

© 2023 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

Citation: K. Acharya, W. Raza, C. Dourado, A. Velasquez and H. H. Song, "Neurosymbolic Reinforcement Learning and Planning: A Survey," in IEEE Transactions on Artificial Intelligence, doi: 10.1109/TAI.2023.3311428.

DOI: <https://doi.org/10.1109/TAI.2023.3311428>

Access to this work was provided by the University of Maryland, Baltimore County (UMBC) ScholarWorks@UMBC digital repository on the Maryland Shared Open Access (MD-SOAR) platform.

Please provide feedback

Please support the ScholarWorks@UMBC repository by emailing scholarworks-group@umbc.edu and telling us what having access to this work means to you and why it's important to you. Thank you.

Neurosymbolic Reinforcement Learning and Planning: A Survey

Kamal Acharya[✉], *Graduate Student Member, IEEE*, Waleed Raza[✉], *Graduate Student Member, IEEE*, Carlos Dourado[✉], *Member, IEEE*, Alvaro Velasquez[✉], *Member, IEEE* and Houbing Song[✉], *Fellow, IEEE*

Abstract—The area of Neurosymbolic Artificial Intelligence (Neurosymbolic AI) is rapidly developing and has become a popular research topic, encompassing sub-fields such as Neurosymbolic Deep Learning (Neurosymbolic DL) and Neurosymbolic Reinforcement Learning (Neurosymbolic RL). Compared to traditional learning methods, Neurosymbolic AI offers significant advantages by simplifying complexity and providing transparency and explainability. Reinforcement Learning, a long-standing AI concept that mimics human behavior using rewards and punishment, is a fundamental component of Neurosymbolic RL, a recent integration of the two fields that has yielded promising results. The aim of this paper is to contribute to the emerging field of Neurosymbolic RL by conducting a literature survey. Our evaluation focuses on the three components that constitute Neurosymbolic RL: neural, symbolic, and RL. We categorize works based on the role played by the neural and symbolic parts into three taxonomies, which are further divided into sub-categories based on their applications. Furthermore, we analyze the RL components of each research work, including the state space, action space, policy module, and RL algorithm. Additionally, we identify research opportunities and challenges in various applications within this dynamic field.

Impact Statement—Neurosymbolic RL has captured the interest of both the academic and industrial communities, as researchers strive to develop a reliable and robust model capable of achieving practical performance. Despite this, there is a lack of a comprehensive documented survey that delves into and scrutinizes the field of Neurosymbolic RL as a whole. While several survey papers devoted to Neurosymbolic AI and many more concerning RL are available, there has been no noteworthy contribution that surveys the intersection of these areas. As a result, the purpose of this article is to bridge this gap by presenting a broad range of relevant papers that have been published, with a focus on the three main elements of Neurosymbolic RL: neural, symbolic, and RL. The article conducts an analysis, identifies potential research opportunities, along with the challenges.

Index Terms—Neurosymbolic, Neurosymbolic reinforcement learning, reinforcement learning

Manuscript received March 10, 2023. This work was supported in part by the U.S. National Science Foundation under Grant No. 2309760 and Grant No. 2317117.

K. Acharya is with the Department of Electrical Engineering and Computer Science, Embry-Riddle Aeronautical University, Daytona Beach, FL 32114 USA (e-mail: acharyk2@my.erau.edu).

W. Raza is with the Department of Electrical Engineering and Computer Science, Embry-Riddle Aeronautical University, Daytona Beach, FL 32114 USA (e-mail: razaw@my.erau.edu).

C. Dourado is with the Department of Electrical Engineering and Computer Science, Embry-Riddle Aeronautical University, Daytona Beach, FL 32114 USA (e-mail: douradoc@erau.edu).

A. Velasquez is with the Department of Computer Science, University of Colorado, Boulder, CO 80309 USA (e-mail: alvaro.velasquez@colorado.edu).

H. Song is with the Department of Information Systems, University of Maryland, Baltimore County, Baltimore, MD 21250 USA (e-mail: songh@umbc.edu).

I. INTRODUCTION

NEUROSYMBOLIC Artificial Intelligence (Neurosymbolic AI), a budding field of Artificial Intelligence, has garnered significant attention in recent times as it combines both neural and symbolic traditions to enhance the performance of neural network models. In this context, the term "neural" pertains to Neural Network (NN) primarily, while "symbolic" refers to the use of various mathematical logic and algorithms for symbolic manipulation. Reinforcement Learning (RL), another emerging area of machine learning, revolves around agents operating in various environments to maximize their rewards. It dates back to the early days of cybernetics and has gained rapid interest in the machine learning and artificial intelligence communities over the last five to ten years. RL involves programming agents by reward and punishment without specifying how to accomplish the task, and it encompasses statistics, psychology, neuroscience, and computer science. However, there are significant computational challenges to overcome[1]. Deep Reinforcement Learning (DRL), which replaces tabular methods of estimating state values with function approximation, has eliminated the need to store all state value pairs in the table, enabling the agent to generalize the value of states that it has never encountered before. DRL has been utilized in programs that have defeated the best human players in game of Go[2]. Additionally, an AI agent named AlphaStar[3] beat the world's best StarCraft II player.

RL has recently drawn much attention in the context of Neurosymbolic AI for policy synthesis and representation. These techniques merge planning-style control-flow instructions with fundamental atomic actions that are learned and represented through (deep) neural networks. The combination of these two approaches enables the efficient use of deep reinforcement learning techniques to improve the interpretability and transparency of an agent's behavior while also leveraging a high-level, symbolic representation of the policies learned by agents. By allowing the neural system to interact with the knowledge base, the reasoning ability is enhanced, and the learning ability is enhanced by interacting with the neural system. This interaction results in better generalization and transfer of knowledge, improved efficiency and robustness, and an increase in explanation and interpretability. Fig.1 illustrates the general idea of combining Neurosymbolic AI with RL.

Further research is needed in Neurosymbolic RL to develop novel approaches, techniques, and their real-time applications that best fit real-world use cases such as computer networks,

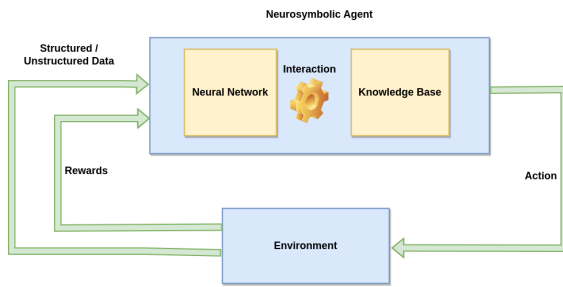


Fig. 1. An Overview of Neurosymbolic RL Process

healthcare, IoT devices, finance, and other industrial domains. In this paper, we have analyzed notable works in Neurosymbolic RL to date. We have examined the neural and symbolic component used in each of the research work. Further we analyzed RL components used in the architecture: RL-algorithms, state space, action space, and policy module used so that we can have transparent view of the working of the model. We have classified these research works into three main categories: Learning for Reasoning RL model, Reasoning for Learning RL model, and Learning-Reasoning RL model, which are further sub-divided according to the significance or role of the model. After that We move on to provide a comprehensive summary of the Neurosymbolic RL approaches related to their specific application cases that need to be developed to meet the needs of AI. Finally, we have presented certain challenges specific to each application case to employ them in real-world scenarios.

In reviewing the history of AI surveys, no significant work has been found that specifically focused on the combination of Neurosymbolic AI with RL. Earlier works are either focused solely on Neurosymbolic AI or on the field of RL. Negligible surveys can be found on Neurosymbolic AI[4], [5] few other provides insight on recent advances[6] and application[7]. A large number of surveys are available on RL on various aspects:

- RL in general[1], DRL[8], [9], Causal RL[10]
- Safety and Security in RL [11], [12]
- Environment [13]
- Agent [14], [15], [16]
- Application like Natural Language Processing [17], [18], Communication Network [19], Robotics [20], Healthcare [21], Transportation [22]

This survey is the first of its kind and the first attempt to evaluate the combination of these two popular areas as one (Neurosymbolic RL). In this survey, we provide insights into all the relevant work done in the past under various taxonomies, along with possible opportunities to address the challenges.

The following is the document's structure: Section II provides an overview of milestones in the AI field from its inception to the present day. In Section III, we present an overview of Neurosymbolic AI and RL, covering relevant literature and significant research findings. Section IV is dedicated to Neurosymbolic Reinforcement Learning, including workable architectures and requirements. In Section V, we summarize notable research in Neurosymbolic RL under vari-

ous headings. In Section VI, we discuss opportunities that have emerged from Neurosymbolic RL. Section VII is devoted to the challenges of implementing proposed Neurosymbolic RL applications. We identify the obstacles and challenges that may arise. Finally, in Section VIII, we offer concluding remarks on our survey paper.

II. MILESTONE IN REINFORCEMENT LEARNING

Reinforcement Learning (RL) has a rich history dating back to the 1940s when B.F. Skinner introduced the concept of operant conditioning in psychology, while Walter Pitts and Warren McCulloch[23] presented a computational model based on the functioning of the human brain. Donald Hebb's Hebbian Learning Rule[24] also formed the basis for modern neural networks. In the late 1950s, Frank Rosenblatt developed the perceptron, which could learn based on associationism. Grigoryevich[25] used complex polynomial equations to statistically analyze network elements, selecting the best ones for the next layer, laying the groundwork for what would become deep learning. During this time, Richard Bellman also developed the mathematical formulation of RL and introduced the Bellman equation for dynamic programming[26]. Later, Temporal Difference Learning (TD Learning)[27] was introduced, which enabled agents to learn from delayed rewards and gradually update their value estimates. In the 1970s, the field of reinforcement learning experienced a reduction in funding, leading only a few scientists to continue their work independently. Nonetheless, during this time, significant progress was achieved. Fukushima developed the Neocognitron neural network[24], which utilized a hierarchical multi-layer architecture to enable computers to learn visual patterns. The Neocognitron later served as a basis for the convolutional neural network that is widely used today. Additionally, Paul Werbos introduced the backpropagation algorithm[24], which, although not widely used at the time, raised questions in cognitive psychology regarding the role of symbolic logic in human comprehension.

