

# Scenario-assisted Deep Reinforcement Learning

Raz Yerushalmi<sup>1</sup>, Guy Amir<sup>2</sup><sup>a</sup>, Achiya Elyasaf<sup>3</sup><sup>b</sup>, David Harel<sup>1</sup>, Guy Katz<sup>2</sup> and Assaf Marron<sup>1</sup><sup>c</sup>

<sup>1</sup>Department of Computer Science and Applied Mathematics, Weizmann Institute of Science, Rehovot 76100, Israel

<sup>2</sup>School of Computer Science and Engineering, The Hebrew University of Jerusalem, Givat Ram, Jerusalem 91904, Israel

<sup>3</sup>Software and Information Systems Engineering, Ben-Gurion University of the Negev, Beer-Sheva 8410501, Israel

**Keywords:** Machine Learning, Scenario-based Modeling, Rule-based Specifications, Domain Expertise.

**Abstract:** Deep reinforcement learning has proven remarkably useful in training agents from unstructured data. However, the opacity of the produced agents makes it difficult to ensure that they adhere to various requirements posed by human engineers. In this work-in-progress report, we propose a technique for enhancing the reinforcement learning training process (specifically, its reward calculation), in a way that allows human engineers to directly contribute their expert knowledge, making the agent under training more likely to comply with various relevant constraints. Moreover, our proposed approach allows formulating these constraints using advanced model engineering techniques, such as scenario-based modeling. This mix of black-box learning-based tools with classical modeling approaches could produce systems that are effective and efficient, but are also more transparent and maintainable. We evaluated our technique using a case-study from the domain of internet congestion control, obtaining promising results.


## 1 INTRODUCTION


*Deep neural networks (DNNs)* have proven highly successful in addressing hard-to-specify cognitive tasks. *Deep reinforcement learning (DRL)* is a particular method for producing DNNs, which is applicable in cases where the training data is unstructured — e.g., in games (Ye et al., 2020), in autonomous driving (Kiran et al., 2021), smart city communications (Xia et al., 2021), manufacturing (Li et al., 2022), chat bots (Mohamad Suhaili et al., 2021), context-aware systems (Elyasaf, 2021), and many others. As this trend continues, it is likely that DRL will gain a foothold in many systems of critical importance.


Despite the success of DRL, and in particular its generally superior performance to that of hand-crafted code in many kinds of applications, various problematic aspects of this paradigm have begun to emerge. One issue is that DRL agents are typically trained on some distribution of inputs (described, e.g., using a Markov Decision Process), but this distribution might differ from the distribution

that the agent encounters after deployment (Eliyahu et al., 2021). Another issue is various vulnerabilities that exist in many DNNs, and in particular in DRL agents (Szegedy et al., 2013), such as sensitivity to adversarial perturbations. When such issues are discovered, the DRL agent often needs to be modified or enhanced; but unfortunately, such routine actions are known to be extremely difficult for DRL agents, which are largely considered “black boxes” (Goodfellow et al., 2016). Specifically, their underlying DNNs are opaque to the human eye, making them hard to interpret; and the DRL training process itself is computationally expensive and time-consuming, making it infeasible to retrain the agent whenever circumstances, or requirements, change. Much research is being conducted on DNN interpretability and explainability (Ribeiro et al., 2016; Samek et al., 2018) as well as on using formal methods to facilitate reasoning about DRL agents (Kazak et al., 2019), but these efforts are still nascent, and typically suffer from limited scalability.

In this paper, we advocate a direction of work for addressing this important gap, by integrating *modeling techniques* into the DRL training process. The idea is to leverage the strengths of classical specification approaches, which are normally applied in procedural or rule-based modeling, and carry these ad-

<sup>a</sup> <https://orcid.org/0000-0002-7951-7795>

<sup>b</sup> <https://orcid.org/0000-0002-4009-5353>

<sup>c</sup> <https://orcid.org/0000-0001-5904-5105>

vantages over to the training process of DRL agents. More specifically, we propose to sometimes override the computation of the DRL agent’s *reward function* (Sutton and Barto, 2018), which is then reflected in the *return* that the agent is being trained to maximize. By creating a connection that will allow modelers to formulate a specification in their modeling formalism of choice, and then have this specification affect the computed reward so that it reflects how well the specification is satisfied, we seek to generate a DRL agent that better conforms to the system’s requirements.

Our proposed approach is general, in the sense that numerous modeling formalisms could be integrated into the DRL process. As a proof-of-concept, and for evaluation purposes, we focus here on a particular brand of modeling schemes, collectively referred to as *scenario-based modeling (SBM)* (Damm and Harel, 2001; Harel et al., 2012b). In SBM, a modeler creates small, stand-alone scenarios, each reflecting a certain desirable or undesirable behavior of the system under development. These scenarios are fully executable, and when interwoven together bring about an executable model of the desired global system behavior. A key feature of SBM is the ability of each scenario to specify *forbidden* behavior, which the system as a whole should avoid. SBM has been shown to be quite effective in modeling systems from varied domains, such as web-servers (Harel and Katz, 2014), cache coherence protocols (Katz et al., 2015), games (Harel et al., 2011), production control (Harel et al., 2005), biological systems (Kugler et al., 2008), transportation (Greenyer et al., 2016a) and others.

During DRL training, the agent may be regarded as a reactive system: it receives an input from the environment, reacts, observes the computed reward that its actions have produced, and optionally adjusts its behavior for the future. In order to integrate SBM and DRL, we suggest to execute the scenario-based model in parallel to the DRL agent’s training. Then, whenever the agent performs an action that the SBM model forbids, we propose to reflect this by penalizing the agent, through its reward values. We argue that this process would increase the likelihood that the agent learns, in addition to its original goals, the constraints expressed through the scenario-based model.

