-Perez, T. (2015). Symbol acquisition for probabilistic high-level planning. In *IJCAI*.
- Konidaris, G., Kaelbling, L. P., & Lozano-Perez, T. (2018). From skills to symbols: Learning symbolic representations for abstract high-level planning. *JAIR*, 61, 215–289.
- Koul, A., Greydanus, S., & Fern, A. (2019). Learning finite state representations of recurrent policy networks. In *ICLR*.

- Kulick, J., Toussaint, M., & Lang, T. et al (2013). Active learning for teaching a robot grounded relational symbols. In *IJCAI*.
- Kunapuli, G., Odom, P., & Shavlik, J. W. et al (2013). Guiding autonomous agents to better behaviors through human advice. In *ICDM*.
- Kwon, M., Xie, S. M., & Bullard, K. et al (2023). Reward design with language models. In *ICLR*.
- Lao, N., & Cohen, W. W. (2010). Relational retrieval using a combination of path-constrained random walks. In *Machine learning*.
- Leonetti, M., Iocchi, L., & Stone, P. (2016). A synthesis of automated planning and reinforcement learning for efficient, robust decision-making. *Artificial Intelligence*, 241, 103–130.
- Leslie, D. (2020). Understanding artificial intelligence ethics and safety: A guide for the responsible design and implementation of AI systems in the public sector. SSRN Electronic Journal
- Levine, S. (2018). Reinforcement Learning and Control as Probabilistic Inference: Tutorial and Review. [arXiv:1805.00909](https://arxiv.org/abs/1805.00909)
- Li, X., Serlin, Z., Yang, G., et al. (2019). A formal methods approach to interpretable reinforcement learning for robotic planning. *Science Robotics*, 4(37), eaay6276.
- Li, X., Vasile, C. I., & Belta, C. (2017a). Reinforcement learning with temporal logic rewards. In *IROS*.
- Li, Y., Sycara, K., & Iyer, R. (2017b). Object-sensitive deep reinforcement learning. In *Global conference on AI*.
- Li, Y., Tarlow, D., Brockschmidt, M. et al (2017c). Gated graph sequence neural networks. In *ICLR*.
- Likmeta, A., Metelli, A. M., Tirinzoni, A., et al. (2020). Combining reinforcement learning with rule-based controllers for transparent and general decision-making in autonomous driving. *Robotics and Autonomous Systems*, 131, 103568.
- Lim, B. Y., Yang, Q., & Abdul, A. et al (2019). Why these explanations? Selecting intelligibility types for explanation goals. In *IUI workshops*.
- Lipton, Z. C. (2017). The mythos of model interpretability. [arXiv:1606.03490](https://arxiv.org/abs/1606.03490)
- Littman, M. L., Topcu, U., & Fu, J. et al (2017). Environment-independent task specifications via GLTL. [arXiv:1704.04341](https://arxiv.org/abs/1704.04341)
- Liu, G., Schulte, O., & Zhu, W. et al (2018). Toward interpretable deep reinforcement learning with linear model U-trees. In *ECML*.
- Liu, Y., Han, T., Ma, S., et al. (2023). Summary of chatgpt-related research and perspective towards the future of large language models. *Meta-Radiology*, 1(2), 100017.
- Lo Piano, S. (2020). Ethical principles in machine learning and artificial intelligence: Cases from the field and possible ways forward. *Humanities and Social Sciences Communications*, 7(1), 1–7.
- Lu, K., Zhang, S., & Stone, P. et al (2018). Robot representation and reasoning with knowledge from reinforcement learning. [arXiv:1809.11074](https://arxiv.org/abs/1809.11074)
- Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. In *NeurIPS*.
- Lyu, D., Yang, F., & Liu, B. et al (2019). SDRL: Interpretable and data-efficient deep reinforcement learning leveraging symbolic planning. In *AAAI*.
- Ma, Z., Zhuang, Y., & Weng, P. et al (2020). Interpretable reinforcement learning with neural symbolic logic. [arXiv:2103.08228](https://arxiv.org/abs/2103.08228)
- Maclin, R., & Shavlik, J. W. (1996). Creating advice-taking reinforcement learners. *Machine Learning*, 22, 251–282.