The 1980s saw the emergence of the field of RL with the introduction of Actor-Critic Algorithms and Q-learning[28]. In the following decade, the field continued to evolve with the introduction of core algorithms such as REINFORCE and SARSA. A pivotal breakthrough occurred in 1999 with the invention of the Graphics Processing Unit (GPU), which enabled RL to tackle more complex environments. This was further enhanced by the parallel computing power of NVIDIA's Compute Unified Device Architecture (CUDA) on GPUs.

The 2010s proved to be a remarkable decade for RL. In 2012, the introduction of the Arcade Learning Environment (ALE) opened the gateway to the use of RL in gaming environments. Deep RL, which combines neural networks with RL to learn high-dimensional state-action value functions, was introduced, leading to breakthroughs in game playing and robotics, such as the Deep Q-Network (DQN). Many researchers became active in modifying existing algorithms, resulting in the development of numerous new algorithms in the RL domain, such as Trust Region Policy Optimization (TRPO), Deep Deterministic Policy Gradient (DDPG)[29],



Fig. 2. RL Milestone Timeline

Soft Actor Critic(SAC)[30], Quality Value Iteration Optimization (QT-Opt)[31] and various variations of DQN, including Double DQN, Dueling DQN[32], and Rainbow DQN[33]. The introduction of OpenAI Gym, an open-source toolkit for developing and comparing reinforcement learning algorithms, opened the door for exploring RL algorithms among the RL community members. In 2016, RL achieved a significant milestone in the history of AI by defeating the world champion in the game of Go. RL continued to conquer many other games against humans, such as Dota 2 and Starcraft II. In 2017, AlphaZero was introduced, which empowered RL models and beat humans in chess, shogi, and Go with a significant margin without human training[34].

As we have progressed into the 2020s, RL continues to be a dynamic field of research, with researchers exploring novel algorithms and applications, such as multi-agent RL, meta-RL, and RL for safety-critical systems. In addition, AlphaZero has been successful in discovering faster multiplication algorithms[35] which showed that RL can also contribute in other field also as a superhuman. Recently, in late 2022, OpenAI released ChatGPT, a chatbot that utilizes RL techniques to be trained and generate diverse responses to various inquiries and concerns. Other major technology companies are also in the race to develop their own innovative AI solutions. A summary of significant milestones in the history of RL is presented in Fig.2.

III. BACKGROUND AND PRELIMINARIES

In this section, we keep focus on the background information related to the topics Neurosymbolic AI followed by RL.

A. Neurosymbolic AI

The field of artificial intelligence has been centered around the goal of developing machines that can achieve human-like levels of intelligence. Two major approaches have been pursued in this effort. The first, symbolic AI, is a rule-based approach that was prevalent from the 1950s to the 1980s. The second approach is a data-based approach known as connectionist AI. While symbolic AI requires a large amount of information to be supplied, it can learn from this information on its own. The primary disadvantage of connectionist AI is its inability to explain the reasoning or logical processes behind the model, leading to these models being referred to as black boxes. Symbolic reasoning provides an explainable inference process and employs powerful declarations to represent knowledge, as well as offering benefits such as fast initial coding, explicit method control, and abstraction of knowledge[36]. However, this approach is limited in its ability to handle vast amounts of incomplete data and to generalize from such data. Psychologist Daniel Kahneman has distinguished between two human cognitive processes, system 1 and system 2. System 1 is fast, automatic, and unconscious, akin to deep learning, while system 2 is slow, effortful, and conscious, similar to symbolic AI[37]. In the context of AI, there have been discussions of ways to combine these two approaches, as the authors of a study[38] conclude that only a combination of both fields is likely to enable the development of human-like intelligence.

Neurosymbolic AI is a subfield of AI that combines two historically prominent approaches: connectionist AI and symbolic AI. This integration enables more efficient derivation of knowledge and general concepts from data, focusing on learning from experience and reasoning about what has been learned from uncertain environments. Hybrid Neurosymbolic systems require less training data and are capable of tracking the steps required to draw conclusions and make inferences which is the reason Neurosymbolic AI has been regarded as the 3rd wave of AI[39]. By combining symbolic reasoning with deep learning, ideal results can be obtained with a limited number of datasets, error correction with recoveries, and enhanced explanatory capabilities that are not possible with deep learning alone[40]. Numerous applications necessitate both learning and reasoning abilities. On the neural aspect, models learn from data provided to them, while the symbolic aspect seeks to retain the innate explanatory power of these systems. The Neurosymbolic AI domain, as previously discussed, can be employed to develop various applications across different fields, such as medical diagnostic systems, recommender systems, and text mining[41]. By incorporating deep human expert knowledge into the system's design and function, Neurosymbolic AI can be leveraged to its fullest potential in creating such applications. Fig.3 depicts the evaluation of the Neurosymbolic AI process within a model design that integrates neural network and symbolic artificial intelligence, harnessing the full strength of both fields in hybrid models.

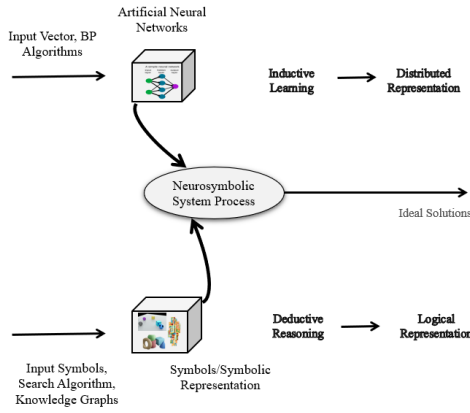


Fig. 3. General Representation Neurosymbolic AI

Numerous researchers have provided insights into how the fields of neural and symbolic AI can be combined practically. Three noteworthy works have contributed significantly to organizing the research in Neurosymbolic systems. The first notable work was a survey paper published in 2005 by Sebastian and Pascal [42]. They identified three main axes of Neurosymbolic integration: Interrelation, Language, and Usage. Each of these axes was further divided into several sub-divisions. Fig.4 provides a simplified visualization of the eight dimensions along with their axes. Another researcher, Henry Kautz[43], proposed a way to classify Neurosymbolic systems into six different categories. He gave them distinctive names, which are detailed in Table I. In a recent survey[6], Neurosymbolic systems were analyzed based on three parameters: efficiency, generalization, and interpretability. The authors proposed a novel taxonomy consisting of three different classes: learning for reasoning, reasoning for learning, and learning-reasoning. Table II provides detailed descriptions of each of these classes.

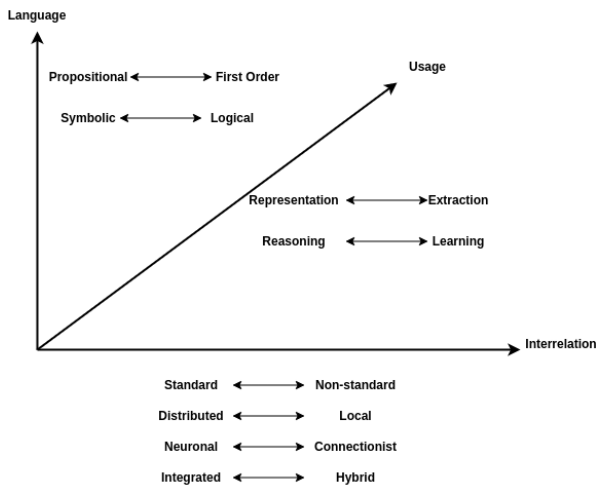


Fig. 4. Classification by Sebastian and Pascal[42]

TABLE I
CLASSIFICATION BY HENRY KAUTZ[43]

Classification	Characteristic Features
Symbolic Neuro symbolic	Symbolic input is converted to feature vectors for the neural networks which give final results in the symbolic form
Symbolic[Neuro]	Neural pattern recognition subroutine within a symbolic problem solver
Neuro Symbolic	A cascade from neural network into symbolic reasoner
Neuro: Symbolic \rightarrow Neuro	Symbolic rules are input which are compiled so that their knowledge end up in the neural network
Neuro_{Symbolic}	Uses direct encodings of logical statements into neural structures
Neuro[Symbolic]	Embed symbolic reasoning inside neural engine to enable both superhuman and super combinatorial reasoning

TABLE II
CLASSIFICATION BY D. YU AND ET AL. [6]

Classification	Characteristic Features
Learning for reasoning	Neural network play the role of the helper, it extracts the important symbols and information so that the search space of the symbolic system narrowed down
Reasoning for learning	Symbolic system act as a helper, it provides symbolic knowledge to the neural network from where the final decision is made
Learning-reasoning	Uses symbolic and neural systems as an alternate process. They both complement each other to give the final results