As a case study, we chose to focus on the Aurora system (Jay et al., 2019), which is a DRL-based Internet congestion control algorithm. The algorithm is deployed at the sender node of a communication system, and controls the sending rate of that node, with the goal of optimizing the communication system’s throughput. As part of our evaluation, we demonstrate how a scenario-based specification can be used to in-

fluence the training of the Aurora DRL agent, in order to improve its *fairness* by preventing it from repeatedly increasing its sending rate at the possible expense of other senders on the same link. Our experiments show that an agent trained this way is far less likely to exhibit the unwanted behavior, when compared to an agent trained by DRL alone, thus highlighting the potential of our approach.

The rest of the paper is organized as follows. In Section 2 we provide background on scenario-based modeling and deep reinforcement learning. In Sections 3 and 4 we describe the integration between SBM and DRL, first conceptually and then technically. In Section 5 we describe our case study. We follow with a discussion of related work in Section 6 and conclude in Section 7.

## 2 BACKGROUND

### 2.1 Scenario-based Modeling

Scenario-based modeling (SBM) (Damm and Harel, 2001; Harel and Marelly, 2003; Harel et al., 2012b) is a modeling paradigm, designed to facilitate the development of reactive systems from components that are aligned with the way humans perceive and describe system behavior. The focus of SBM is on inter-object, system-wide behaviors, thus differing from the more conventional, object-centric paradigms. In SBM, a system is comprised of components called *scenario objects*, each of which describes a single desired or undesired behavior of the system. This behavior is formalized as a sequence of events. A scenario-based model is fully executable: when it runs, all its scenario objects are composed and run in parallel, in a synchronized fashion, resulting in cohesive system behavior. The resulting model thus complies with the requirements and constraints of each of the participating scenario objects (Harel and Marelly, 2003; Harel et al., 2012b).

More concretely, each scenario object in a scenario-based model can be regarded as a transition system, whose states are referred to as *synchronization points*. The scenario object transitions between synchronization points according to the triggering of *events* by a global *event selection mechanism*. At each synchronization point, the scenario object affects the triggering of the next event by declaring events that it *requests* and events that it *blocks*. These declarations encode, respectively, desirable and forbidden actions, as seen from the perspective of that particular scenario object. Scenario objects can also declare events that they passively *wait-for*, thus asking to be notified

when these occur. After making its event declarations, the scenario object is suspended until an event that it requested or waited-for is triggered, at which point the scenario resumes and may transition to another synchronization point.

At execution time, all scenario objects are run in parallel, until they all reach a synchronization point. Then, the declarations of all requested and blocked events are collected by the event selection mechanism, which first selects, and then triggers one of the events that is requested by at least one scenario object and is blocked by none.

Fig. 1 (borrowed from (Harel et al., 2012a)) depicts a scenario-based model of a simple system for controlling the temperature and fluid level in a water tank. Each scenario object is depicted as a transition system, in which the nodes represent synchronization points. The transition edges are associated with the requested or waited-for events in the preceding node. The scenarios ADDHOTWATER and ADDCOLDWATER repeatedly wait for WATERLOW events and then request three times the event ADDHOT or ADDCOLD, respectively. Since, by default, these six events may be triggered in any order, a new scenario STABILITY is introduced, with the intent of keeping the temperature more stable. It enforces the interleaving of ADDHOT and ADDCOLD events by alternately blocking them. The resulting execution trace is depicted in the event log.

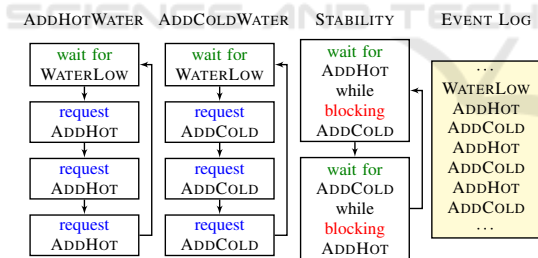


Figure 1: (Borrowed from (Harel et al., 2012a)) A scenario-based model for controlling water temperature and level.

Selecting the next event to be triggered, from among all events that are requested and not blocked, can be customized to fit the system at hand. Common policies include arbitrary selection, randomized selection, a selection based on predefined priorities, or a selection based on look-ahead for achieving certain outcomes (Harel et al., 2002).

For our present purposes, it is convenient to think of scenario objects in terms of transition systems. In practice, SBM is supported in a variety of textual and visual frameworks. Notable examples include the language of *live sequence charts* (LSC), where SBM concepts are applied to produce modal sequence dia-

grams (Damm and Harel, 2001; Harel and Marelly, 2003); implementations in various high level languages, such as Java, C, C++, Erlang, JavaScript, and Python (see, e.g., (Harel et al., 2012b)); and various domain specific languages (Greenyer et al., 2016a) and extensions (Harel et al., 2020; Katz et al., 2019).

A particularly useful trait of SBM is that the resulting models are amenable to model checking, and facilitate compositional verification (Harel et al., 2011; Harel et al., 2013b; Harel et al., 2015c; Katz et al., 2015; Katz, 2013; Harel et al., 2015b). Thus, it is often possible to apply formal verification to ensure that a scenario-based model satisfies various criteria, either as a stand-alone model or as a component within a larger system. Automated analysis techniques can also be used to execute scenario-based models in distributed architectures (Harel et al., 2015a; Steinberg et al., 2018; Steinberg et al., 2017; Greenyer et al., 2016b; Harel et al., 2013a), to automatically repair these models (Harel et al., 2014; Harel et al., 2012a; Katz, 2021b), and to augment them in various ways, e.g., as part of the Wise Computing initiative (Harel et al., 2018; Marron et al., 2016; Harel et al., 2016a; Harel et al., 2016b).

