

Falsification of Cyber-Physical Systems Using Deep Reinforcement Learning

Takumi Akazaki^{1,2}, Shuang Liu³, Yoriyuki Yamagata⁴, Yihai Duan³, and Jianye Hao³

¹ The University of Tokyo

² Japan Society for the Promotion of Science

³ School of Software, Tianjin University

⁴ National Institute of Advanced Industrial Science and Technology (AIST)

Abstract. With the rapid development of software and distributed computing, *Cyber-Physical Systems* (CPS) are widely adopted in many application areas, e.g., smart grid, autonomous automobile. It is difficult to detect defects in CPS models due to the complexities involved in the software and physical systems. To find defects in CPS models efficiently, robustness guided falsification of CPS is introduced. Existing methods use several optimization techniques to generate counterexamples, which falsify the given properties of a CPS. However those methods may require a large number of simulation runs to find the counterexample and is far from practical. In this work, we explore state-of-the-art *Deep Reinforcement Learning* (DRL) techniques to reduce the number of simulation runs required to find such counterexamples. We report our method and the preliminary evaluation results.

1 Introduction

Cyber-Physical Systems (CPS) are more and more widely adopted in safety-critical domains, which makes it extremely important to guarantee the correctness of CPS systems. Testing and verification on *models* of CPS are common methods to guarantee the correctness. However, it is hard for testing to achieve a high coverage; verification techniques are usually expensive and undecidable [3] due to the infinite state space of CPS models. Therefore, robustness guided falsification [2,6] method is introduced to detect defects efficiently. In robustness guided falsification, *Signal Temporal Logic* (STL) [10] formulas are usually used to specify properties which must be satisfied by a CPS model. Robustness of an STL formula, which is a numeric measure of how “robust” a property holds in the given CPS model, is defined. The state space of the CPS model is explored and a trajectory which minimizes the robustness value is identified as a good candidate for testing. In this way, robustness guided falsification aids to generate defect-leading inputs (counterexamples), which enables more efficient, yet automatic detection of defects. Although non-termination of robustness guided falsification does not mean the absence of counterexamples, it suggests the correctness of the CPS model to some extent.

Existing approaches adopt various kinds of stochastic global optimization algorithms e.g., simulated annealing [3] and cross-entropy [28], to minimize robustness.

These methods take a full trajectory (a sequence of actions) as input, and adjusting input during the simulation is not supported. As a result, a large number of simulation runs are required in the falsification process. Existing methods cannot guarantee finding a counterexample of practical CPS models in a limited time window because the simulation would then be tremendous.

In this paper, we adopt *deep reinforcement learning* (DRL) [26] algorithms to solve the problem of falsification of STL properties for CPS models. Reinforcement learning techniques can observe feedbacks from the environment, and adjust the input action immediately. In this way, we are able to converge faster towards minimum robustness value. In particular, we adopt two state-of-the-art DRL techniques, i.e., *Asynchronous Advanced Actor Critic* (A3C) and *Double Deep-Q Network* (DDQN). Our contributions are two folds: (1) we show how to transform the problem of falsifying CPS models into a reinforcement learning problem; and (2) we implement our method and conduct preliminary evaluations to show DRL technology can help reduce the number of simulation runs required to find a falsifying input for CPS models. Reducing the number of simulation runs is important because during falsification, the majority of execution time is spent for simulation runs if CPS models are complex.

Related Work There are two kinds of works, i.e., robustness guided falsification and controller synthesis, which are most related to our approach.

In robustness guided falsification methods, quantitative semantics over *Metric Interval Temporal Logic* (MITL) and its variants STL [24,17] are employed. Then the fault detection problem is translated into the numerical minimization problem. Several tools e.g., S-TaLiRo [6,20] and Breach [16] are developed to realize this approach. Moreover, various kind of numerical optimization techniques, e.g., simulated annealing [3], cross-entropy [28], and Gaussian process optimization [8,9,4,29], are studied to solve the falsification problem efficiently. All these methods optimize the whole output trajectory of a CPS by changing the whole input trajectory. As stated above, we use reinforcement learning which can observe feedbacks from a CPS and adjust the input immediately. Thus, our method can be expected to arrive the falsifying input faster.