B. Reinforcement Learning

The fast-learning algorithms and wide-ranging applications of RL have made it increasingly popular in academia and industry, thanks to significant technological advancements [44], [45]. In earlier literature, RL was described as a class of problems that an agent encounters in a dynamic, unpredictable environment and solves through trial and error. Nowadays, RL is viewed as a machine learning paradigm that trains an agent to make decisions based on its immediate surroundings to optimize rewards. The training process involves a loop of interaction with the environment, including observing, receiving rewards, making decisions, and obtaining feedback signals [46]. RL has proven its ability to solve complex real-world problems, such as natural language processing, image classification, speech recognition, and decision-making, which has improved planning and perception in various applications [47]. RL is an essential component of autonomous driving cars and robots, which can perform tasks such as food preparation without human intervention or specific programming. RL-based strategies could play a crucial role in enabling fully autonomous systems in the future [48]. RL employs algorithms and methods to enable an agent to obtain optimal control in an environment, and the agents in RL can range from a game player to a stock trading bot. Another field which is similar to RL is Intrinsic Motivation(IM) but it lack feedback mechanism. Many research papers in RL have utilized IM to address complex problems in sparse reward platforms [49],

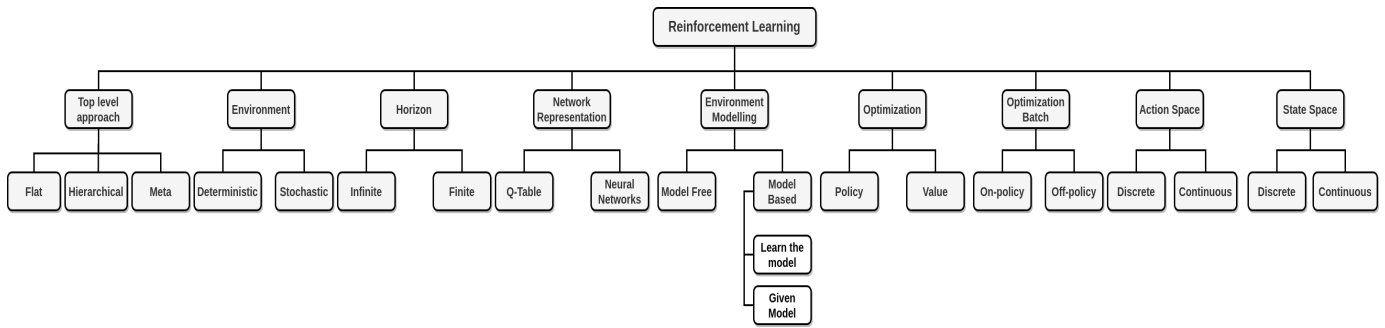


Fig. 5. RL Methods and Techniques

[50]. The interactions that occur between an agent and its environment are typically modeled as a Markov Decision Problem (MDP)[51], or a Partially Observable Markov Decision Process (POMDP). An MDP is a framework used for sequential decision-making in Markovian dynamical systems, which extends the Multi-Armed Bandits (MAB) framework by allowing the system state to stochastically change based on the actions taken and their resulting outcomes. On the other hand, a POMDP is a newer version of MDP, where the system state is not directly observable. In certain cases, MDPs can be solved analytically, while in many cases, they can be solved iteratively through the use of dynamic or linear programming. When no model is present, RL methods can be employed to obtain sample trajectories and directly interact with the system[52]. As the number of computing devices continues to increase rapidly, it is expected that the number of devices capable of handling complex and dynamic systems with minimal programming will grow exponentially, potentially reaching billions.

Looking at the bigger picture, RL can be classified as a type of sample-based approach for solving MDP problems. The RL technique uses sample trajectories and the agent's interaction with the system, which can be obtained from a simulation. This approach is quite common in practical applications, where a simulation is available, and a clear transition-probability model is not required. In such scenarios, dynamic or linear programming may not be suitable, making the RL method a more practical option[53]. Main Components of Reinforcement Learning are:

- **Policy:** It refers to the way an agent behaves at a given time, which is generally a mapping from perceived states to the action that needs to be taken when in those states of the environment. The goal of the policy is to maximize the expected cumulative reward received by the agent over time. There are various types of policies, including deterministic policies and stochastic policies, which define the agent's behavior in different ways.
- **Reward Signal:** It refers to the objective or goals of the problem and is a number delivered to the agent by the environment at each time step. The reward signal is used to train the agent to learn a behavior that maximizes the cumulative reward over time. It is a crucial component as it guides the agent to take actions that lead to achieving the desired goals.

- **Value Function:** It represents the expected long-term cumulative reward that an agent can obtain by following a specific policy. It estimates the value of each state or state-action pair, which allows the agent to choose the best action in each state. The value function can be expressed mathematically as the expected sum of discounted future rewards starting from a given state or state-action pair. The estimation of the value function can be done through various methods, such as Monte Carlo methods, Temporal Difference learning, or Bellman equations.
- **Model of the Environment:** It refers to the representation of how the environment behaves in response to the actions taken by the agent. It allows the agent to predict the next state and reward given the current state and action. The model can be either known or unknown, and the goal is to use it to optimize the agent's behavior. In cases where the model is known, dynamic programming techniques can be used to find the optimal policy. The model can be represented in different forms, including transition probabilities, state-transition diagrams, or function approximators.

Reinforcement Learning (RL) can be categorized into various types based on different parameters, such as the environment, policy, model, and others. These categories provide a framework to understand and classify different RL approaches. Fig.5 summarizes the various types of RL based on these parameters, including environment type, horizon type, optimization type, and more.

The success of RL largely depends on the quality of its algorithm. Numerous RL algorithms have been developed, tailored to specific contexts, and based on various parameters such as environment type, action space type, and model type. These algorithms are continuously modified to improve their performance and expand their scope of applications[54]. Fig.6 provides a brief overview of the various types of RL algorithms that have been used to date.

IV. NEUROSymbolic Reinforcement Learning

RL, a long-standing topic in the field of AI, has faced the curse of dimensionality, but the introduction of DRL solved this problem. However, DRL has several limitations. For instance, DRL can be extremely data-inefficient. In a paper by Deepmind [33], they demonstrated that the Rainbow DQN

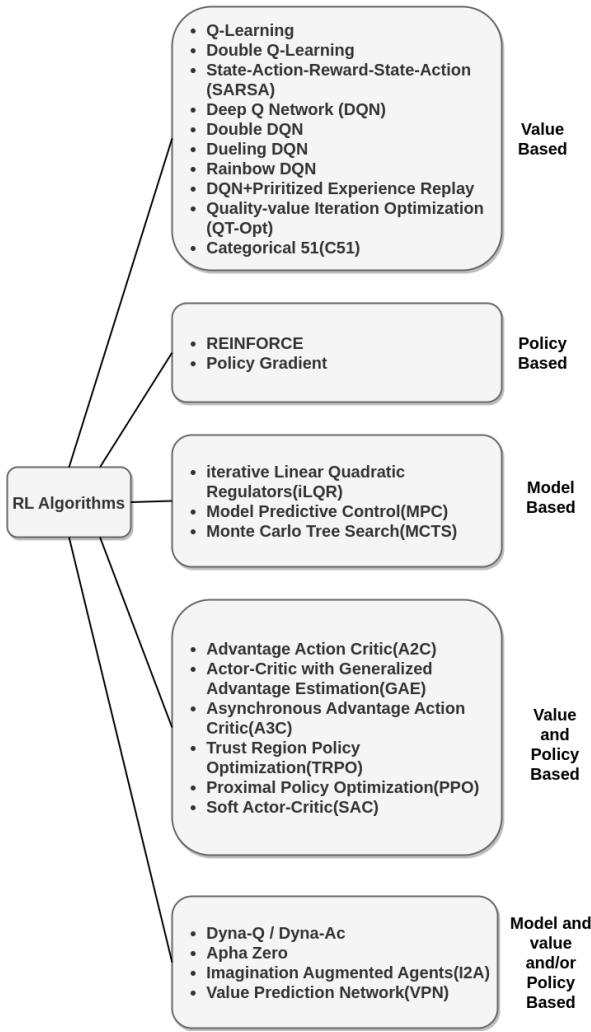


Fig. 6. Overview of RL Algorithms

method can achieve state-of-the-art performance in terms of both performance and data efficiency. Nevertheless, it required almost 83 hours (about 3 and a half days) of playtime in addition to the training time. Conversely, many people can achieve this level of performance in just a few minutes. Another issue with DRL is that, except for rare scenarios, domain-specific algorithms work better than DRL. In the field of robotics, Boston Dynamics¹ is a leading research institution that focuses mainly on classical robotics techniques such as time-varying Linear Quadratic Regulator(LQR), Quadratic Programming(QP) solvers, and convex optimization. Another main issue with RL is the reward system, which can be easily functionalized, but the challenge arises when trying to encourage appropriate behavior while still making it learnable. Sparse rewards are problematic because they only supply rewards in the goal state, making them difficult to shape. Shaped rewards are easier to learn because they provide positive feedback even when the whole solution has not yet been figured out. However, the problem with shaped rewards is that they are biased. The agent becomes focused on maximizing the reward

¹<https://www.bostondynamics.com/research>

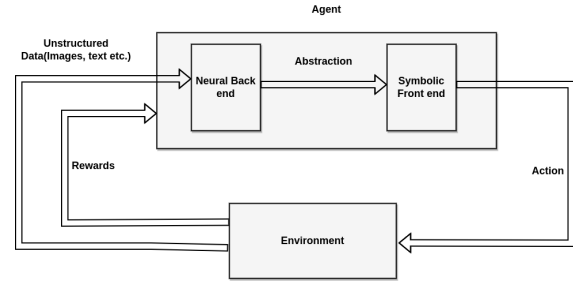


Fig. 7. Learning for Reasoning RL model

instead of finding the complete solution. For example, in a study on text summarization [55], the RL model focused on increasing the ROUGE score, which it succeeded in doing. However, it failed to achieve the actual task of generating readable summaries. In contrast, summarized text generated by the model with lower ROUGE scores was found to be more readable and efficient.