- Madumal, P., Miller, T., & Sonenberg, L. et al (2020a). Distal explanations for model-free explainable reinforcement learning. [arXiv:2001.10284](https://arxiv.org/abs/2001.10284)
- Madumal, P., Miller, T., & Sonenberg, L. et al (2020b). Explainable reinforcement learning through a causal lens. In *AAAI*.
- Maes, F., Fonteneau, R., & Wehenkel, L. et al (2012a). Policy search in a space of simple closed-form formulas: towards interpretability of reinforcement learning. In *Discovery science*.
- Maes, F., Wehenkel, L., & Ernst, D. (2012b). Automatic discovery of ranking formulas for playing with multi-armed bandits. In *Recent advances in reinforcement learning*.
- Maes, P., Mataric, M. J., & Meyer, J. A. et al (1996). Learning to use selective attention and short-term memory in sequential tasks. In *International conference on simulation of adaptive behavior*.
- Mania, H., Guy, A., & Recht, B. (2018). Simple random search of static linear policies is competitive for reinforcement learning. In *NeurIPS*.
- Marom, O., & Rosman, B. (2018). Zero-shot transfer with deictic object-oriented representation in reinforcement learning. In *NeurIPS*.
- Martínez, D., Alenyà, G., Torras, C. et al (2016). Learning relational dynamics of stochastic domains for planning. In *ICAPS*.
- Martínez, D., Alenyà, G., Ribeiro, T., et al. (2017). Relational reinforcement learning for planning with exogenous effects. *Journal of Machine Learning Research*, 18(78), 1–44.

- Martínez, D., Alenyà, G., & Torras, C. (2017). Relational reinforcement learning with guided demonstrations. *Artificial Intelligence*, 247, 295–312.
- Mehrabani, N., Morstatter, F., & Saxena, N., et al. (2019). A survey on bias and fairness in machine learning. [arXiv:1908.09635](https://arxiv.org/abs/1908.09635)
- Metzen, J. H. (2013). Learning graph-based representations for continuous reinforcement learning domains. In *ECML*.
- Michels, J., Saxena, A., & Ng, A. Y. (2005). High speed obstacle avoidance using monocular vision and reinforcement learning. In *ICML*.
- Miller, T. (2019). Explanation in artificial intelligence: Insights from the social sciences. *Artificial Intelligence*, 267, 1–38.
- Minervini, P., Demeester, T., & Rocktäschel, T., et al. (2017). Adversarial sets for regularising neural link predictors. In *UAI*.
- Mittelstadt, B., Russell, C., & Wachter, S. (2019). Explaining explanations in AI. In *Conference on fairness, accountability, and transparency*.
- Mnih, V., Kavukcuoglu, K., Silver, D., et al. (2015). Human-level control through deep reinforcement learning. *Nature*, 518, 529–533.
- Mohseni, S., Zarei, N., & Ragan, E. D. (2020). A multidisciplinary survey and framework for design and evaluation of explainable AI systems. [arXiv:1811.11839](https://arxiv.org/abs/1811.11839)
- Molnar, C. (2019). Interpretable machine learning: A guide for making black box models explainable.
- Morley, J., Floridi, L., Kinsey, L., et al. (2020). From what to how: An initial review of publicly available AI ethics tools, methods and research to translate principles into practices. *Science and Engineering Ethics*, 26(4), 2141–68.
- Mott, A., Zoran, D., & Chrzanowski, M., et al. (2019). Towards interpretable reinforcement learning using attention augmented agents. In *NeurIPS*.
- Munzer, T., Piot, B., & Geist, M., et al. (2015). Inverse reinforcement learning in relational domains. In *IJCAI*.
- Nagesh Rao, S., Costa, B., & Filev, D. (2019). Interpretable approximation of a deep reinforcement learning agent as a set of if-then rules. In *ICMLA*.
- Natarajan, S., Joshi, S., & Tadepalli, P., et al. (2011). Imitation learning in relational domains: A functional-gradient boosting approach. In *IJCAI*.
- Ng, A. Y., & Russell, S. (2000). Algorithms for inverse reinforcement learning. In *ICML*.
- OpenAI, Akkaya, I., & Andrychowicz, M., et al. (2019). Solving Rubik's Cube with a Robot Hand. [arXiv:1910.07113](https://arxiv.org/abs/1910.07113)
- OpenAI, & Achiam, J., et al. (2023). Gpt-4 technical report. [arXiv:2303.08774](https://arxiv.org/abs/2303.08774)
- Osa, T., Pajarinen, J., Neumann, G., et al. (2018). Algorithmic perspective on imitation learning. *Foundations and Trends in Robotics*, 7(1–2), 1–179.