In contrast to robustness guided falsification, controller synthesis techniques enable choosing the input signal at a certain step based on observations of output signals. There are works that synthesize the controller to enforce the *Markov decision process* to satisfy a given LTL formula [27,23,15,30,14]. The most closest related works [22,21] apply reinforcement learning techniques to enforce the small robotic system to satisfy the given LTL formula. Our work is different from those works in two aspects: (1) we falsify the properties while the control synthesis methods try to satisfy the properties; and (2) with DRL, we could employ complex non-linear functions to learn and model the environment, which is suitable to analyze the complex dynamics of CPS.

2 Preliminary

Robustness guided falsification In this paper, we employ a variant of *Signal Temporal Logic* (STL) defined in [10]. The syntax is defined in the equation (1),

$$\varphi ::= v \sim c \mid p \mid \neg\varphi \mid \varphi_1 \vee \varphi_2 \mid \varphi_1 \mathcal{U}_I \varphi_2 \mid \varphi_1 \mathcal{S}_I \varphi_2 \quad (1)$$

where v is *real* variable, c is a rational number, p is atomic formula, $\sim \in \{<, \leq\}$ and I is an interval over non-negative real numbers. If I is $[0, \infty]$, I is omitted. We also use other common abbreviations, e.g., $\Box_I \varphi \equiv \text{True} \mathcal{U}_I \varphi$ and $\Xi_I \varphi \equiv \text{True} \mathcal{S}_I \varphi$. For a given formula φ , an output signal \mathbf{x} and time t , we adopt the notation of work [10] and denote the *robustness degree* of output signal \mathbf{x} satisfying φ at time t by $\rho(\varphi, \mathbf{x}, t)$.

We also adopt the notion of *future-reach* $\text{fr}(\varphi)$ and *past-reach* $\text{pr}(\varphi)$ following [19]. Intuitively, $\text{fr}(\varphi)$ is the time in future which is required to determine the truth value of formula φ , and $\text{pr}(\varphi)$ is the time in past. For example, $\text{fr}(p) = 0$, $\text{fr}(\Box_{[0,3]} p) = 3$ and $\text{fr}(\Xi_{[0,3]} p) = 0$. Similarly, for past-reach, $\text{pr}(p) = 0$, $\text{pr}(\Box_{[0,3]} p) = 0$, $\text{pr}(\Xi_{[0,3]} p) = 3$.

In this paper, we focus on a specific class of the temporal logic formula called *life-long property*.

Definition 1 (life-long property). A life-long property is an STL formula $\psi \equiv \Box \varphi$ where $\text{fr}(\varphi)$, $\text{pr}(\varphi)$ are finite. If $\text{fr}(\varphi) = 0$, we call ψ past-dependent life-long property.

Reinforcement Learning Reinforcement learning is one of machine learning techniques in which an agent learns the structure of the environment based on observations, and maximizes the rewards by acting according to the learnt knowledge. The standard setting of a reinforcement learning problem consists of an agent and an environment. The agent observes the current state and reward from the environment, and returns the next action to the environment. The goal of reinforcement learning is for each step n , given the sequence of previous states x_0, \dots, x_{n-1} , rewards r_1, \dots, r_n and actions a_0, \dots, a_{n-1} , generate an action a_n , which maximizes expected value of the sum of rewards: $r = \sum_{k=n}^{\infty} \gamma^k r_{k+1}$, where $0 < \gamma \leq 1$ is a discount factor. Deep reinforcement learning is a reinforcement learning technique which uses a *deep neural network* for learning. In this work, we particularly adopted 2 state-of-the-art deep reinforcement learning algorithms, i.e., *Asynchronous Advantage Actor-Critic* (A3C) [25] and *Double Deep Q Network* (DDQN) [18].