The combination of Neurosymbolic systems and RL appears to be a solution to many of the issues identified in previous DRL methods. This approach not only adds reasoning and explaining capabilities to DRL but also provides a breakthrough in the field of RL. There are multiple ways in which the Neurosymbolic counterpart can be combined with RL, each with its own unique features. In this context, we will discuss three different approaches.

A. Learning for Reasoning RL model

The Learning for Reasoning RL model combines a neural component with a symbolic system to improve reasoning capabilities. The neural component functions as a co-actor, while the symbolic system handles the problem of reasoning. The DNNs in the model help to reduce the symbolic space, leading to faster convergence and improved performance. In cases where the data presented to the model are unstructured, DNNs can transform them into a symbolic form that the symbolic system can utilize. Furthermore, DNNs can also distill the learning policy to the symbolic system, which enhances verifiability. Serialization characterizes the neural and symbolic counterparts in this model, as shown in Fig.7.

B. Reasoning for Learning RL model

The Reasoning for Learning RL model is a different approach that utilizes symbolic models to guide the output of the neural network. By incorporating structured knowledge from the symbolic system, the performance and interpretability of the DNNs can be improved. The symbolic model can also help with reward shaping to enable faster convergence and improved performance of the DNNs agent. Additionally, the symbolic system can aid in generating the programmatic policy, making the RL model more interpretable and explainable. This type of model is characterized by parallelization, as shown in Fig.8.

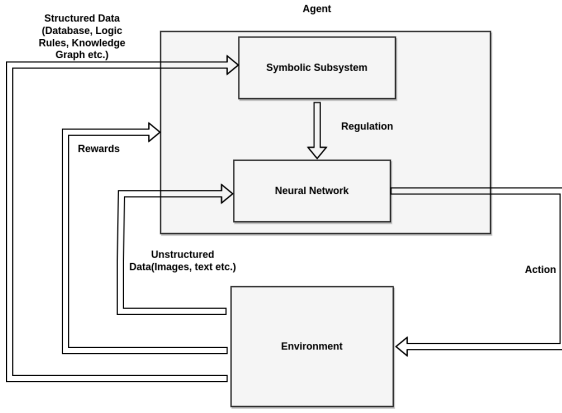


Fig. 8. Reasoning for Learning RL model

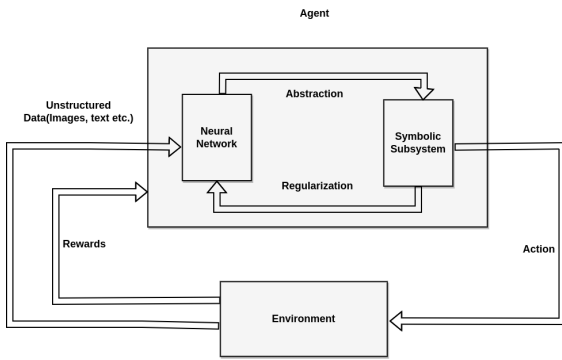


Fig. 9. Learning-Reasoning RL model

C. Learning-Reasoning RL model

In the Learning-Reasoning RL model, the neural and symbolic components work bidirectionally, where the output of one can be the input of the other. This approach combines the benefits of both the Learning for Reasoning RL and Reasoning for Learning RL models, resulting in a balanced combination of interpretability and reasoning capacity. The symbolic part provides the structured knowledge to DNNs to enhance their interpretability and performance, while the neural component reduces the symbolic space, enabling the symbolic counterpart to achieve faster convergence. Moreover, the two parts work together which make the RL model more interpretable and explainable. This type of model is characterized by bidirectional communication, as depicted in Fig.9.

V. RELATED WORKS

Neurosymbolic Reinforcement Learning (RL) is an emerging area of AI that is currently lacking in literature. This section aims to provide an overview of the implementation stages of Neurosymbolic RL, including the state of the art, current trends, and proposed research studies. Various notable works have been analyzed, detailing the neural and symbolic components used in Table III, and information about the RL algorithm, reward space, action space, and policy module in Table IV. These works have been classified into three main RL models: Learning for Reasoning, Reasoning for Learning, and Learning-Reasoning, and further sub-divided according to

their areas of application. A summary of this classification can be found in Table V. The classification and analysis of these works provides insights into the use of Neurosymbolic RL models and their potential for further development in the future.

A. Learning for Reasoning RL model

This particular Neurosymbolic RL model involves the use of a neural network as an auxiliary tool to extract crucial symbols and information, which helps to reduce the search space of the symbolic system. This results in a faster problem-solving process, making it especially useful for problems that require reasoning. This architecture has been used for the following primary goals:

1) *Transforming unstructured data into a symbolic representation*: Symbolic systems, which rely on logical rules and representations of symbolic data, are often limited in their ability to process unstructured data such as images, videos, and natural language text. Most of the real world data are inherently present in the unstructured form so there must be a model to transform them to the symbolic form before being processed by the symbolic models. DNNs have been shown to be very effective in processing and generating such unstructured data and can be used to generate structured data that can be used as input to a symbolic system.

Deep Symbolic Reinforcement Learning (DSRL)[56] consists of two main components: a deep neural network that learns a low-level continuous representation of the state space and map it to low-dimensional symbolic space, and a symbolic model that distills the learned policy into a more interpretable form by mapping symbolic representation to action. The authors use this framework to learn policies for a range of environments, and demonstrate that their system outperforms traditional reinforcement learning algorithms in terms of interpretability, generalization, and efficiency. Symbolic Reinforcement Learning with Common Sense (SRL+CS)[57] a novel extension of DSRL where the authors create a meaningful symbolic representation of the world using sub-states before applying learning and decision-making algorithms. They have two modification in Q-value function: restricting the updates to specific sub-state and assigning importance on the basis of distance of the objects. This approach provides better generalization and explainability. Neural Symbolic Reinforcement Learning (NSRL)[58], includes a reasoning module based on neural attention networks, which performs relational reasoning on symbolic states and induces the RL policy, enabling end-to-end learning with prior symbolic knowledge. It can extract the logical rules selected by the attention modules instead of storing all the rules, saving memory budget and improving scalability. Deep Symbolic Policy[59], uses an autoregressive recurrent neural network to generate symbolic policies, which are optimized using a risk-seeking policy gradient. To scale to environments with multi-dimensional action spaces, the authors propose an "anchoring" algorithm that distills pre-trained neural network-based policies into fully symbolic policies. The authors also introduce two novel methods to improve exploration in DRL-based combinatorial optimization which are hierarchical entropy regularizer and a soft length prior.

Detect, Understand, Act (DUA) [60], composed of three components: Detect, which consists of a traditional computer vision object detector and tracker, Understand, which provides an answer set programming (ASP) paradigm for symbolically implementing a meta-policy over options, and Act, which houses a set of options that are high-level actions enacted by pre-trained DRL policies. The paper evaluates the DUA framework on the Animal-AI (AAI) competition testbed and achieves state-of-the-art results in multiple categories. It is modular approach, allowing for straightforward generalization and transfer to other complex tasks. Another study, Symbolic Options for Reinforcement Learning (SORL)[61], proposes a method for automatically discovering and learning symbolic options, which are higher-level actions with specified preconditions and postconditions, to assist deep reinforcement learning (DRL) agents in complex environments. It was successful in mitigating the problem of sparse and delayed reward along with improving efficiency. Neurosymbolic Logic Neural Network (LNN) for RL algorithm[62], supplies fast convergence and interpretability for RL policies in text-based interaction games by extracting first-order logical facts from text observation using semantic parser(ConceptNet) and history, then trains the symbolic rules with logical functions in the neural networks.

2) *Knowledge Graph Reasoning*: Knowledge Graph (KG) reasoning is the task of inferring new information from a given KG, which consists of a set of entities and their relationships. KG reasoning is important for a wide range of applications, including question answering, information retrieval, and recommender systems.

DeepPath[63], uses a policy-based agent with continuous states based on knowledge graph embeddings to sample the most promising relation and extend the multi-hop relational path. The authors demonstrate that their method outperforms path-ranking based algorithms and knowledge graph embedding methods on two standard reasoning tasks on Freebase and Never-Ending Language Learning datasets. Meandering In Networks of Entities to Reach Verisimilar Answers (MINERVA)[64]outperforms DeepPath, as DeepPath cannot be applied to query answering tasks where the second entity is unknown. It uses neural reinforcement learning to learn how to navigate the knowledge graph conditioned on the input query to find predictive paths.

3) *Verification*: The process of verification aims to determine if a model meets a particular desired property, and can play a key role in enhancing quality and safety. In the context of reinforcement learning (RL), it is important to verify a model's convergence, correctness, and robustness to ensure it functions effectively.

Verifiability via Iterative Policy Extraction(VIPER)[65], is a modification of Q-DAGGER algorithm used for verifying the correctness of deep reinforcement learning policies. The approach involves extracting a small, interpretable model from a deep neural network policy, which can then be verified using existing verification techniques. The extracted model approximates the original policy well and can be used to analyze the policy's convergence, correctness, and robustness. Another model, Reinforcement Learning with Verified Exploration

TABLE III
NEURAL AND SYMBOLIC COMPONENTS IN RELATED WORKS

Research	NN	Knowledge Base
[56]	CNN	First Order Logic
[57]	CNN	First Order Logic
[58]	Transformer	First Order Logic
[59]	RNN	Decision Tree
[60]	DNN	Answer Set Programming
[61]	DNN	Propositional Logic
[62]	LNN	First Order Logic
[63]	DNN	Knowledge Graph
[64]	LSTM	Knowledge Graph
[65]	DNN	Decision Tree
[66]	DNN	Symbolic Policies
[67]	DNN	Decision Tree
[68]	CNN	Finite Trace Linear Temporal Logic
[69]	CNN	Finite Trace Linear Temporal Logic
[70]	DNN	Omega Regular Language
[71]	DNN	Programmatic Policy
[72]	DNN	Programmatic Policy
[73]	NN	State Machines
[74]	RNN	Programmatic Policy
[75]	DNN	Deterministic Finite Automaton
[76]	DNN	Propositional Logic
[77]	CNN	First Order Logic

(REVEL)[66], incorporates a differentiable symbolic planner that generates a set of safe exploration actions, which the RL agent executes to find optimal policies while avoiding potentially unsafe states. The proposed framework is formally verified using the Coq proof assistant, which ensures that the system is free from runtime errors and satisfies desired safety properties.