- Pace, A., Chan, A., & van der Schaar, M. (2022). POETREE: Interpretable policy learning with adaptive decision trees. In *ICLR*.
- Páez, A. (2019). The pragmatic turn in explainable artificial intelligence (XAI). *Minds and Machines*, 29(3), 441–459.
- Paischer, F., Adler, T., & Hofmarcher, M., et al. (2023). Semantic helm: A human-readable memory for reinforcement learning. In *NeurIPS*.
- Pasula, H. M., Zettlemoyer, L. S., & Kaelbling, L. P. (2007). Learning symbolic models of stochastic domains. In *JAIR*.
- Payani, A., & Fekri, F. (2019a). Inductive logic programming via differentiable deep neural logic networks. [arXiv:1906.03523](https://arxiv.org/abs/1906.03523)
- Payani, A., & Fekri, F. (2019b). Learning algorithms via neural logic networks. [arXiv:1904.01554](https://arxiv.org/abs/1904.01554)
- Payani, A., & Fekri, F. (2020). Incorporating Relational Background Knowledge into Reinforcement Learning via Differentiable Inductive Logic Programming. [arXiv:2003.10386](https://arxiv.org/abs/2003.10386)
- Penkov, S., & Ramamoorthy, S. (2019). Learning programmatically structured representations with perceptor gradients. In *ICLR*.
- Plumb, G., Al-Shedivat, M., & Cabrera, AA., et al. (2020). Regularizing black-box models for improved interpretability. [arXiv:1902.06787](https://arxiv.org/abs/1902.06787)
- Pomerleau, D. (1989). *Alvinn: An autonomous land vehicle in a neural network*. In *NeurIPS*.
- Puiutta, E., & Veith, E. M. (2020). Explainable reinforcement learning: A survey. In *LNCS*.
- Puterman, M. (1994). *Markov decision processes: Discrete stochastic dynamic programming*. Wiley.
- Qiu, W., & Zhu, H. (2022). Programmatic reinforcement learning without oracles. In *ICLR*.
- Rafailov, R., Sharma, A., & Mitchell, E., et al. (2023). Direct preference optimization: Your language model is secretly a reward model. In *NeurIPS*.

- Raji, I. D., Smart, A., & White, R. N., et al. (2020). Closing the AI accountability gap: defining an end-to-end framework for internal algorithmic auditing. [arXiv:2001.00973](https://arxiv.org/abs/2001.00973)
- Ramesh, A., Pavlov, M., & Goh, G., et al. (2021). Zero-shot text-to-image generation. [arXiv:2102.12092](https://arxiv.org/abs/2102.12092)
- Randlov, J., & Alstrom, P. (1998). Learning to drive a bicycle using reinforcement learning and shaping. In *ICML*.
- Redmon, J., Divvala, S., & Girshick, R., et al. (2016). You only look once: Unified, real-time object detection. In *CVPR*.
- Ribeiro, M. T., Singh, S., & Guestrin, C. (2016a). Model-Agnostic Interpretability of Machine Learning. In *ICML workshop on human interpretability in ML*.
- Ribeiro, M. T., Singh, S., & Guestrin, C. (2016b). "Why Should I Trust You?": Explaining the Predictions of Any Classifier. In *KDD*.
- Rocktäschel, T., Singh, S., & Riedel, S. (2015). Injecting logical background knowledge into embeddings for relation extraction. In *Human language technologies*.
- Rombach, R., Blattmann, A., & Lorenz, D., et al. (2022). High-resolution image synthesis with latent diffusion models. In *CVPR*.
- Ross, S., Gordon, G. J., & Bagnell, J. A. (2011). A reduction of imitation learning and structured prediction to no-regret online learning. In *AISTATS*.
- Roth, A. M., Topin, N., & Jamshidi, P., et al. (2019). Conservative Q-Improvement: Reinforcement Learning for an Interpretable Decision-Tree Policy. [arXiv:1907.01180](https://arxiv.org/abs/1907.01180)
- Rothkopf, C. A., & Dimitrakakis, C. (2011). Preference elicitation and inverse reinforcement learning. In *ECML*.