3 Our Approach

3.1 Overview of our algorithm

Let us consider the falsification problem to find a counterexample of the life-long property $\psi \equiv \Box \varphi$. If the output signal is infinitely long to past and future directions, ψ is logically equivalent to a past-dependent life-long property $\Box \Xi_{[\text{fr}(\varphi), \text{fr}(\varphi)]} \varphi$. In general, the output signal is not infinitely long to some direction but using this conversion we convert all life-long properties to past-dependent life-long properties. Our evaluation in Section 4 suggests that this approximation does not adversely affect the performance. Therefore, assume ψ is a past-dependent life-long property, we generate an input signal \mathbf{u} for system \mathcal{M} , such that the corresponding output signal $\mathcal{M}(\mathbf{u})$ does not satisfy ψ .

In our algorithm, we fix the simulation time to be T_{end} and call one simulation until time T_{end} an *episode* in conformance with the reinforcement learning terminology. We fix the discretization of time to a positive real number Δ_T . The agent \mathcal{A} generates the piecewise-constant input signal $\mathbf{u} = [(0, u_0), (\Delta_T, u_1), (2\Delta_T, u_2), \dots]$ by iterating the following steps:

Algorithm 1 Falsification for $\psi = \Box\varphi$ by reinforcement learning

input: A past-dependent life-long property $\psi = \Box\varphi$, a system \mathcal{M} , an agent \mathcal{A}
output: A counterexample input signal \mathbf{u} if exists
parameters: A step time Δ_T , the end time T_{end} , the maximum number of the episode N

```

1: for numEpisode  $\leftarrow 1$  to  $N$  do
2:    $i \leftarrow 0, r \leftarrow 0, x$  be the initial (output) state of  $\mathcal{M}$ 
3:    $\mathbf{u}$  be the empty input signal sequence
4:   while  $i\Delta_T < T_{\text{end}}$  do
5:      $u \leftarrow \mathcal{A}.\text{step}(x, r), \mathbf{u} \leftarrow \text{append}(\mathbf{u}, (i\Delta_T, u))$   $\triangleright$  choose the next input by the agent
6:      $\mathbf{x} \leftarrow \mathcal{M}(\mathbf{u}), x \leftarrow \mathbf{x}((i+1)\Delta_T)$   $\triangleright$  simulate, observe the new output state
7:      $r \leftarrow \text{reward}(\mathbf{x}, \psi)$ 
8:      $i \leftarrow i + 1$   $\triangleright$  calculate the reward by following eq. (2)
9:   end while
10:  if  $\mathbf{x} \not\models \psi$  then return  $\mathbf{u}$  as a falsifying input
11:  end if
12:   $\mathcal{A}.\text{reset}(x, r)$ 
13: end for
  
```

(1) At time $i\Delta_T$ ($i = 0, 1, \dots$), the agent \mathcal{A} chooses the next input value u_i . The generated input signal is extended to $\mathbf{u} = [(0, u_0), \dots, (i\Delta_T, u_i)]$.

(2) Our algorithm obtains the corresponding output signal $\mathbf{x} = \mathcal{M}(\mathbf{u})$ by stepping forward one simulation on the model \mathcal{M} from time $i\Delta_T$ to $(i+1)\Delta_T$ with input u_i .

(3) Let $x_{i+1} = \mathbf{x}((i+1)\Delta_T)$ be the new (observed) state (i.e., output) of the system.

(4) We compute reward r_{i+1} by $\text{reward}(\varphi, \mathbf{x}, (i+1)\Delta_T)$ (defined in Section 3.2).

(5) The agent \mathcal{A} updates its action based on the new state x_{i+1} and reward r_{i+1} .

At the end of each episode, we obtain the output signal trajectory \mathbf{x} , and check whether it satisfies the property $\psi = \Box\varphi$ or not. If it is falsified, return the current input signal \mathbf{u} as a counterexample. Otherwise, we discard the current generated signal input and restart the episode from the beginning.