4) *Gaming*: Neurosymbolic RL is a promising approach in the field of game playing, where the goal is to develop agents that can learn to play games at a human-like level or beyond. It involves combining neural networks and symbolic systems to develop agents that can learn the rules of the game and develop strategies to play the game effectively. By combining symbolic reasoning with deep learning, Neurosymbolic RL models can provide explanations for the decisions made by the agent, making it easier for humans to understand and evaluate the agent's performance. AlphaGo Zero[67], achieves superhuman performance in the game of Go without using human data or domain knowledge. The system is based on a combination of deep neural networks and Monte Carlo tree search, with the neural networks trained through a reinforcement learning process using self-play. It beats the older AlphaGo version in game Go with score of 100-0.

B. Reasoning for Learning RL model

In this type of Neurosymbolic RL model, symbolic system acts as a helper, it provides symbolic knowledge to the neural network from where the final decision is made. This approach is particularly useful in complex applications, such as robotics, where the environment is often uncertain and dynamic, and where the use of symbolic knowledge can facilitate high-level reasoning and decision-making. It has been used in the following application area:

1) *Reward Shaping*: In the field of RL, one of the primary challenges is dealing with sparse rewards. One solution to

this problem is reward shaping, which involves incorporating domain knowledge. Rather than relying on a single, final reward, intermediate rewards are provided to the agent for exhibiting desirable behavior. This encourages the agent to take effective actions early on in the learning process, leading to faster convergence.

Monte Carlo Tree Search with Automaton-Guided Reward Shaping (MCTS-A)[68], helps to improve the learning process and final performance of the agent in domains with sparse rewards. The method introduces an automaton that guides the reward shaping process, allowing for a dynamic and flexible approach. The automaton's states correspond to the different learning phases of the agent, and each state has its own shaping rules that change over time. Transfer learning between the two different environments with the same objective is also eased with this approach. Authors in[69], extended the work by introducing Multiagent Tree Search Algorithm with reward shaping(MATS-A) so that it can be applied to multi-agent scenario and can handle both stochastic and deterministic transition in Multi-agent Non-Markovian Reward Decision Process. They prove that sharing the same search tree and DFA objective can be used to develop competitive and cooperative behavior among the agents, within and across the team. Research work [70] first converts an omega-regular specification into a Buchi automaton. It is then used to construct an average reward objective which can then be optimized by standard RL algorithms. The authors prove that the learned policy converges to the optimal policy and demonstrate the effectiveness of the method.

2) *Programmatic Policy Design*: A programmatic policy refers to a decision-making algorithm that governs the behavior of an agent that can make decisions. The programmatic policy takes inputs from the environment and computes a set of actions that the agent should take in response. The design of programmatic policies can vary based on the complexity of the task and available data, including decision trees, state machines, and programs.

Imitation-Projected Programmatic Reinforcement Learning (PROPEL)[71], proposes a new approach for combining imitation learning with reinforcement learning, called imitation-projected programmatic reinforcement learning (IP-PRL). The approach uses a programmatic policy to encode a priori knowledge about the task and trains an agent through a combination of imitation learning and reinforcement learning. The agent first learns to imitate an expert's behavior through supervised learning, and then the agent's policy is updated through reinforcement learning while being constrained to stay close to the expert's policy. IP-PRL outperforms both pure imitation learning and pure reinforcement learning in terms of sample efficiency and final performance. Another work Programmatically Interpretable Reinforcement Learning (PIRL)[72], uses Neurally Directed Program Synthesis (NDPS) algorithm to generate interpretable neural policies which can be verified through the symbolic approach. It first learn a neural policy network using deep reinforcement learning and then performing a local search over programmatic policies that seeks to minimize the distance from this neural oracle. [73] proposes a novel approach for synthesizing policies

for automated decision-making systems that can generalize to new situations. The authors use inductive programming techniques, specifically "program synthesis by example", to generate policies that satisfy a set of example-based specifications. Framework[74] synthesize programmatic policies that are more interpretable and generalizable than neural network policies produced by deep reinforcement learning methods. It uses a program representation and only requires minimal supervision compared to prior programmatic reinforcement learning and program synthesis works. It learns a program embedding space that parameterizes diverse behaviors in an unsupervised manner and then searches over this space to find a program that maximizes the return for a given task.

3) *Task Segmentation*: Main goal or task is broken down into the smaller task with their own set of rewards so that the task become more generalizable and reasonable. DeepSynth[75], uses automata synthesis to automatically segment a task. A task was broken down into smaller subgoals, each with its own reward. The proposed approach learns a model of an automaton that represents the state machine of a task and uses it to segment the task into subgoals. The learned automaton is then used to guide the agent in finding the optimal policy for the task.

4) *Knowledge Initialized Model*: Researches has found that the model give higher convergence rate, reasoning ability if initialized with knowledge base. Before starting the learning process, the knowledge base of the agent is initialized with some prior information instead of starting from zero. Propositional Logic Nets (PROLONETS)[76], enables warm start of learning process by efficient initialization of RL agents using human-specified policies, without requiring an Imitation Learning (IL) phase. This approach helps RL agents to navigate complex environments that pose challenges to randomly initialized models, and allows for greater exploration. It outperforms baseline RL approaches such as IL and knowledge-based techniques.

C. Learning-Reasoning RL model

This Neurosymbolic RL model uses symbolic and neural systems as an alternate process. They both complement each other by performing abstraction and regularization to give the final results. Symbolic Deep Reinforcement Learning (SDRL) [77], uses planner-controller-meta-controller architecture where planner uses prior symbolic knowledge for long term planning, controller uses DRL algorithms for intrinsic rewards and meta-controller evaluate training performance of controller based on extrinsic rewards along with proposing new intrinsic goals to the planner.

VI. OPPORTUNITIES

Real-world applications strive to minimize errors that can arise from risky exploration and exploitation, whereas Neurosymbolic RL methods employ a trial-and-error mechanism. However, to reconcile this contradiction between Neurosymbolic RL and real-world applications, a viable approach is to create an authentic simulator using real data and domain knowledge of the model dynamics. Subsequently, objectives

TABLE IV
RL COMPONENTS OF THE RELATED WORKS

Research	RL-Algorithm	State Space	Action Space	Policy Module
[56]	Q-Learning	Multi-dimensional vector	Multi-dimensional vector	Tabular Q Learning
[57]	Q-Learning	Multi-dimensional vector	Multi-dimensional vector	Q-Table
[58]	Double DQN	Set of Predicates	Set of Predicates	Multi-layer Perceptron
[59]	Policy Gradients	Multi-dimensional vector	Multi-dimensional vector	RNN
[60]	PPO	Multi-dimensional vector	Discrete action space	NN
[61]	Double Q-Learning	High level dimension sapce	5-dimensional vector	Option Set
[62]	DQN	Multi-dimensional vector	Discrete set of 10 different actions	LNN
[63]	Policy Gradients	Entities in Knowledge Graph	Relations in Knowledge Graph	Fully connected NN
[64]	REINFORCE	Entities in Knowledge Graph	Relations in Knowledge Graph	LSTM
[65]	VIPER	Multi-dimensional vector	Leaf nodes in Decision Tree	Decision Tree
[66]	Policy Gradients	Real vector space	Real vector space	NN
[67]	Policy Gradients	Muti-dimensional vector	Muti-dimensional vector	DNN
[68]	MCTS-A	Nodes of Automata	Transitions of Automata	CNN
[69]	MATS-A	Nodes of Automata	Transitions of Automata	CNN
[70]	Differential Q-Learning	Nodes in GFM automaton	Relation in GFM automaton	DNN
[71]	Policy Gradients	Continuous	Continuous	Programatic and Neural
[72]	NDPS	Unconstrained Policy Space	Continuous	Programatic and Deterministic
[73]	Gradient Based Optimization	Continuous	Continuous	State Machine
[74]	REINFORCE	Program Embedding Space	Program Execution Trace	RNN
[75]	DQN	Multi-dimensional vector	Multi-dimensional vector	DNN
[76]	PPO	193D and 37D	44D and 10D	Decision Tree
[77]	Double Q-Learning	High dimensional	Set of Primitive Action	DNN

TABLE V
SUMMARY OF CLASSIFICATION OF RELATED RESEARCHES

RL model	Areas of Application	Related Researches
Learning for Reasoning	Transforming unstructured data into a symbolic representation	[56][57][58][59][60][61][62]
	Knowledge Graph Reasoning	[63][64]
	Verification	[65][66]
	Gaming	[67]
Reasoning for Learning	Reward Shaping	[68][69] [70]
	Programatic Policy Design	[71][72][73][74]
	Task Segmentation	[75]
	Knowledge Initialized Model	[76]
Learning-Reasoning	Task Segmentation	[77]

can be designed for the agent, and the policy network can be trained in the simulator. Finally, the trained policy can be deployed in the real world with further enhancements. Though Neurosymbolic RL is in its early stage but it has started contributing to other RL areas as well. Causal Reinforcement Learning[78] has been able to produce the significant result since its collaboration with Neurosymbolic RL model. In this section, we examine the opportunities of Neurosymbolic RL methods in various fields.