- Rudin, C. (2019). Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5), 206–215.
- Rudin, C., & Carlson, D. (2019). The secrets of machine learning: ten things you wish you had known earlier to be more effective at data analysis. In *Operations research & management science in the age of analytics* (pp. 44–72).
- Russell, S. (1998). Learning agents for uncertain environments. In *COLT*.
- Rusu, A. A., Colmenarejo, S. G., Gülçehre, Ç., et al. (2016). Policy distillation. In *ICLR*.
- Sanchez-Gonzalez, A., Heess, N., & Springenberg, J. T., et al. (2018). Graph networks as learnable physics engines for inference and control. In *ICML*.
- Sanner, S. (2005). Simultaneous learning of structure and value in relational reinforcement learning. In *ICML workshop on rich representations for RL*.
- Sanner, S. (2011). Relational dynamic influence diagram language (RDDDL): Language description. In *International planning competition*.
- Santoro, A., Raposo, D., Barrett, D. G. T., et al. (2017). A simple neural network module for relational reasoning. In *NeurIPS*.
- Scarselli, F., Gori, M., Tsoi, A. C., et al. (2009). The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1), 61–80.
- Scholz, J., Levihn, M., & Isbell, C. L., et al. (2014). A physics-based model prior for object-oriented MDPs. In *ICML*.
- Schulman, J., Wolski, F., & Dhariwal, P., et al. (2017). Proximal policy optimization algorithms. [arXiv:1707.06347](https://arxiv.org/abs/1707.06347)
- Sequeira, P., & Gervasio, M. (2020). Interestingness elements for explainable reinforcement learning: Understanding agents' capabilities and limitations. *Artificial Intelligence*, 288, 103367.
- Serafini, L., & d'Avila Garcez, A. (2016). Logic tensor networks: Deep learning and logical reasoning from data and knowledge. In *CEUR workshop*.
- Shi, W., Huang, G., & Song, S., et al. (2020). Self-supervised discovering of interpretable features for reinforcement learning. [arXiv:2003.07069](https://arxiv.org/abs/2003.07069)
- Shu, T., Xiong, C., & Socher, R. (2018). Hierarchical and interpretable skill acquisition in multi-task reinforcement learning. In *ICLR*.
- Silva, A., & Gombolay, M. (2020). Neural-encoding Human Experts' Domain Knowledge to Warm Start Reinforcement Learning. [arXiv:1902.06007](https://arxiv.org/abs/1902.06007)
- Silva, A., Gombolay, M., & Killian, T., et al. (2020). Optimization methods for interpretable differentiable decision trees applied to reinforcement learning. In *AISTATS*.
- Silver, D., Schrittwieser, J., Simonyan, K., et al. (2017). Mastering the game of Go without human knowledge. *Nature*, 550, 354–359.
- Singh, C., Askari, A., Caruana, R., et al. (2023). Augmenting interpretable models with large language models during training. *Nature Communications*, 14, 7913.
- Slaney, J., & Thiébaux, S. (2001). Blocks world revisited. *Artificial Intelligence*, 125(1–2), 119–153.

- Sridharan, M., Gelfond, M., Zhang, S., et al. (2019). REBA: A refinement-based architecture for knowledge representation and reasoning in robotics. *JAIR*, 65, 87–180.
- Srinivasan, S., & Doshi-Velez, F. (2020). Interpretable batch IRL to extract clinician goals in ICU hypotension management. In *AMIA joint summits on translational science*.
- Sun, S. H., Wu, T. L., & Lim, J. J. (2020). Program guided agent. In *ICLR*.
- Sutton, R. S., & Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT Press
- Sutton, R. S., Precup, D., & Singh, S. (1999). Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artificial Intelligence*, 112(1–2), 181–211.
- Swain, M. (2013). Knowledge Representation. In *Encyclopedia of Systems Biology* (pp. 1082–1084).
- Tang, Y., Nguyen, D., & Ha, D. (2020). Neuroevolution of self-interpretable agents. In *GECCO*.
- Tasse, G. N., James, S., & Rosman, B. (2020). A boolean task algebra for reinforcement learning. In *NeurIPS*.
- Tasse, G. N., James, S., & Rosman, B. (2022). Generalisation in lifelong reinforcement learning through logical composition. In *ICLR*.