The complete algorithm of our approach is shown in Algorithm 1. The method call $\mathcal{A}.\text{step}(x, r)$ represents the agent \mathcal{A} push the current state reward pair (x, r) into its memory and returns the next action u (the input signal in the next step). The method call $\mathcal{A}.\text{reset}(x, r)$ notifies the agent that the current episode is completed, and returns the current state and reward. Function $\text{reward}(\mathbf{x}, \psi)$ calculates the reward based on Definition 2.

3.2 Reward definition for life-long property falsification

Our goal is to find the input signal \mathbf{u} to the system \mathcal{M} which minimizes $\rho(\psi, \mathcal{M}(\mathbf{u}), 0)$ where $\psi = \Box\varphi$ and ρ is a robustness. We determine u_0, u_1, \dots in a greedy way. Assume that u_0, \dots, u_i are determined. u_{i+1} can be determined by

$$u_{i+1} = \arg \min_{u_{i+1}} \min_{u_{i+2}, \dots} \rho(\Box\varphi, \mathcal{M}([(0, u_0), (\Delta_T, u_1), \dots]), 0) \quad (2)$$

$$\sim \arg \max_{u_{i+1}} \max_{u_{i+2}, \dots} \sum_{k=i+1}^{\infty} \{e^{-\rho(\varphi, \mathcal{M}([(0, u_0), \dots, (k\Delta_T, u_k)]), k\Delta_T)} - 1\} \quad (3)$$

The detailed derivation steps can be found in Appendix A.

In our reinforcement learning base approach, we use discounting factor $\gamma = 1$ and reward $r_i = e^{-\rho(\varphi, \mathcal{M}([(0, u_0), \dots, (i\Delta_T, u_i)]), i\Delta_T)} - 1$ to approximately compute action u_{i+1} , from $u_0, \dots, u_i, \mathcal{M}([(0, u_0), \dots, (i\Delta_T, u_i)])$ and r_1, \dots, r_i .

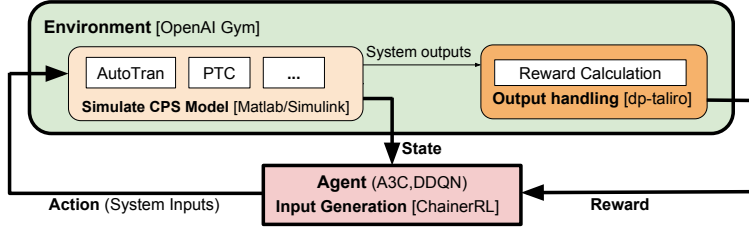


Fig. 1. Architecture of our system

Definition 2 (reward). Let $\psi \equiv \Box\varphi$ be a past-dependent formula and $\mathbf{x} = \mathcal{M}(\mathbf{u})$ be a finite length signal until the time t . We define the reward $\text{reward}(\psi, \mathbf{x})$ as

$$\text{reward}(\psi, \mathbf{x}) = \exp(-\rho(\varphi, \mathbf{x}, t)) - 1 \quad (4)$$

4 Preliminary Results

Implementation The overall architecture of our system is shown in Fig. 1. Our implementation consists of three components, i.e., input generation, output handling and simulation. The input generation component adopts reinforcement learning techniques and is implemented based on the ChainerRL library [1]. We use default hyper-parameters in the library or sample programs without change. The output handling component conducts reward calculation using dp-TaliRo [5]. The simulation is conducted with Matlab/Simulink models, which are encapsulated by the openAI gym library [11].

Evaluation Settings We use a widely adopted CPS model, automatic transmission control system (AT) [7], to evaluate our method. AT has throttle and brake as input ports, and the output ports are the vehicle velocity v , the engine rotation speed ω and the current gear state g . We conduct our evaluation with the formulas in Table 1. Formulas $\varphi_1-\varphi_6$ are rewriting of $\varphi_1^{AT}-\varphi_6^{AT}$ in benchmark [7] into life-long properties in our approach. In addition, we propose three new formulas $\varphi_7-\varphi_9$. For each formula $\varphi_1-\varphi_9$, we compare the performance of our approaches (A3C, DDQN), with the baseline algorithms, i.e., simulated annealing (SA) and cross entropy (CE). For each property, we run the falsification procedure 20 times. For each falsification procedure, we execute simulation episodes up to 200 times and measure the number of simulation episodes required to falsify the property. If the property cannot be falsified within 200 episodes, the procedure fails. We observe that Δ_T may strongly affect the performance of each algorithm. Therefore, we vary Δ_T (among $\{1, 5, 10\}$ except for the cases of A3C and DDQN for $\varphi_7-\varphi_9$ among we use $\{5, 10\}$ ⁵) and report the setting (of Δ_T) which leads to the best performance (the least episode number and highest success rate) for each algorithm.