A. Robotics and Control

Building autonomous embodied robotic systems requires designing suitable policies that ensure the system operates within reasonable mechanical constraints while maintaining safety and data efficiency. RL has been growing in robotics from very old time[79], [80]. The symbolic method has been introduced in robotic motion planning and control in 2007[81] to address these concerns. It is clear that decision-making is a crucial aspect of robotics control, and there have been various approaches to address this challenge. One notable technique is the Neurosymbolic Program Search (NSPS)[82], which produces interpretable and robust Neurosymbolic programs for autonomous driving design. Another approach is the decomposition of decision-making into two levels: what to do and

how to do it[83]. This method utilizes Neurosymbolic skills and has been shown to be effective in various robotics tasks. There have also been efforts to construct robotic platforms for building manipulation environments, such as the open-source platform CausalWorld[84]. Overall, these methods and platforms have been shown to improve policy learning and performance in robotics tasks.

B. Gaming RL

Games are considered as suitable benchmarks with clear-cut rules and boundaries in the RL community. Over the past few years, gaming AI has exhibited extraordinary decision-making abilities, surpassing human-level performance in various decision-making games, including card games[85], board games[2], and video games[3]. Neurosymbolic RL has primarily been utilized in board games and video games, yielding state-of-the-art outcomes. It is expected that these models will also outperform others in card games. However, open-ended games such as Minecraft and XLand remain unexplored.

C. Intelligent Question Answering

Neurosymbolic RL has emerged as a powerful tool in natural language processing for intelligent question answering,

which involves deducing the answer to a given question based on the surrounding context, often comprising both text and images. Although various studies in the context of Neurosymbolic RL have focused on knowledge graph reasoning[63], [64] for this task, combining both text and images remains a relatively unexplored research area. As such, there is considerable potential for future research to investigate and advance this topic.

D. Safe Reinforcement Learning

Reinforcement learning (RL) has become popular due to its ability to learn from experience and make decisions in complex environments. However, in safety-critical settings such as autonomous driving, robotics, or medical diagnosis, the failure of the system can result in severe consequences, including loss of property or human lives. Ensuring the safety of RL agents is crucial in such settings. Neurosymbolic RL has been used for the verification of the RL[65], [66] and it has given some significant result. But, it is still in its infancy period, and there are many opportunities for safely exploring the RL.

E. Optimizing Parameters of RL

The RL framework consists of several components, including the environment, the agent, the policy, and the reward function. Neurosymbolic RL has been applied to different components of the RL framework, combining symbolic reasoning with neural networks to solve complex RL problems. It has been successful in addressing the issue of sparse rewards by formulating reward functions that provide more informative feedback to the agent [68], [69], [70]. It has also been used to learn programmatic policies that are more generalizable and flexible to different environments [71], [72], [73], [74]. Additionally, Neurosymbolic RL has been effective in reducing the symbolic space, resulting in more efficient representations of the policy, and improving the agent's performance [56], [57], [58], [59], [60], [61], [62].

VII. CHALLENGES

Neurosymbolic RL addresses a variety of issues that were previously challenges for DRL and has opened up new opportunities for researchers to develop novel methodologies. In this section, we outline a list of problems that are still prevalent with Neurosymbolic RL, including some that are specific to DRL and others that are more general research gaps.

A. Automated Generation of Symbolic Knowledge

Neurosymbolic RL relies on an environment where the agent can interact and receive rewards. Typically, these environments are represented by symbolic knowledge, as explained in the previous section. Symbolic knowledge encompasses both logic rules and knowledge graphs. While research into the automatic construction of non-logical symbolic part like knowledge graphs is relatively mature[86], [87], [88], the automatic learning of logic rules from data remains an underexplored area. Typically, domain experts manually construct the

logic which is a time-consuming, laborious, and non-scalable process. Additionally, achieving end-to-end learning for rules that describe prior knowledge from data is a challenge for Neurosymbolic systems. Moreover, the inclusion of intricate logic, probabilistic relations, or diverse data sources adds further complexity to the problem. We contend that greater attention should be given to the comprehensive and automatic discovery of symbolic knowledge, not only from increasingly vast data sets but also from networks with rapidly expanding dimensionality.

B. Verification and Validation

Neurosymbolic RL models have gained popularity across multiple industries, with their size increasing rapidly to enable deployment in larger scenarios. These models have achieved state-of-the-art results and have provided a degree of reasoning and explainability. However, due to the relative novelty of this field, there is a lack of validation and verification methods for these models, which need to be addressed. For instance, despite AI surpassing humans in the game of Go in 2016, recent adversarial attacks on the models have exposed their weaknesses and led to humans defeating similar AI models in the game[89]. This highlights the significant gap in the verification field, which requires extensive work to ensure that Neurosymbolic RL models are thoroughly validated and can be deployed without any flaws. Some work[65], [66], [90] has been initiated in this in this direction but prior work [91], [92] also need expansion so that they can be applied to Neurosymbolic RL domain.

C. Neurosymbolic RL Algorithms

The combination of neural, symbolic, reinforcement learning allows for a more comprehensive approach to problem-solving, as it enables the system to work with both numerical and symbolic data for RL. This provides a more powerful and flexible framework for learning, allowing for the integration of different types of knowledge and reasoning techniques for the agent. However, in order to effectively combine these fields, new learning algorithms need to be designed that can take advantage of the strengths of both symbolic and neural learning so as to be implemented in RL. The traditional reinforcement learning algorithms may not be accurate enough for Neurosymbolic learning, as they do not account for the complexities and nuances of symbolic reasoning. Therefore, new algorithms need to be optimized to work under the union of two sets of knowledge, leveraging the strengths of both neural and symbolic learning[93]. By doing so, researchers can develop more accurate and efficient learning algorithms that can be applied to a wide range of problems in fields such as natural language processing, robotics, and healthcare, among others.

D. Balancing Reasoning and Learning in RL

Neurosymbolic RL requires training the neural components using meaningful symbolic constraints and allowing the symbolic components to evolve with high-quality data-driven

rules. However, transitioning between neural and symbolic components can lead to a loss of learning or reasoning power, which presents scalability challenges for the field. One crucial issue in Neurosymbolic RL is how to align symbolic specifications with representations learned using neural methods, known as the symbol grounding problem[94], [95]. This challenge is well-known in AI, but it is particularly complicated in Neurosymbolic RL, where symbolic and neural components can be interwoven in intricate ways.

VIII. CONCLUSION

In recent years, there has been a remarkable growth in the field of Neurosymbolic Reinforcement Learning (RL). This survey provides a comprehensive overview of Neurosymbolic RL, which can be classified into three RL models: Learning for Reasoning, Reasoning for Learning, and Learning-Reasoning. We have examined each category's core area of application and conducted an in-depth analysis of its various components, including neural, symbolic, and RL. Additionally, we have highlighted future opportunities for exploration and the potential challenges that might arise. Our hope is that this survey will inspire the AI community to delve deeper into this area and explore its possibilities.