For our work here, namely the injection of domain-specific knowledge into the DRL training procedure, SBM is an attractive choice, as it is formal, executable, and facilitates incremental development (Gordon et al., 2012; Alexandron et al., 2014). Furthermore, its natural alignment with how experts may describe the specification of the system at hand helps in transparently highlighting important parts of the training procedure. Indeed, using SBM to complement DRL was demonstrated in the past, although the focus so far has been on *guarding* an already-trained DRL agent, rather than affecting what the agent actually learns (Katz, 2020; Katz and Elyasaf, 2021).

## 2.2 Deep Reinforcement Learning

Deep reinforcement learning (Sutton and Barto, 2018) is a method for automatically producing a decision-making agent, whose goal during training is to achieve a high *return* (according to some function) through interactions with its environment.

Fig. 2 depicts the basic DRL learning cycle. The agent and its environment interact at discrete time steps  $t \in \{0, 1, 2, 3, \dots\}$ . At each time step  $t$ , the agent observes the environment's state  $s_t$ , and selects its action  $a_t$  accordingly. In the subsequent time step  $t + 1$ , and as a result of its action  $a_t$  at time  $t$ , the agent receives its reward  $R_t = R(s_t, a_t)$ , the environment moves to state  $s_{t+1}$ , and the process repeats. Through this interaction, the agent gradually learns a

policy function  $f: s_t \rightarrow a_t$  that maximizes its *return*  $G_t$ , the future cumulative discounted reward:

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

where  $R_t$  is the reward at time  $t$ , and  $\gamma$  is a *discount rate* parameter,  $0 \leq \gamma \leq 1$ .

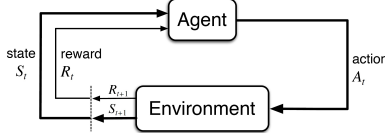


Figure 2: (Borrowed from (Sutton and Barto, 2018)) The agent-environment interaction in reinforcement learning.

It is commonly accepted that specifying an appropriate reward function  $R_t$  (as function of the action  $a_t$  in state  $s_t$ ) is crucial to the success of the DRL training process. Consequently, this topic has received significant attention (Ng et al., 1999; Zou et al., 2019; Sutton and Barto, 2018). As we later explain, the approach that we advocate here is complementary to this line of research: we propose to augment the reward function with constraints and specifications provided by domain experts.

### 3 INTEGRATING SBM INTO THE REWARD FUNCTION

Our proposed approach is to integrate a scenario-based model into the DRL training loop, in order to instruct the agent being trained to follow the constraints and specifications embodied in those scenarios. Specifically, we create a one-to-one mapping between the DRL agent's possible actions and a dedicated subset of the events in the scenario-based model, so that the scenario objects may react to the agent's actions.

We execute the scenario-based model alongside the agent under training, and, if at time step  $t$  we denote the model's state  $\tilde{s}_t$ , the agent's reward function  $R_t$  is computed as follows: (i) at time step  $t$ , the agent selects an action  $a_t$ ; (ii) the environment reacts to  $a_t$ , transitions to a new state  $s_{t+1}$ , and computes a candidate reward value  $\tilde{R}_t$ ; (iii) the scenario-based model also receives  $a_t$ , and transitions to a new state  $\tilde{s}_{t+1}$ ; (iv) if  $a_t$  is blocked in state  $\tilde{s}_t$ , the scenario-based model *penalizes* the agent by decreasing the reward:

$$R_t = \begin{cases} \alpha \cdot \tilde{R}_t - \Delta & ; \text{if } \tilde{s}_t \xrightarrow{a_t} \tilde{s}_{t+1} \text{ is blocked} \\ \tilde{R}_t & ; \text{otherwise} \end{cases} \quad (1)$$

for some constants  $\alpha \in [-1, 1]$  and  $\Delta \geq 0$ ; and (v)  $R_t$  is returned to the agent as the reward value at this step. The new training process is illustrated in Fig. 3. The motivation for these changes in the training process is to allow the DRL agent to learn a policy that satisfies the requirements encoded in the original reward function, while at the same time learning to satisfy the specifications encoded as scenarios. This is done without changing the interface between the agent and its learning environment. Observe that the scenario-based model is currently aware of its own internal state and of the actions selected by the DRL agent; but is unaware of the environment state  $s_t$ . Allowing the scenario objects to view also the environment state, which may be required in more complex systems, is left for future work.

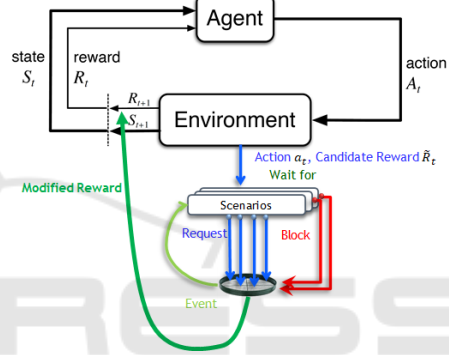


Figure 3: The agent-environment interaction in DRL integrated with SBM: At each time-step, the environment calculates a candidate reward  $\tilde{R}_t$  based on its state  $s_t$ , and the agent action  $a_t$ . The scenario-based model is executed in parallel, and may reduce  $\tilde{R}_t$  if  $a_t$  is a forbidden action in its state  $\tilde{s}_t$ .

Selecting an appropriate penalty policy, i.e. selecting the  $\alpha$  and  $\Delta$  constants, can of course have a significant impact on the learned policy. Consideration should be given to the balance between allowing the agent to learn a policy that solves the original problem, and encouraging it to follow the SBM specifications.