Evaluation Results The preliminary results are presented in Table. 2. The Δ_T columns indicate the best performing Δ_T for each algorithm. The ‘‘Success rate’’ columns in-

⁵ These methods with $\Delta_T = 1$ for $\varphi_7-\varphi_9$ shows bad performance and did not terminate in 5 days.

id	Formula	id	Formula
φ_1	$\Box \omega \leq \bar{\omega}$	φ_6	$\Box(\Box_{[0,t_1]} \omega \leq \bar{\omega} \rightarrow \Box_{[t_1,t_2]} v \leq \bar{v})$
φ_2	$\Box(v \leq \bar{v} \wedge \omega \leq \bar{\omega})$	φ_7	$\Box v \leq \bar{v}$
φ_3	$\Box((g_2 \wedge \diamond_{[0,0.1]} g_1) \rightarrow \Box_{[0.1,1.0]} \neg g_2)$	φ_8	$\Box \diamond_{[0,25]} \neg(v \leq v \leq \bar{v})$
φ_4	$\Box((\neg g_1 \wedge \diamond_{[0,0.1]} g_1) \rightarrow \Box_{[0.1,1.0]} g_1)$	φ_9	$\Box \neg \Box_{[0,20]} (\neg g_4 \wedge \omega \geq \bar{\omega})$
φ_5	$\Box \bigwedge_{i=1}^4 ((\neg g_i \wedge \diamond_{[0,0.1]} g_i) \rightarrow \Box_{[0.1,1.0]} g_i)$		

Table 1. The list of the evaluated properties on AT.

id	Δ_T				Success rate				numEpisode			
	A3C	DDQN	SA	CE	A3C	DDQN	SA	CE	A3C	DDQN	SA	CE
φ_1	5	1	10	5	100%*	100%*	65.0%	10.0%	16.5**	24.5	118.5	200.0
φ_2	5	1	10	5	100%*	100%*	65.0%	10.0%	11.5**	27.5	118.5	200.0
φ_3	1	1	1	1	75.0	5.0%	20.0%	85.0%	44.0	200.0	200.0	26.5
φ_4	1	1	1	1	75.0	10.0%	20.0%	85.0%	67.5	200.0	200.0	26.5*
φ_5	1	1	1	1	100%	100%	100%	100%	1.0	2.0	1.0	1.0
φ_6	10	10	10	10	100%*	100%*	70.0%	50.0%	3.5**	3.5**	160.5	119.0
φ_7	5	5	1	1	65.0%	100%**	0.0%	0.0%	125.0	63.0**	200.0	200.0
φ_8	10	10	10	1	80.0%	95.0%	90.0%	75.0%	72.0	52.0	83.0	21.0
φ_9	10	10	10	10	95.0%	100%**	15.0%	5.0%	46.0	12.0**	200.0	200.0

Table 2. The experimental result on AT.

dicating the success rate of the falsification process. The “numEpisode” columns show the median (among the 20 procedures) of the number of simulation episodes required to falsify the formula. If the falsification procedure fails, we consider the number of simulation episodes to be the maximum allowed episodes (200). We use the median since the distribution of the number of simulation episodes tends to be skewed.

The best results (success rate and numEpisode) of each formula are highlighted in bold. If the difference between the best entry of our methods and the best entry of the baseline methods is statistically significant by Fisher’s exact test and the Mann-Whitney U-test [13], we mark the best entry with * ($p < 0.05$) or ** ($p < 0.001$), respectively.