REFERENCES

- [1] L. P. Kaelbling, M. L. Littman, and A. W. Moore. "Reinforcement learning: A survey." *Journal of artificial intelligence research*, vol. 4, pp. 237-285, 1996.
- [2] D. Silver et al., "Mastering the game of Go with deep neural networks and tree search", *nature*, vol. 529, no. 7587, pp. 484-489, 2016.
- [3] R. Liu et al. "An Introduction of mini-AlphaStar.", 2021, arXiv:2104.06890.
- [4] T. R. Besold et al., "Neural-symbolic learning and reasoning: A survey and interpretation.", 2017, arXiv:1711.03902.
- [5] W. Wang and Y. Yang., "Towards Data-and Knowledge-Driven Artificial Intelligence: A Survey on Neuro-Symbolic Computing.", 2022, arXiv:2210.15889.
- [6] D. Yu et al. "Recent Advances in Neural-symbolic Systems: A Survey" 2021, arXiv:2111.08164.
- [7] D. Bouneffouf and C. C. Aggarwal., "Survey on Applications of Neurosymbolic Artificial Intelligence." ,2022, arXiv:2209.12618.
- [8] K. Arulkumaran et al., "Deep reinforcement learning: A brief survey." *IEEE Signal Processing Magazine*, vol.34, no. 6, pp. 26-38, 2017.
- [9] Y. Li, "Deep reinforcement learning: An overview.", 2017, arXiv:1701.07274.
- [10] Y. Zeng et al., "A Survey on Causal Reinforcement Learning.", 2023, arXiv:2302.05209.
- [11] J. Garcia and F. Fernández. "A comprehensive survey on safe reinforcement learning." *Journal of Machine Learning Research*, vol.16, no. 1, pp.1437-1480, 2015.
- [12] A. Upreti and D. B. Rawat., "Reinforcement learning for iot security: A comprehensive survey.", *IEEE Internet of Things Journal*, vol. 8, no. 11, pp.8693-8706, 2020.
- [13] S. Padakandla "A survey of reinforcement learning algorithms for dynamically varying environments." *ACM Computing Surveys (CSUR)*, vol.54, no. 6, pp.1-25,2021.
- [14] L. Buşoniu, R. Babuška, and B. D. Schutter. "Multi-agent reinforcement learning: An overview.", *Innovations in multi-agent systems and applications-1*, pp.183-221,2010.
- [15] L. Canese et al., "Multi-agent reinforcement learning: A review of challenges and applications." *Applied Sciences*, vol.11, no. 11, pp.4948, 2021.
- [16] S. Gronauer and K. Diepold. "Multi-agent deep reinforcement learning: a survey." *Artificial Intelligence Review*, pp. 1-49, 2022.
- [17] V. Uc-Cetina et al., "Survey on reinforcement learning for language processing." *Artificial Intelligence Review*, pp. 1-33, 2022.
- [18] J. Luketina et al. "A survey of reinforcement learning informed by natural language.", 2019, arXiv:1906.03926.
- [19] Y. Qian et al. "Survey on reinforcement learning applications in communication networks." *Journal of Communications and Information Networks*, vol.4, no. 2, pp.30-39, 2019.
- [20] J. Kober, J. A. Bagnell, and J. Peters. "Reinforcement learning in robotics: A survey." *The International Journal of Robotics Research*, vol.32, no. 11, pp.1238-1274, 2013.
- [21] C. Yu et al., "Reinforcement learning in healthcare: A survey.", *ACM Computing Surveys (CSUR)*, vol.55, no. 1, pp.1-36, 2021.
- [22] A. Haydari and Y. Yilmaz., "Deep reinforcement learning for intelligent transportation systems: A survey.", *IEEE Transactions on Intelligent Transportation Systems*, vol.23, no. 1, pp.11-32, 2020.
- [23] W. S. MCCULLOCH and W. PITTS "A Logical Calculus of Ideas Immanent in Nervous Activity", *Bulletin of Mathematical Biophysics*, vol. 5, p127-147, 1943.
- [24] H. WANG and B. RAJ, "On the Origin of Deep Learning", 2017, arXiv:1702.07800v4.
- [25] A. G. Ivakhnenko, "Polynomial Theory of Complex Systems," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. SMC-1, no. 4, pp. 364-378, Oct. 1971, doi: 10.1109/TSMC.1971.4308320.
- [26] R. Bellman, "Dynamic programming.", *science*, vol.153, no. 3731, pp. 34-37, 1966.
- [27] G. Tesauro. "Temporal Difference Learning and TD-Gammon", *Communications of the ACM*, vol. 38, no. 3, pp. 58-68, 1995.
- [28] C. J. Watkins and P. Dayan "Q-learning. *Machine learning*", vol.8, pp.279-292, 1992.
- [29] T. P. Lillicrap et al. "Continuous Control with Deep Reinforcement Learning." 2015, arxiv: 1509.02971.
- [30] T. Haarnoja et al., "Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor." 2018, arxiv: 1801.01290.
- [31] D. Kalashnikov et al., "QT-Opt: Scalable Deep Reinforcement Learning for Vision-Based Robotic Manipulation." 2018, arxiv: 1806.10293.
- [32] Z. Wang et al. "Dueling Network Architectures for Deep Reinforcement Learning." 2015, arxiv: 1511.06581.
- [33] M. Hessel and et al., "Rainbow: Combining improvements in deep reinforcement learning." in *Proc of the AAAI conference on artificial intelligence*, 2018, vol. 32, no. 1.
- [34] D. Silver et al. "Mastering Chess and Shogi by Self-Play with a General Reinforcement Learning Algorithm.", 2017, arxiv: 1712.01815.
- [35] A. Fawzi, et al., "Discovering faster matrix multiplication algorithms with reinforcement learning", *Nature*, vol.610, pp. 47-53, 2022.
- [36] J. Mao et al. "The neuro-symbolic concept learner: Interpreting scenes, words, and sentences from natural supervision." 2019, arXiv:1904.12584.
- [37] D. Kahneman, *Thinking, fast and slow*. Macmillan, 2011.
- [38] G. Booch et al., "Thinking Fast and Slow in AI", 2020, arXiv:2010.06002.
- [39] A. D. A. Garcez and L. C. Lamb, "Neurosymbolic AI: the 3rd wave", 2020 , arXiv:2012.05876.
- [40] L. C. Lamb et al. "Graph neural networks meet neural-symbolic computing: A survey and perspective." 2020, arXiv:2003.00330.
- [41] M. K. Sarker et al. "Neuro-symbolic artificial intelligence: Current trends." 2021, arXiv:2105.05330.
- [42] S. Bader and P. Hitzler, "Dimensions of neural-symbolic integration-a structured survey", 2005, arXiv:cs/0511042.
- [43] H. A. Kautz, 2022. "The third AI summer: AAAI Robert S. Engelmore Memorial Lecture" *AI Magazine*, vol. 43, no.1, pp.105-125, 31 March 2022. [Online]. Available: <https://onlinelibrary.wiley.com/doi/full/10.1002/aaai.12036>
- [44] J. Salvador, J. Oliveira, and M. Breternitz, "REINFORCEMENT LEARNING: A LITERATURE REVIEW", Sep. 2020.
- [45] S. Balhara et al. "A survey on deep reinforcement learning architectures, applications and emerging trends" *IET Communication*, 2022.[Online]. Available: <https://doi.org/10.1049/cmu2.12447>
- [46] A. Gosavi, "Reinforcement learning: A tutorial survey and recent advances." *INFORMS Journal on Computing*, vol.21, no. 2, pp.178-192, 2009.
- [47] H. Wang et al. "Deep reinforcement learning: a survey" *Front Inform Technol Electron Eng*, vol. 21, pp. 1726-1744, 2020.
- [48] S. Milani and et al. "A Survey of Explainable Reinforcement Learning", 2022, arXiv:2202.08434.
- [49] A. Aubret, L. Maignon and S. Hassas, S "A survey on intrinsic motivation in reinforcement learning", 2019, arXiv:1908.06976.
- [50] P. Yadav et al., "A Survey on Deep Reinforcement Learning-based Approaches for Adaptation and Generalization", 2022, arXiv:2202.08444.
- [51] M. L. Puterman, "Markov decision processes" in *Stochastic Models*, Elsevier, 1990, ch.8, pp.331-434

- [52] Z. Qin, H. Zhu and J. Ye "Reinforcement learning for ridesharing: An extended survey." *Transportation Research Part C: Emerging Technologies*, vol. 144, pp.103852, 2022.
- [53] A. Kai et al., "Deep reinforcement learning: A brief survey", *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 26-38, 2017
- [54] J. Zhu, F. Wu, and J. Zhao. "An Overview of the Action Space for Deep Reinforcement Learning.", in 4th International Conference on Algorithms, Computing and Artificial Intelligence, 2021, pp. 1-10.
- [55] P. Romain, C. Xiong, and R. Socher. "A deep reinforced model for abstractive summarization.", 2017, arXiv:1705.04304
- [56] M. Garnelo, K. Arulkumaran and M. Shanahan "Towards deep symbolic reinforcement learning", 2016, arXiv:1609.05518.
- [57] A. D. Garcez et al. "Towards symbolic reinforcement learning with common sense.", 2018, arXiv:1804.08597.
- [58] Z. Ma et al. "Learning Symbolic Rules for Interpretable Deep Reinforcement Learning", 2021, arXiv:2103.08228.
- [59] M. Landajuela et al. "Discovering symbolic policies with deep reinforcement learning." in *Proc of the 38th International Conference on Machine Learning*, 2021, vol. 139, pp. 5979-5989.
- [60] L. Mitchener et al. "Detect, Understand, Act: A Neuro-symbolic Hierarchical Reinforcement Learning Framework." *Machine Learning*, vol.111, no. 4, pp.1523-1549, 2022.
- [61] M. Jin et al. "Creativity of ai: Automatic symbolic option discovery for facilitating deep reinforcement learning." in *Proc. of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 6, pp. 7042-7050, 2022.
- [62] D. Kimura et al. "Neuro-Symbolic Reinforcement Learning with First-Order Logic.", 2021, arXiv:2110.10963
- [63] W. Xiong, T. Hoang, and W. Y. Wang. "DeepPath: A reinforcement learning method for knowledge graph reasoning.", 2017, arXiv:1707.06690.
- [64] R. Das et al. "Go for a walk and arrive at the answer: Reasoning over paths in knowledge bases using reinforcement learning", 2017, arXiv:1711.05851.
- [65] O. Bastani, Y. Pu and A. Solar-Lezama. "Verifiable Reinforcement Learning via Policy Extraction", 2018, arXiv:1805.08328
- [66] G. Anderson et al. "Neurosymbolic reinforcement learning with formally verified exploration." *Advances in neural information processing systems*, vol. 33, pp. 6172-6183, 2020.
- [67] D. Silver et al. "Mastering the game of go without human knowledge." *nature* vol. 550, no. 7676, pp. 354-359, 2017.
- [68] A. Velasquez et al., "Dynamic Automaton-Guided Reward Shaping for Monte Carlo Tree Search" in *Proc. of the AAAI Conference on Artificial Intelligence*, vol.35, no.13, pp.12015-12023, 2021
- [69] A. Velasquez et al., "Multi-Agent Tree Search with Dynamic Reward Shaping" in *Proc. of the International Conference on Automated Planning and Scheduling*, vol.32, pp.52-61, 2022
- [70] M. Kazemi et al. "Translating omega-regular specifications to average objectives for model-free reinforcement learning." in *Proc. of the 21st international conference on autonomous agents and multiagent systems*, pp. 732-741, 2022.
- [71] A. Verma et al. "Imitation-Projected Programmatic Reinforcement Learning", 2019, arXiv:1907.05431
- [72] A. Verma et al. "Programmatically Interpretable Reinforcement Learning", 2018, arXiv:1804.02477
- [73] J. P. Inala et al. "Synthesizing Programmatic Policies that Inductively Generalize", in *Proc. of the International Conference on Learning Representations*, 2020
- [74] D. Trivedi et al. "Learning to synthesize programs as interpretable and generalizable policies." in *Advances in neural information processing systems*, vol.34, pp. 25146-25163, 2021.
- [75] M. Hasanbeig et al. "DeepSynth: Program synthesis for automatic task segmentation in deep reinforcement learning.", 2019, arXiv:1911.10244.
- [76] A. Silva and M. Gombolay, "Encoding human domain knowledge to warm start reinforcement learning," in *AAAI*, vol. 35, no. 6, pp. 5042-5050, 2021.
- [77] D. Lyu et al. "SDRL: interpretable and data-efficient deep reinforcement learning leveraging symbolic planning." in *Proc. of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, pp. 2970-2977, 2019.
- [78] S. Zhu, I. Ng, and Z. Chen., "Causal discovery with reinforcement learning.", 2019, arXiv:1906.04477.
- [79] J. Kober, J. A. Bagnell, and J. Peters, "Reinforcement learning in robotics: A survey," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1238-1274, 2013.
- [80] M. P. Deisenroth, G. Neumann, J. Peters et al., "A survey on policy search for robotics," *Foundations and Trends® in Robotics*, vol. 2, no. 1-2, pp. 1-142, 2013.
- [81] C. Belta, A. Bicchi, M. Egerstedt, E. Frazzoli, E. Klavins, and G. Pappas, "Symbolic planning and control of robot motion," *Robotics & Automation Magazine*, IEEE, vol. 14, pp. 61 – 70, 2007.
- [82] J. Sun, H. Sun, T. Han, and B. Zhou, "Neuro-symbolic program search for autonomous driving decision module design," in *Conference on Robot Learning*, pp.21-30, 2021.
- [83] T. Silver et al., "Learning neuro-symbolic skills for bilevel planning," in *Conference on Robot Learning*, 2022.
- [84] O. Ahmed et al, "Causalworld: A robotic manipulation benchmark for causal structure and transfer learning," in *International Conference on Learning Representations*, 2020.
- [85] D. Zha et al., "Douzero: Mastering douzidzhu with self-play deep reinforcement learning," *international conference on machine learning*, 2021.
- [86] L. Qiao et al., "Knowledge graph construction techniques," *Journal of computer research and development*, vol. 53, no. 3, p. 582, 2016.
- [87] J. L. Martinez-Rodriguez et al., "Openie-based approach for knowledge graph construction from text," in *Expert Syst. Appl.*, vol. 113, pp. 339-355, 2018.
- [88] S. Ji et al., "A survey on knowledge graphs: Representation, acquisition, and applications," *IEEE Trans. Neural Netw. Learning Sys.*, vol. 33, no. 2, pp. 494-514, 2021.
- [89] R. Waters, "Man beats machine at Go in human victory over AI", *FINANCIAL TIMES*, 19 February 2023. [Online]. Available: <https://arstechnica.com/information-technology/2023/02/man-beats-machine-at-go-in-human-victory-over-ai/>
- [90] A. Karimi and P. S. Duggirala, "Formalizing traffic rules for uncontrolled intersections," *2020 ACM/IEEE 11th International Conference on Cyber-Physical Systems (ICCPs)*, Sydney, NSW, Australia, 2020, pp. 41-50, doi: 10.1109/ICCPs48487.2020.00012.
- [91] J. Barnat et al., "DiVinE: Parallel Distributed Model Checker," *2010 Ninth International Workshop on Parallel and Distributed Methods in Verification, and Second International Workshop on High Performance Computational Systems Biology*, Enschede, Netherlands, 2010, pp. 4-7, doi: 10.1109/PDMC-HiBi.2010.9.
- [92] C. Barrett, A. Stump, and C. Tinelli. "The SMT-LIB Standard - Version 2.0". in *Proc. of the 8th International Workshop on Satisfiability Modulo Theories (Edinburgh, England)*, 2010.
- [93] S. Chaudhuri et al., "Neurosymbolic programming", *Foundations and Trends® in Programming Languages*, vol. 7, no. 3, pp. 158-243, 2021.
- [94] S. Harnad, "The symbol grounding problem", *Physica D: Nonlinear Phenomena*, vol.42, no.1, pp.335-346, 1990
- [95] R. J. Mooney, "Learning to Connect Language and Perception." in *AAAI*, pp. 1598-1601, 2008.