Let us consider again the “hot/cold” example of Fig. 1, this time from a DRL perspective. Suppose we want to train a DRL agent whose goal is to keep the water temperature steady, by mixing hot and cold water. There are many strategies that would achieve this goal; but suppose we also wish to introduce a constraint that the agent should avoid two consecutive additions of hot water or of cold water. We can achieve this by creating a simple scenario-based model comprised of the *Stability* scenario (Fig. 1), and integrating it into the training process. This scenario object would then penalize the DRL agent whenever it performs the undesirable sequence of actions, pushing it



towards learning the desired policy.

## 4 PROOF-OF-CONCEPT IMPLEMENTATION

We have created a simple proof-of-concept implementation of our approach, by extending the AI-Gym Python framework for training DRL agents (Brockman et al., 2016). Specifically, we connected the AI-Gym environment object (*env*) to the BP-Py framework for specifying scenario-based models in Python (Yaacov, 2020), in a way that allows scenario objects to affect the computation of the reward function by the *env*, as previously described.

The main connection point is AI-Gym *env*'s *step* method, which is invoked in every iteration of the DRL training process, and which eventually computes the reward value. We altered the method so that it communicates with the SBM core, to inform the scenario objects of the agent's selected action, and in turn to receive instructions on how to modify the reward value, if needed.

Recall that normally, a scenario-based model will, in each iteration: (i) trigger an event that is requested and not blocked; and (ii) wake up all scenario objects that requested or waited-for the triggered event, allowing them to react by transitioning to a new state and updating their respective declarations of requested, blocked and waited-for events. This high-level loop appears in Fig. 4. The execution terminates once there are no enabled events, e.g., events that are requested and not blocked.

---

```
# Main loop
while True:
    if noEnabledEvents():
        terminateExecution()

    event = selectEnabledEvent()
    advanceAllBThreads(event)
```

---

Figure 4: Pseudo code of the main execution loop in a regular scenario-based model. The execution will terminate once there are no enabled events.

In order to support integration with AI-Gym, and allow the scenario-based model to execute in parallel to the training of the DRL agent, we modified the execution scheme to run in *super steps* (Harel et al., 2002): the SBM model runs until it has no additional enabled events, and then, instead of terminating, it waits for a new action event to be selected by the DRL agent. Once such an event is triggered by the agent,

it is processed by the scenario objects (like any other triggered event). However, if the event triggered by the agent happens to be blocked by the scenarios, information is passed back to the AI-Gym *env* to penalize the agent's reward value. The scenario objects then carry out another super step, and the process repeats. Fig. 5 shows a pseudo code of the modified, high-level loop.

---

```
# Modified main loop, with super steps
while True:
    # Perform super step
    while haveEnabledEvents():
        event = selectEnabledEvent()
        advanceAllBThreads(event)

    # Handle an agent action
    action = waitForAgentAction()
    if isBlocked(action):
        penalizeAgentReward()
    else:
        keepOriginalAgentReward()
    advanceAllBThreads(action)
```

---

Figure 5: Pseudo code of the main loop of the scenario-based execution integrated with the DRL training. After each super step, the model waits for an action from the agent, and penalizes the agent if that action was blocked.

## 5 CASE STUDY: THE AURORA CONGESTION CONTROLLER

As a case study, we chose to focus on the Aurora system (Jay et al., 2019) — a DRL-based Internet congestion control algorithm. The algorithm is deployed at the sender node of a communication system, and controls the sending rate of that node, with the goal of optimizing the communication system's throughput. The selection of sending rate is based on various parameters, such as previously observed throughput, the link's latency, and the percent of previously lost packets. Aurora is intended to replace earlier hand-crafted algorithms for obtaining similar goals, and was shown to achieve excellent performance (Jay et al., 2019).

The authors of Aurora raised an interesting point regarding its fairness, asking: *Can our RL agent be trained to "play well" with other protocols (TCP, PCC, BBR, Copa)?* (Jay et al., 2019). Indeed, in the case of Aurora, and more generally in DRL, it is often hard to train the agent to comply with various fairness properties. In our case study, we set out to add specific fairness constraints to the Aurora agent, using SBM. Specifically, we attempted to avoid the situa-

tion where the algorithm becomes a “bandwidth hog” — i.e., to train it not to increase its sending rate continuously, thus providing other senders on the same link with a fair share of the bandwidth.

## 5.1 Evaluation Setup

Using the BP-Py environment (Yaacov, 2020), we created a simple scenario-based model. This model, comprised of a single scenario, is designed to penalize the Aurora agent for  $k$  consecutive increases in sending rate. The SBM code for  $k = 3$  appears in Fig. 6.

```
def SBP_avoid_k_in_a_row():
    k = 3
    counter = 0
    blockedEvList = []
    waitForEvList = [BEvent("IncreaseRate"),
                     BEvent("DecreaseRate"),
                     BEvent("KeepRate")]

    while True:
        lastEv = yield{ waitFor:waitForEvList,
                       block:blockedEvList }
        if lastEv != None:
            if lastEv == BEvent("DecreaseRate")
            or lastEv == BEvent("KeepRate"):
                counter = 0
                blockedEvList = []
            else:
                if counter == k - 1:
                    #Blocking!
                    blockedEvList.append(
                        BEvent("IncreaseRate"))
                else:
                    counter += 1
```

Figure 6: The Python implementation of a scenario that blocks the *IncreaseRate* event after  $k - 1$  consecutive occurrences.