As shown in Table 2, RL-based methods almost always outperform baseline methods on success rate, which means RL-based methods are more likely to find falsified inputs with a limited number of episodes. This is because RL-based methods learn knowledge from the environment and generate input signals adaptively during the simulations. Among the statistically significant results of numEpisode, our methods are best for five cases ($\varphi_1, \varphi_2, \varphi_6, \varphi_7, \varphi_9$), while the baseline methods are best for one case (φ_4). For the case of φ_4 , it is likely because all variables in this formula take discrete values, thus, reinforcement learning is less effective. Further, DDQN tends to return extreme values as actions, which are not solutions to falsify φ_3 and φ_4 . This explains the poor performance of DDQN for the case of φ_3 and φ_4 .

5 Conclusion and Future Work

In this paper, we report an approach which adopts reinforcement learning algorithms to solve the problem of robustness-guided falsification of CPS systems. We implement our approach in a prototype tool and conduct preliminary evaluations with a widely adopted CPS system. The evaluation results show that our method can reduce the number of episodes to find the falsifying input. As a future work, we plan to extend the current work to explore more reinforcement learning algorithms and evaluate our methods on more CPS benchmarks.

References

1. *The ChainerRL Library*. <https://github.com/chainer/chainerrl>.
2. H. Abbas, G. Fainekos, S. Sankaranarayanan, F. Ivančić, and A. Gupta. Probabilistic temporal logic falsification of cyber-physical systems. *ACM Trans. Embed. Comput. Syst.*, 12(2s):95:1–95:30, May 2013.
3. H. Abbas and G. E. Fainekos. Convergence proofs for simulated annealing falsification of safety properties. In *50th Annual Allerton Conference on Communication, Control, and Computing, Allerton 2012, Allerton Park & Retreat Center, Monticello, IL, USA, October 1-5, 2012*, pages 1594–1601. IEEE, 2012.
4. T. Akazaki. Falsification of conditional safety properties for cyber-physical systems with gaussian process regression. In Y. Falcone and C. Sánchez, editors, *Runtime Verification - 16th International Conference, RV 2016, Madrid, Spain, September 23-30, 2016, Proceedings*, volume 10012 of *Lecture Notes in Computer Science*, pages 439–446. Springer, 2016.
5. Y. Annpureddy, C. Liu, G. Fainekos, and S. Sankaranarayanan. *S-TaLiRo: A Tool for Temporal Logic Falsification for Hybrid Systems*, pages 254–257. Springer Berlin Heidelberg, Berlin, Heidelberg, 2011.
6. Y. Annpureddy, C. Liu, G. E. Fainekos, and S. Sankaranarayanan. S-taliro: A tool for temporal logic falsification for hybrid systems. In P. A. Abdulla and K. R. M. Leino, editors, *Tools and Algorithms for the Construction and Analysis of Systems - 17th International Conference, TACAS 2011, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2011, Saarbrücken, Germany, March 26-April 3, 2011. Proceedings*, volume 6605 of *Lecture Notes in Computer Science*, pages 254–257. Springer, 2011.
7. H. A. Bardh Hoxha and G. Fainekos. Benchmarks for temporal logic requirements for automotive systems. *Proc. of Applied Verification for Continuous and Hybrid Systems*, 2014.
8. E. Bartocci, L. Bortolussi, L. Nenzi, and G. Sanguinetti. On the robustness of temporal properties for stochastic models. In T. Dang and C. Piazza, editors, *Proceedings Second International Workshop on Hybrid Systems and Biology, HSB 2013, Taormina, Italy, 2nd September 2013.*, volume 125 of *EPTCS*, pages 3–19, 2013.
9. E. Bartocci, L. Bortolussi, L. Nenzi, and G. Sanguinetti. System design of stochastic models using robustness of temporal properties. *Theor. Comput. Sci.*, 587:3–25, 2015.
10. E. Bartocci, J. Deshmukh, A. Donzé, G. Fainekos, O. Maler, D. Nickovic, and S. Sankaranarayanan. Specification-based monitoring of cyber-physical systems: A survey on theory, tools and applications. *The Handbook of Runtime Verification*, 2018.
11. G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. Openai gym, 2016.
12. J. D. Cook. Basic properties of the soft maximum. 2011.
13. G. W. Corder and D. I. Foreman. *Nonparametric statistics: A step-by-step approach*. John Wiley & Sons, 2014.
14. X. C. Ding, S. L. Smith, C. Belta, and D. Rus. MDP optimal control under temporal logic constraints. In *Proceedings of the 50th IEEE Conference on Decision and Control and European Control Conference, CDC-ECC 2011, Orlando, FL, USA, December 12-15, 2011*, pages 532–538. IEEE, 2011.
15. X. C. Ding, S. L. Smith, C. Belta, and D. Rus. Optimal control of markov decision processes with linear temporal logic constraints. *IEEE Trans. Automat. Contr.*, 59(5):1244–1257, 2014.
16. A. Donzé. Breach, A toolbox for verification and parameter synthesis of hybrid systems. In T. Touili, B. Cook, and P. B. Jackson, editors, *Computer Aided Verification, 22nd International Conference, CAV 2010, Edinburgh, UK, July 15-19, 2010. Proceedings*, volume 6174 of *Lecture Notes in Computer Science*, pages 167–170. Springer, 2010.