Kamal Acharya (Graduate Student Member, IEEE) received his Engineering degree in Electronics and Communication Engineering from Tribhuvan University, Kathmandu, Nepal in 2011 and Masters degree in Information System Engineering from Purbanchal University, Kathmandu, Nepal in 2019. Currently, he is pursuing PhD. in Electrical Engineering and Computer Science from Embry-Riddle Aeronautical University, Daytona Beach, FL.

He has been involved in teaching profession for about 7 years in the various universities of Nepal, Tribhuvan University and Purbanchal University were among few of them. He is mainly associated with the courses like programming (C++, Python), Computer Networks and Computer Architecture. He is working as Graduate Research Assistant in Embry-Riddle Aeronautical University. He is also serving as an reviewer for IEEE Transactions on Artificial Intelligence (TAI) and IEEE Transactions on Intelligent Transportation Systems. His preferred areas of research are Natural Language Processing (NLP), Deep Learning and Reinforcement Learning.



Waleed Raza received a B.E. degree in Electronic Engineering from the Department of Electronic Engineering, Dawood University of Engineering and Technology Karachi, Pakistan in 2017. He received an M.S. degree in underwater acoustic communication engineering from the College of underwater acoustic engineering, Harbin Engineering University, Harbin China where he researched OFDM communication for underwater technology. He holds the editorial board member for Engineering, Technology, and Applied Science Research (ETASR) (2021-

present). He is an active reviewer of a few journals including IEEE Sensor Journal, IEEE Access, and International Journal of Electronics and Communications (2018-2022), he has recently joined the IEEE as a student member. Currently, he is pursuing a Ph.D. degree at Embry Riddle Aeronautical University in Electrical and Computer Engineering. His research area of interest includes underwater acoustic OFDM communication, underwater acoustic target detection, artificial intelligence, machine learning for communication engineering, and autonomous unmanned systems such as UAVs and their characteristics.



Carlos M. J. M. Dourado Jr received the Ph.D. degree in Informatics from the University of Fortaleza, Ceara, Brazil in February 2019 and MSc in Teleinformatics Engineering from the PPGETI/UFC (UFC, 2008). He completed a BSe in Electronics Engineering at the University of Fortaleza (Unifor, 2004). He is an associate professor and researcher at the Department of Telematics (DTEL)/Graduate Program in Computer Engineering (PPGCC) at the Federal Institute of Ceara (IFCE), Brazil and Post-Doctoral Research Associate in Embry-Riddle Aeronautical University. His main research areas include Internet of Things and Artificial Intelligence.



Alvaro Velasquez is a program manager in the Innovation Information Office (I2O) of the Defense Advanced Research Projects Agency (DARPA), where he currently leads the Assured Neuro-Symbolic Learning and Reasoning (ANSR) program. Before that, Alvaro oversaw the machine intelligence portfolio of investments for the Information Directorate of the Air Force Research Laboratory (AFRL). Alvaro received his PhD in Computer Science from the University of Central Florida and is a recipient of the National Science Foundation Graduate Research

Fellowship Program (NSF GRFP) award, the University of Central Florida 30 Under 30 award, and best paper and patent awards from AFRL. He has co-authored 60 papers and two patents and serves as Associate Editor of the IEEE Transactions on Artificial Intelligence and his research has been funded by the Air Force Office of Scientific Research.



Houbing Song (M'12–SM'14–F'23) received the Ph.D. degree in electrical engineering from the University of Virginia, Charlottesville, VA, in August 2012.

He is currently a Tenured Associate Professor, the Director of NSF Center for Aviation Big Data Analytics (Planning), the Associate Director for Leadership of the DOT Transportation Cybersecurity Center for Advanced Research and Education (Tier 1 Center), and the Director of the Security and Optimization for Networked Globe Laboratory

(SONG Lab, www.SONGLab.us), University of Maryland, Baltimore County (UMBC), Baltimore, MD. Prior to joining UMBC, he was a Tenured Associate Professor of Electrical Engineering and Computer Science at Embry-Riddle Aeronautical University, Daytona Beach, FL. He serves as an Associate Editor for IEEE Transactions on Artificial Intelligence (TAI) (2023-present), IEEE Internet of Things Journal (2020-present), IEEE Transactions on Intelligent Transportation Systems (2021-present), and IEEE Journal on Miniaturization for Air and Space Systems (J-MASS) (2020-present). He was an Associate Technical Editor for IEEE Communications Magazine (2017-2020). He is the editor of eight books, the author of more than 100 articles and the inventor of 2 patents. His research interests include cyber-physical systems/internet of things, cybersecurity and privacy, and AI/machine learning/big data analytics. His research has been sponsored by federal agencies (including National Science Foundation, US Department of Transportation, and Federal Aviation

Administration, among others) and industry. His research has been featured by popular news media outlets, including IEEE GlobalSpec's Engineering360, Association for Uncrewed Vehicle Systems International (AUVSI), Security Magazine, CXOTech Magazine, Fox News, U.S. News & World Report, The Washington Times, and New Atlas.

Dr. Song is an IEEE Fellow, an ACM Distinguished Member, and an ACM Distinguished Speaker. Dr. Song is a Highly Cited Researcher identified by Clarivate™ (2021, 2022) and a Top 1000 Computer Scientist identified by Research.com. He received Research.com Rising Star of Science Award in 2022 (World Ranking: 82; US Ranking: 16). Dr. Song was a recipient of 10+ Best Paper Awards from major international conferences, including IEEE CPSCOM-2019, IEEE ICII 2019, IEEE/AIAA ICNS 2019, IEEE CBDCOM 2020, WASA 2020, AIAA/ IEEE DASC 2021, IEEE GLOBECOM 2021 and IEEE INFOCOM 2022.