- Todorov, E. (2009). Compositionality of optimal control laws. In *NeurIPS*.
- Topin, N., & Veloso, M. (2019). Generation of policy-level explanations for reinforcement learning. In *AAAI*.
- Topin, N., Milani, S., & Fang, F., et al. (2021). Iterative bounding MDPs: Learning interpretable policies via non-interpretable methods. In *AAAI*.
- Toro Icarte, R., Klassen, T., & Valenzano, R., et al. (2018a). Using reward machines for high-level task specification and decomposition in reinforcement learning. In *ICML*.
- Toro Icarte, R., Klassen, T. Q., & Valenzano, R., et al. (2018b). Teaching multiple tasks to an rl agent using LTL. In *AAMAS*.
- Toro Icarte, R., Waldie, E., & Klassen, T., et al. (2019). Learning reward machines for partially observable reinforcement learning. In *NeurIPS*.
- Torrey, L., & Taylor, M. E. (2013). Teaching on a budget: Agents advising agents in reinforcement learning. In *AAMAS*.
- van der Maaten, L., & Hinton, G. (2008). Visualizing data using t-SNE. *JMLR Sci* 9(86), 2579–2605.
- van der Waa, J., van Diggelen, J., van den Bosch, K., et al. (2018). Contrastive explanations for reinforcement learning in terms of expected consequences. In *IJCAI workshop on XAI*.
- van Otterlo, M. (2005). *A survey of reinforcement learning in relational domains*. CTIT Technical Report Series: Tech. rep.
- van Otterlo, M. (2009). *The logic of adaptive behavior: Knowledge representation and algorithms for adaptive sequential decision making under uncertainty in first-order and relational domains*. IOS Press.
- van Otterlo, M. (2012). Solving relational and first-order logical markov decision processes: A Survey. In M. Wiering & M. van Otterlo (Eds.), *Reinforcement learning* (Vol. 12, pp. 253–292). Berlin Heidelberg: Springer.
- Vasic, M., Petrovic, A., & Wang, K., et al. (2019). MoET: Interpretable and verifiable reinforcement learning via mixture of expert trees. [arXiv:1906.06717](https://arxiv.org/abs/1906.06717)
- Vaswani, A., Shazeer, N., & Parmar, N., et al. (2017). Attention is all you need. In *NeurIPS*.
- Veerapaneni, R., Co-Reyes, J. D., & Chang, M., et al. (2020). Entity abstraction in visual model-based reinforcement learning. In *CoRL*.
- Verma, A., Murali, V., & Singh, R., et al. (2018). Programmatically interpretable reinforcement learning. In *ICML*.
- Verma, A., M. Le, H., & Yue, Y., et al. (2019). Imitation-projected programmatic reinforcement learning. In *NeurIPS*.
- Vinyals, O., Ewalds, T., & Bartunov, S., et al. (2017). StarCraft II: A new challenge for reinforcement learning. [arXiv:1708.04782](https://arxiv.org/abs/1708.04782)
- Vinyals, O., Babuschkin, I., Czarnecki, W. M., et al. (2019). Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature*, 575(7782), 350–354.
- Viola, P., & Jones, M. (2001). Robust real-time object detection. In *International journal of computer vision*.
- Walker, T., Shavlik, J., & Maclin, R. (2004). Relational reinforcement learning via sampling the space of first-order conjunctive features. In *ICML workshop on relational reinforcement learning*.
- Walker, T., Torrey, L., & Shavlik, J., et al. (2008). Building relational world models for reinforcement learning. In *LNCS*.
- Walsh, J. (2010). Efficient learning of relational models for sequential decision making. PhD thesis, Rutgers.
- Wang, T., Liao, R., & Fidler, S. (2018). NerveNet: Learning Structured Policy with Graph Neural Networks. In: *ICLR*
- Wang, W., & Pan, S. J. (2019). Integrating deep learning with logic fusion for information extraction. In *AAAI*.

- Wang, Y., Mase, M., & Egi, M. (2020). Attribution-based salience method towards interpretable reinforcement learning. In *Spring symposium on combining ml and knowledge engineering in practice*.
- Weng, P., Busa-Fekete, R., Hüllermeier, E. (2013). Interactive Q-learning with ordinal rewards and unreliable tutor. In *ECML workshop on RL with generalized feedback*.