17. A. Donzé and O. Maler. Robust satisfaction of temporal logic over real-valued signals. In K. Chatterjee and T. A. Henzinger, editors, *Formal Modeling and Analysis of Timed Systems - 8th International Conference, FORMATS 2010, Klosterneuburg, Austria, September 8-10, 2010. Proceedings*, volume 6246 of *Lecture Notes in Computer Science*, pages 92–106. Springer, 2010.
18. S. Gu, T. Lillicrap, I. Sutskever, and S. Levine. Continuous deep q-learning with model-based acceleration. In M. F. Balcan and K. Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 2829–2838, New York, New York, USA, 20–22 Jun 2016. PMLR.
19. H. Ho, J. Ouaknine, and J. Worrell. Online monitoring of metric temporal logic. In B. Bonakdarpour and S. A. Smolka, editors, *Runtime Verification - 5th International Conference, RV 2014, Toronto, ON, Canada, September 22-25, 2014. Proceedings*, volume 8734 of *Lecture Notes in Computer Science*, pages 178–192. Springer, 2014.
20. B. Hoxha, H. Abbas, and G. E. Fainekos. Using s-taliro on industrial size automotive models. In G. Frehse and M. Althoff, editors, *1st and 2nd International Workshop on Applied Verification for Continuous and Hybrid Systems, ARCH@CPSWeek 2014, Berlin, Germany, April 14, 2014 / ARCH@CPSWeek 2015, Seattle, WA, USA, April 13, 2015.*, volume 34 of *EPIc Series in Computing*, pages 113–119. EasyChair, 2014.
21. X. Li, Y. Ma, and C. Belta. A policy search method for temporal logic specified reinforcement learning tasks. *CoRR*, abs/1709.09611, 2017.
22. X. Li, C. I. Vasile, and C. Belta. Reinforcement learning with temporal logic rewards. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2017, Vancouver, BC, Canada, September 24-28, 2017*, pages 3834–3839. IEEE, 2017.
23. R. Luna, M. Lahijanian, M. Moll, and L. E. Kavragi. Asymptotically optimal stochastic motion planning with temporal goals. In H. L. Akin, N. M. Amato, V. Isler, and A. F. van der Stappen, editors, *Algorithmic Foundations of Robotics XI - Selected Contributions of the Eleventh International Workshop on the Algorithmic Foundations of Robotics, WAFR 2014, 3-5 August 2014, Boğaziçi University, İstanbul, Turkey*, volume 107 of *Springer Tracts in Advanced Robotics*, pages 335–352. Springer, 2014.
24. O. Maler and D. Nickovic. Monitoring temporal properties of continuous signals. In Y. Lakhnech and S. Yovine, editors, *Formal Techniques, Modelling and Analysis of Timed and Fault-Tolerant Systems, Joint International Conferences on Formal Modelling and Analysis of Timed Systems, FORMATS 2004 and Formal Techniques in Real-Time and Fault-Tolerant Systems, FTRTFT 2004, Grenoble, France, September 22-24, 2004, Proceedings*, volume 3253 of *Lecture Notes in Computer Science*, pages 152–166. Springer, 2004.
25. V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. P. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu. Asynchronous Methods for Deep Reinforcement Learning. 48, 2016.
26. V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
27. D. Sadigh, E. S. Kim, S. Coogan, S. S. Sastry, and S. A. Seshia. A learning based approach to control synthesis of markov decision processes for linear temporal logic specifications. In *53rd IEEE Conference on Decision and Control, CDC 2014, Los Angeles, CA, USA, December 15-17, 2014*, pages 1091–1096. IEEE, 2014.
28. S. Sankaranarayanan and G. E. Fainekos. Falsification of temporal properties of hybrid systems using the cross-entropy method. In T. Dang and I. M. Mitchell, editors, *Hybrid Systems: Computation and Control (part of CPS Week 2012), HSCC'12, Beijing, China, April 17-19, 2012*, pages 125–134. ACM, 2012.
29. S. Silveti, A. Policriti, and L. Bortolussi. An active learning approach to the falsification of black box cyber-physical systems. In N. Polikarpova and S. Schneider, editors, *Integrated*