The scenario waits for three possible events that represent the three different actions of the Aurora agent: *IncreaseRate*, *DecreaseRate*, and *KeepRate*, representing the agent’s decision to increase the sending rate, decrease it or keep it steady, respectively. Whenever the *IncreaseRate* event is triggered  $k - 1$  times consecutively, the scenario will block it, until a different action is selected by the agent. Once the execution environment detects that a requested agent action maps to a blocked event, it will override the reward with a penalty, thus signaling to the agent that it is not a desirable behavior.

For training the Aurora agent, we used the original training framework provided in (Jay et al., 2019). We compared an agent trained using the original

framework,  $\mathcal{A}_O$ , to an agent trained using our SBM-enhanced framework,  $\mathcal{A}_E$ . All Aurora-related configurable parameters were identical between the two agents.

For computing the penalty when training  $\mathcal{A}_E$ , we empirically selected the penalty function parameters (as described by equation 1) to be  $\alpha = 0, \Delta = -4.5$ . These parameters allowed the agent to effectively learn both the main goals and the additional constraints specified by the scenario-based model. Values of  $\Delta$  in the range  $(-2, 0]$  proved to be too small (the agent failed to address the SBM constraints), whereas values in the range  $(-\infty, -10]$  resulted in a disruption to the process of the agent learning its main goals.

## 5.2 Evaluation Results

We begin by comparing the two agents employing the evaluation metrics used by (Jay et al., 2019). Fig. 7 shows the training performance of both  $\mathcal{A}_O$  and  $\mathcal{A}_E$  as a function of the time-step (in log scale).

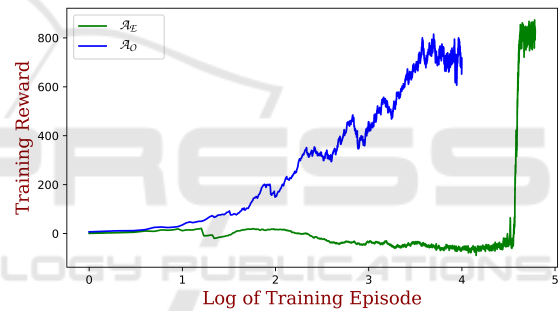


Figure 7: The average training reward obtained by  $\mathcal{A}_O$  and  $\mathcal{A}_E$ , as a function of time (log scaled).

These results reveal a significant difference in the training times of the two agents. Specifically, it takes  $\mathcal{A}_E$  significantly more epochs to learn an “adequate” policy, i.e. to reach a similar reward level to that obtained by  $\mathcal{A}_O$ . We observe that  $\mathcal{A}_E$  converges to a good policy after about 40,000 epochs (a little after 4.5 in the logarithmic time-step scale), compared to about 3000 epochs of  $\mathcal{A}_O$  (around 3.5 in the logarithmic time-step scale).

Next, we compare the frequency of the two agents choosing to increase their sending rate three consecutive times or more (performing a “violation”). For  $\mathcal{A}_O$ , the average frequency of such violations is 9%-11%, as can be seen in Fig. 8.

In contrast, for  $\mathcal{A}_E$ , the average frequency of performing a violation is about 0.34%, as can be seen in Fig. 9.

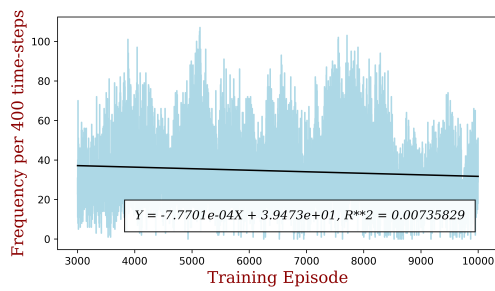


Figure 8: Each line shows the number of times  $\mathcal{A}_O$  chose to increase the sending rate three consecutive times or more, during the training process. Each training episode is comprised of 400 time steps, and the agent performed on average 35-45 violations per episode, which translates to a frequency of about 9%-11%. The resulting linear regression line is  $y = -0.00077701x + 39.47343569$ .

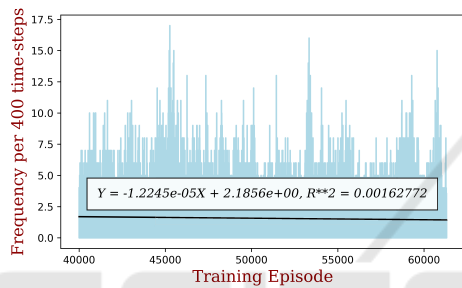


Figure 9: Each line shows the number of times  $\mathcal{A}_E$  chose to increase the sending rate three consecutive times or more, during the training process. This time the agent performed only 1-2 violations on average per training episode, which translates to a frequency of about 0.34%. The resulting linear regression line is  $y = -0.000012245x + 2.1856$ .

**Summary.** The results above demonstrate a highly significant change in behavior between  $\mathcal{A}_O$  and  $\mathcal{A}_E$  when it comes to the frequency of performing a violation: whereas  $\mathcal{A}_O$  would perform a violation about 9-11% of the time, for  $\mathcal{A}_E$  this rate drops to 0.34%. Both agents, however, achieve a similar overall reward level, indicating that they both learned an adequate policy with respect to the main goal’s of the system. While the enhanced agent took longer to converge to this policy (as it had to learn additional constraints), these results showcase the feasibility and potential of our approach.

## 6 RELATED WORK

Several approaches have been proposed in recent years for using hand-crafted software components to enhance the run-time functionality and performance of DNNs, or to improve their training process. One notable family of approaches calls for *composing*

DNNs and hand-crafted components. See for example (Shalev-Shwartz et al., 2017) where the decisions of an autonomous driving systems can be overridden by rules. This composition can be parallel, where a DNN and hand-crafted code run side by side, each handling different tasks; sequential, where a DNN’s output feeds as input into hand-crafted code, or vice-versa; or ensemble-based, where DNNs and hand-crafted code attempt to solve the same problem and agree on an output.