- Whittlestone, J., Arulkumaran, K., & Crosby, M. (2021). The societal implications of deep reinforcement learning. *JAIR*, 70, 1003–1030.
- Wiegrefe, S., & Pinter, Y. (2019). Attention is not not Explanation. In *EMNLP*.
- Wiener, N. (1954). The human use of human beings. Houghton Mifflin
- Wojke, N., Bewley, A., & Paulus, D. (2017). Simple online and realtime tracking with a deep association metric. In *2017 IEEE international conference on image processing (ICIP)*.
- Wu, B., Gupta, J. K., & Kochenderfer, M. J. (2019a). Model primitive hierarchical lifelong reinforcement learning. In *AAMAS*.
- Wu, M., Parbhoo, S., & Hughes, M. C., et al. (2019b). Optimizing for interpretability in deep neural networks with tree regularization. [arXiv:1908.05254](https://arxiv.org/abs/1908.05254)
- Wu, Z., Geiger, A., & Potts, C., et al. (2023). Interpretability at scale: Identifying causal mechanisms in alpaca. In *NeurIPS*.
- Xu, J., Zhang, Z., & Friedman, T., et al. (2018). A semantic loss function for deep learning with symbolic knowledge. In *ICML*.
- Xu, Z., Gavran, I., & Ahmad, Y., et al. (2020). Joint inference of reward machines and policies for reinforcement learning. In *ICAPS*.
- Yang, F., Yang, Z., & Cohen, W. W. (2017). Differentiable learning of logical rules for knowledge base reasoning. In *NeurIPS*.
- Yang, F., Lyu, D., Liu, B., et al. (2018a). PEORL: Integrating symbolic planning and hierarchical reinforcement learning for robust decision-making. In *IJCAI*.
- Yang, Y., & Song, L. (2019). Learn to explain efficiently via neural logic inductive learning. In *ICLR*.
- Yang, Y., Morillo, I. G., & Hospedales, T. M. (2018b). Deep neural decision trees. In *ICML workshop on human interpretability in ML*.
- Younes, L. (2004). PPDDL1.0: The language for the probabilistic part of IPC-4.
- Yu, H., Shen, Z., & Miao, C., et al. (2018). Building ethics into artificial intelligence. In *IJCAI*.
- Zahavy, T., Ben-Zrihem, N., & Mannor, S. (2016). Graying the black box: Understanding DQNs. In *ICML*.
- Zambaldi, V., Raposo, D., & Santoro, A., et al. (2019). Deep reinforcement learning with relational inductive biases. In *ICLR*.
- Zhang, A., Sukhbaatar, S., & Lerer, A., et al. (2018a). Composable planning with attributes. In *ICML*.
- Zhang, C., Vinyals, O., & Munos, R., et al. (2018b). A Study on Overfitting in Deep Reinforcement Learning. [arXiv:1804.06893](https://arxiv.org/abs/1804.06893)
- Zhang, H., Gao, Z., & Zhou, Y., et al. (2019). Faster and Safer Training by Embedding High-Level Knowledge into Deep Reinforcement Learning. [arXiv:1910.09986](https://arxiv.org/abs/1910.09986)
- Zhang, S., & Sridharan, M. (2020). A Survey of Knowledge-based Sequential Decision Making under Uncertainty. [arXiv:2008.08548](https://arxiv.org/abs/2008.08548)
- Zhang, Y., Lee, J. D., & Jordan, M. I. (2016). L1-regularized neural networks are improperly learnable in polynomial time. In *ICML*.
- Zhu, G., Huang, Z., & Zhang, C. (2018). Object-oriented dynamics predictor. In *NeurIPS*.
- Zhu, G., Wang, J., & Ren, Z., et al. (2020). Object-oriented dynamics learning through multi-level abstraction. In *AAAI*.
- Zhu, H., Magill, S., & Xiong, Z., et al. (2019). An inductive synthesis framework for verifiable reinforcement learning. In *ACM SIGPLAN conference on PLDI*.
- Zimmer, M., Viappiani, P., & Weng, P. (2014). Teacher-student framework: A reinforcement learning approach. In *AAMAS workshop on autonomous robots and multirobot systems*.
- Zimmer, M., Feng, X., & Glanois, C., et al. (2021). Differentiable logic machines. [arXiv:2102.11529](https://arxiv.org/abs/2102.11529)

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.