Formal Methods - 13th International Conference, IFM 2017, Turin, Italy, September 20-22, 2017, Proceedings, volume 10510 of *Lecture Notes in Computer Science*, pages 3–17. Springer, 2017.

30. S. E. Z. Soudjani and R. Majumdar. Controller synthesis for reward collecting markov processes in continuous space. In G. Frehse and S. Mitra, editors, *Proceedings of the 20th International Conference on Hybrid Systems: Computation and Control, HSCC 2017, Pittsburgh, PA, USA, April 18-20, 2017*, pages 45–54. ACM, 2017.

A Derivation of Equation (3) in Section 3

In Section 3, the derivation of (3) from (2) is omitted. In this Appendix, the omitted derivation is presented.

$$\arg \min_{u_{i+1}} \min_{u_{i+2}, \dots} \rho(\Box\varphi, \mathcal{M}([(0, u_0), (\Delta_T, u_1), \dots]), 0) \quad (2)$$

$$= \arg \min_{u_{i+1}} \min_{u_{i+2}, \dots} \min_{t \in \mathbb{R}} \rho(\varphi, \mathcal{M}([(0, u_0), (\Delta_T, u_1), \dots]), t) \quad (A1)$$

$$\sim \arg \min_{u_{i+1}} \min_{u_{i+2}, \dots} \min_{k=i+1, i+2, \dots} \rho(\varphi, \mathcal{M}([(0, u_0), \dots, (k\Delta_T, u_k)]), k\Delta_T) \quad (A2)$$

$$\sim \arg \min_{u_{i+1}} \min_{u_{i+2}, \dots} \left[-\log \left\{ 1 + \sum_{k=i+1}^{\infty} \{e^{-\rho(\varphi, \mathcal{M}([(0, u_0), \dots, (k\Delta_T, u_k)]), k\Delta_T)} - 1\} \right\} \right] \quad (A3)$$

$$= \arg \max_{u_{i+1}} \max_{u_{i+2}, \dots} \sum_{k=i+1}^{\infty} \{e^{-\rho(\varphi, \mathcal{M}([(0, u_0), \dots, (k\Delta_T, u_k)]), k\Delta_T)} - 1\} \quad (3)$$

(2) is the equation which derives the next input u_{i+1} . (A2) uses the fact φ is past-dependent and (A3) uses an approximation of minimum by the log-sum-exp function [12]. Finally, the last equation (3) is the same equation (3) in Section 3. This completes the derivation.