Another notable family of approaches contains those that are *reflection-based*, where both hand-crafted code and DNN gradually adjust themselves according to their past performance (Kang et al., 2017; Milan et al., 2017; Ray and Chakrabarti, 2020). The approach we propose here can be viewed as a bridge between composition-based and reflection-based approaches: a hand-crafted model is run alongside a DRL agent, with the purpose of improving the latter’s training process.

Prior work has explored the potential synergies between SBM and DRL. In one attempt, Elyasaf et al. (Elyasaf et al., 2019) used DRL to fine-tune the execution strategy of an existing scenario-based model. Using a game of RoboSoccer as a case-study, they demonstrated how the DRL agent could learn to guide a scenario-based player to more effectively grab the soccer ball. In a separate attempt, Katz (Katz, 2020; Katz, 2021a) focused on using SBM models to *guard* an existing DRL agent — that is, to override decisions made by the DRL agent that violate the scenario-based model. The technique that we propose here is different from and complementary to both of these approaches: instead of using DRL to guide a hand-crafted model or using SBM to guard an existing DRL agent, we propose to use SBM to improve the DRL agent, a-priori, so that it better abides by a scenario-based specification.

## 7 CONCLUSION AND FUTURE WORK

Deep reinforcement learning is an excellent tool for addressing many real-world problems; but it is lacking, in the sense that it does not naturally lend itself to the integration of expert knowledge. Through our proposed approach, we seek to bridge this gap, and allow the integration of classical modeling techniques into the DRL training loop. The resulting agents, as we have demonstrated, are more likely to adhere to the policies, goals and restrictions defined by the domain experts. Apart from improving performance, this approach increases the transparency and explain-

ability of DRL agents, and can be regarded as documenting and explaining these agents' behavior.

Turning to the future, we intend to apply the technique to more complex case-studies, involving more intricate agents and more elaborate scenario-based specifications — and use them to also study the scalability of the approach. Another angle we intend to pursue is to enhance the DRL-SBM interface, either by changing the SBM semantics or by defining an event-based protocol so that it can allow manipulating the agent's reward function in more subtle ways. Thus, the SBM feedback will be able to distinguish actions that are slightly undesirable from those that are extremely undesirable. We also intend to explore adding constructs for encouraging an agent to take desirable actions, in addition to penalizing it for taking undesirable ones. As discussed earlier, we are also interested in exploring generalizable criteria for choosing penalty values that will be conducive to learning the desired properties, while preserving the overall learning of the task at hand.

Another angle for future research is to measure additional effects that scenario-assisted training may have on DRL, such as accelerating the learning of properties that could, in principle, be learned without such assistance. It is also interesting to see on a variety of case studies if SBM-specified expert advice that is aimed at improving system performance (as opposed to complying with new requirements) indeed accomplishes such improvement over what the system could learn on its own.

## ACKNOWLEDGEMENTS

The work of R. Yerushalmi, G. Amir, A. Elyasaf and G. Katz was partially supported by a grant from the Israeli Smart Transportation Research Center (ISTRC). The work of D. Harel, A. Marron and R. Yerushalmi was partially supported by a research grant from the Estate of Ursula Levy, the Midwest Electron Microscope Project, Dr. Elias and Esta Schwartz Instrumentation Fund in Memory of Uri and Atida Litauer, and the Crown Family Philanthropies.

## REFERENCES

- Alexandron, G., Armoni, M., Gordon, M., and Harel, D. (2014). Scenario-Based Programming: Reducing the Cognitive Load, Fostering Abstract Thinking. In *Proc. 36th Int. Conf. on Software Engineering (ICSE)*, pages 311–320.
- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W. (2016). OpenAI Gym. Technical Report. <http://arxiv.org/abs/1606.01540>.
- Damm, W. and Harel, D. (2001). LSCs: Breathing Life into Message Sequence Charts. *Journal on Formal Methods in System Design (FMSD)*, 19(1):45–80.
- Eliyahu, T., Kazak, Y., Katz, G., and Schapira, M. (2021). Verifying Learning-Augmented Systems. In *Proc. Conf. of the ACM Special Interest Group on Data Communication on the Applications, Technologies, Architectures, and Protocols for Computer Communication (SIGCOMM)*, pages 305–318.
- Elyasaf, A. (2021). Context-Oriented Behavioral Programming. *Information and Software Technology*, 133.
- Elyasaf, A., Sadon, A., Weiss, G., and Yaacov, T. (2019). Using Behavioural Programming with Solver, Context, and Deep Reinforcement Learning for Playing a Simplified RoboCup-Type Game. In *Proc. 22nd ACM/IEEE Int. Conf. on Model Driven Engineering Languages and Systems Companion (MODELS-C)*, pages 243–251.
- Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep Learning*. MIT Press.
- Gordon, M., Marron, A., and Meerbaum-Salant, O. (2012). Spaghetti for the Main Course? Observations on the Naturalness of Scenario-Based Programming. In *Proc. 17th ACM Annual Conf. on Innovation and Technology in Computer Science Education (ITCSE)*, pages 198–203.
- Greenyer, J., Gritzner, D., Katz, G., and Marron, A. (2016a). Scenario-Based Modeling and Synthesis for Reactive Systems with Dynamic System Structure in ScenarioTools. In *Proc. 19th ACM/IEEE Int. Conf. on Model Driven Engineering Languages and Systems (MODELS)*, pages 16–23.
- Greenyer, J., Gritzner, D., Katz, G., Marron, A., Glade, N., Gutjahr, T., and König, F. (2016b). Distributed Execution of Scenario-Based Specifications of Structurally Dynamic Cyber-Physical Systems. In *Proc. 3rd Int. Conf. on System-Integrated Intelligence: New Challenges for Product and Production Engineering (SYSINT)*, pages 552–559.
- Harel, D., Kantor, A., and Katz, G. (2013a). Relaxing Synchronization Constraints in Behavioral Programs. In *Proc. 19th Int. Conf. on Logic for Programming, Artificial Intelligence and Reasoning (LPAR)*, pages 355–372.
- Harel, D., Kantor, A., Katz, G., Marron, A., Mizrahi, L., and Weiss, G. (2013b). On Composing and Proving the Correctness of Reactive Behavior. In *Proc. 13th Int. Conf. on Embedded Software (EMSOFT)*, pages 1–10.
- Harel, D., Kantor, A., Katz, G., Marron, A., Weiss, G., and Wiener, G. (2015a). Towards Behavioral Programming in Distributed Architectures. *Journal of Science of Computer Programming (J. SCP)*, 98:233–267.
- Harel, D. and Katz, G. (2014). Scaling-Up Behavioral Programming: Steps from Basic Principles to Application Architectures. In *Proc. 4th SPLASH Workshop on Programming based on Actors, Agents and Decentralized Control (AGERE!)*, pages 95–108.



- Harel, D., Katz, G., Lampert, R., Marron, A., and Weiss, G. (2015b). On the Succinctness of Idioms for Concurrent Programming. In *Proc. 26th Int. Conf. on Concurrency Theory (CONCUR)*, pages 85–99.
- Harel, D., Katz, G., Marelly, R., and Marron, A. (2016a). An Initial Wise Development Environment for Behavioral Models. In *Proc. 4th Int. Conf. on Model-Driven Engineering and Software Development (MODELSWARD)*, pages 600–612.
- Harel, D., Katz, G., Marelly, R., and Marron, A. (2016b). First Steps Towards a Wise Development Environment for Behavioral Models. *Int. Journal of Information System Modeling and Design (IJISMD)*, 7(3):1–22.
- Harel, D., Katz, G., Marelly, R., and Marron, A. (2018). Wise Computing: Toward Endowing System Development with Proactive Wisdom. *IEEE Computer*, 51(2):14–26.
- Harel, D., Katz, G., Marron, A., Sadon, A., and Weiss, G. (2020). Executing Scenario-Based Specification with Dynamic Generation of Rich Events. *Communications in Computer and Information Science (CCIS)*, 1161.
- Harel, D., Katz, G., Marron, A., and Weiss, G. (2012a). Non-Intrusive Repair of Reactive Programs. In *Proc. 17th IEEE Int. Conf. on Engineering of Complex Computer Systems (ICECCS)*, pages 3–12.
- Harel, D., Katz, G., Marron, A., and Weiss, G. (2014). Non-Intrusive Repair of Safety and Liveness Violations in Reactive Programs. *Transactions on Computational Collective Intelligence (TCCI)*, 16:1–33.
- Harel, D., Katz, G., Marron, A., and Weiss, G. (2015c). The Effect of Concurrent Programming Idioms on Verification. In *Proc. 3rd Int. Conf. on Model-Driven Engineering and Software Development (MODELSWARD)*, pages 363–369.
- Harel, D., Kugler, H., Marelly, R., and Pnueli, A. (2002). Smart Play-Out of Behavioral Requirements. In *Proc. 4th Int. Conf. on Formal Methods in Computer-Aided Design (FMCAD)*, pages 378–398.
- Harel, D., Kugler, H., and Weiss, G. (2005). Some Methodological Observations resulting from Experience using LSCs and the Play-In/Play-Out Approach. In *Scenarios: Models, Transformations and Tools*, pages 26–42. Springer.
- Harel, D., Lampert, R., Marron, A., and Weiss, G. (2011). Model-checking behavioral programs. In *Proc. 9th ACM Int. Conf. on Embedded Software (EMSOFT)*, pages 279–288.
- Harel, D. and Marelly, R. (2003). *Come, Let's Play: Scenario-Based Programming using LSCs and the Play-Engine*, volume 1. Springer Science & Business Media.
- Harel, D., Marron, A., and Weiss, G. (2012b). Behavioral Programming. *Communications of the ACM (CACM)*, 55(7):90–100.
- Jay, N., Rotman, N., Godfrey, B., Schapira, M., and Tamar, A. (2019). A Deep Reinforcement Learning Perspective on Internet Congestion Control. In *Proc. 36th Int. Conf. on Machine Learning (ICML)*, pages 3050–3059.
- Kang, C., Kim, G., and Yoo, S.-I. (2017). Detection and Recognition of Text Embedded in Online Images via Neural Context Models. In *Proc. 31st AAAI Conf. on Artificial Intelligence (AAAI)*.
- Katz, G. (2013). On Module-Based Abstraction and Repair of Behavioral Programs. In *Proc. 19th Int. Conf. on Logic for Programming, Artificial Intelligence and Reasoning (LPAR)*, pages 518–535.
- Katz, G. (2020). Guarded Deep Learning using Scenario-Based Modeling. In *Proc. 8th Int. Conf. on Model-Driven Engineering and Software Development (MODELSWARD)*, pages 126–136.
- Katz, G. (2021a). Augmenting Deep Neural Networks with Scenario-Based Guard Rules. *Communications in Computer and Information Science (CCIS)*, 1361:147–172.
- Katz, G. (2021b). Towards Repairing Scenario-Based Models with Rich Events. In *Proc. 9th Int. Conf. on Model-Driven Engineering and Software Development (MODELSWARD)*, pages 362–372.
- Katz, G., Barrett, C., and Harel, D. (2015). Theory-Aided Model Checking of Concurrent Transition Systems. In *Proc. 15th Int. Conf. on Formal Methods in Computer-Aided Design (FMCAD)*, pages 81–88.
- Katz, G. and Elyasaf, A. (2021). Towards Combining Deep Learning, Verification, and Scenario-Based Programming. In *Proc. 1st Workshop on Verification of Autonomous and Robotic Systems (VARS)*, pages 1–3.
- Katz, G., Marron, A., Sadon, A., and Weiss, G. (2019). On-the-Fly Construction of Composite Events in Scenario-Based Modeling Using Constraint Solvers. In *Proc. 7th Int. Conf. on Model-Driven Engineering and Software Development (MODELSWARD)*, pages 143–156.
- Kazak, Y., Barrett, C., Katz, G., and Schapira, M. (2019). Verifying Deep-RL-Driven Systems. In *Proc. 1st ACM SIGCOMM Workshop on Network Meets AI & ML (NetAI)*, pages 83–89.
- Kiran, B., Sobh, I., Talpaert, V., Mannion, P., Sallab, A., Yogamani, S., and Perez, P. (2021). Deep Reinforcement Learning for Autonomous Driving: A Survey. *IEEE Transactions on Intelligent Transportation Systems*, pages 1–18.
- Kugler, H., Marelly, R., Appleby, L., Fisher, J., Pnueli, A., Harel, D., Stern, M., Hubbard, J., et al. (2008). A Scenario-Based Approach to Modeling Development: A Prototype Model of *C. Elegans* Vulval Fate Specification. *Developmental biology*, 323(1):1–5.
- Li, J., Pang, D., Zheng, Y., Guan, X., and Le, X. (2022). A Flexible Manufacturing Assembly System with Deep Reinforcement Learning. *Control Engineering Practice*, 118.
- Marron, A., Arnon, B., Elyasaf, A., Gordon, M., Katz, G., Lapid, H., Marelly, R., Sherman, D., Szekely, S., Weiss, G., and Harel, D. (2016). Six (Im)possible Things before Breakfast: Building-Blocks and Design-Principles for Wise Computing. In

- Proc. 19th ACM/IEEE Int. Conf. on Model Driven Engineering Languages and Systems (MODELS)*, pages 94–100.
- Milan, A., Rezafofighi, H., Dick, A., Reid, I., and Schindler, K. (2017). Online Multi-Target Tracking using Recurrent Neural Networks. In *Proc. 31st AAAI Conf. on Artificial Intelligence (AAAI)*.
- Mohamad Suhaili, S., Salim, N., and Jambli, M. (2021). Service Chatbots: A Systematic Review. *Expert Systems with Applications*, 184:115461.
- Ng, A., Harada, D., and Russell, S. (1999). Policy Invariance Under Reward Transformations: Theory and Application to Reward Shaping. In *Proc. 16th Int. Conf. on Machine Learning (ICML)*, pages 278–287.
- Ray, P. and Chakrabarti, A. (2020). A Mixed Approach of Deep Learning Method and Rule-Based Method to Improve Aspect Level Sentiment Analysis. *Applied Computing and Informatics*.
- Ribeiro, M., Singh, S., and Guestrin, C. (2016). Why Should I Trust You?: Explaining the Predictions of any Classifier. In *Proc. 22nd ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining*, pages 1135–1144.
- Samek, W., Wiegand, T., and Müller, K. (2018). Explainable Artificial Intelligence: Understanding, Visualizing and Interpreting Deep Learning Models. *ITU Journal: The Impact of Artificial Intelligence (AI) on Communication Networks and Services*, 1(1):39–48.
- Shalev-Shwartz, S., Shammah, S., and Shashua, A. (2017). On a Formal Model of Safe and Scalable Self-Driving Cars. Technical Report. <http://arxiv.org/abs/1708.06374>.
- Steinberg, S., Greenyer, J., Gritzner, D., Harel, D., Katz, G., and Marron, A. (2017). Distributing Scenario-Based Models: A Replicate-and-Project Approach. In *Proc. 5th Int. Conf. on Model-Driven Engineering and Software Development (MODELSWARD)*, pages 182–195.
- Steinberg, S., Greenyer, J., Gritzner, D., Harel, D., Katz, G., and Marron, A. (2018). Efficient Distributed Execution of Multi-Component Scenario-Based Models. *Communications in Computer and Information Science (CCIS)*, 880:449–483.
- Sutton, R. and Barto, A. (2018). *Introduction to Reinforcement Learning*. MIT press Cambridge.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. (2013). Intriguing Properties of Neural Networks. Technical Report. <http://arxiv.org/abs/1312.6199>.
- Xia, Z., Xue, S., Wu, J., Chen, Y., Chen, J., and Wu, L. (2021). Deep Reinforcement Learning for Smart City Communication Networks. *IEEE Transactions on Industrial Informatics*, 17(6):4188–4196.
- Yaacov, T. (2020). BPPy: Behavioral Programming in Python. <https://github.com/bThink-BGU/BPPy>.
- Ye, D., Liu, Z., Sun, M., Shi, B., Zhao, P., Wu, H., Yu, H., Yang, S., Wu, X., Guo, Q., Chen, Q., Yin, Y., Zhang, H., Shi, T., Wang, L., Fu, Q., Yang, W., and Huang, L. (2020). Mastering Complex Control in MOBA Games with Deep Reinforcement Learning. In *Proc. 34th AAAI Conf. on Artificial Intelligence (AAAI)*, pages 6672–6679.
- Zou, H., Ren, T., Yan, D., Su, H., and Zhu, J. (2019). Reward Shaping via Meta-Learning. Technical Report. <http://arxiv.org/abs/1901.09330>.