

Does Neuron Coverage Matter for Deep Reinforcement Learning? A Preliminary Study

Miller Trujillo¹, Mario Linares-Vásquez¹, Camilo Escobar-Velásquez¹, Ivana Dusparic², Nicolás Cardozo¹

¹Systems and Computing Engineering, Universidad de los Andes, Bogotá, Colombia

²School of Computer Science and Statistics, Trinity College Dublin, Ireland

{ma.trujillo10,m.linaresv,ca.escobar234}@uniandes.edu.co, ivana.dusparic@scss.tcd.ie, n.cardozo@uniandes.edu.co

ABSTRACT

Deep Learning (DL) is powerful family of algorithms used for a wide variety of problems and systems, including safety critical systems. As a consequence, analyzing, understanding, and testing DL models is attracting more practitioners and researchers with the purpose of implementing DL systems that are robust, reliable, efficient, and accurate. **First software testing approaches for DL systems have focused on black-box testing, white-box testing, and test cases generation, in particular for deep neural networks** (CNNs and RNNs). However, Deep Reinforcement Learning (DRL), which is a branch of DL extending reinforcement learning, is still out of the scope of research providing testing techniques for DL systems. In this paper, we present a first step towards testing of DRL systems. In particular, we investigate whether neuron coverage (a widely used metric for white-box testing of DNNs) could be used also for DRL systems, by analyzing coverage evolutionary patterns, and the correlation with RL rewards.

CCS CONCEPTS

• **Theory of computation** → **Reinforcement learning**; • **Software and its engineering** → **Software testing and debugging**.

KEYWORDS

Deep networks, Reinforcement learning, Coverage analysis, Testing

1 INTRODUCTION

With the recent advancement of Deep Neural Networks (DNNs) [10], and in particular their success in image and large data set processing, Deep Learning (DL) techniques are being applied to a wide variety of problems. Notwithstanding, the inner working of algorithms developed using DNNs remains to be largely unknown. As a consequence, in recent years there has been an increasing interest in measuring, analyzing, and testing DNNs [19]. However, we identify two main limitations in existing proposals. First, existing approaches focus on testing standard DNNs. Deep Reinforcement Learning (DeepRL) techniques differ from standard DNNs as they merge DNNs and RL, and are therefore not contemplated by current efforts. Second, existing approaches focus mostly on black-box testing of the DNN algorithms, which does not offer insight into the internals of the network, just about its output.

Both black-box and white-box testing approaches have been studied focusing on testing CNNs and RNNs [19]. However, black-box testing is not directly applicable to DeepRL because DeepRL applications could have a range of correct behaviors, differing in optimality (e.g., suitable value range rather than the exact value) [2]. When testing DNNs, there is a “ground truth” that serves as baseline for assessing the accuracy (e.g., true positives rate). Identifying the rate of correct actions that were learned with DeepRL, i.e., assessing how close the system is to optimal values, is not straight-forward, particularly when optimizing the long-term performance.

Conversely, white-box testing is a more suitable approach for DeepRL, because the internals of the system (i.e., neurons in a DeepRL system) can be inspected to measure coverage. This type of metrics have been already proposed to measure the percentage of neurons that are activated in a DNN system and different activation patterns. Some examples are neuron coverage [12], MC/DC [11], SS/Vs/SV/VV coverage [13], and neuron boundary coverage, multi-section neuron coverage, strong neuron coverage, and top neuron coverage [5]. However, as of today there is no similar study that is performed on DeepRL systems.

In this paper, we explore testing techniques for DeepRL approaches (Section 2). We use white-box testing to study the evolutionary coverage of deep networks, for the specific case of DeepRL (Section 3). To achieve this, we analyze neuron coverage and rewards obtained by two different models of Deep Q-Network (DQN) implemented for the Mountain Car Problem (MPC) (Section 4). The results of our study (Section 5) show that “good” neuron coverage does not necessarily mean success in a RL task.

Future work should replicate this study on a larger scale and with different DeepRL problems. In addition, more studies should be focused on analyzing the impact of coverage on the effectiveness of DeepRL systems, and defining metrics that measure effectiveness by combining internal behavior and outputs of DeepRL systems.

2 RELATED WORK

DeepRL is an extension of RL that takes advantage of the capabilities of DNNs to capture larger state-action spaces. As little work has been done in the field of testing RL-based and DeepRL-based software systems, in this section we present existing approaches for testing of DNN systems and their relation with our work.

Symbolic execution has been used for identifying system’s component activation based on inputs to: (1) generate more effective tests [8], (2) find pixel attacks that have the same activation pattern as the original image [3], and (3) generate test inputs [1]. DeepExplore [7] proposes neuron coverage as a test adequacy metric for DNNs. Although no existing work uses concrete execution,

recent approaches [12] use concolic execution to resolve coverage requirements, resulting in a 10% higher neuron coverage than DeepExplore.

Dynamic analysis has been also used, for instance, Wicker et al. [17] present a feature-guided black-box safety testing approach (for image classifiers) that relies on generating adversarial examples. Wang et al. [15], propose a way to test the robustness of *Natural language inference* models by swapping the order of the provided fragments. Ma et al. [4] present an adaptation of combinatorial testing for DNNs, along with a set of coverage criteria for DL systems.

Other approaches combine white- and black-box testing. For example, *DeepHunter* [18] uses metamorphic mutation to generate semantically preserved tests. Pei et al. [7] propose a white-box differential testing technique to generate test inputs for a deep learning system. Such approach is based on comparing multiple DL systems by cross-checking each other instead of using an oracle. DeepTest [14] systematically explores different parts of the DNN logic by generating test inputs that maximize the number of activated neurons.

The analysis of the related work suggests different approximations to test DNNs adequacy. The results from these proposals are promising in the evaluation of the models. However, a commonality in these models is that they do not incorporate a notion of optimality with respect to the system's goal (i.e., reward models), required for DeepRL. Therefore, we must inquiry the appropriateness of such techniques to assess DeepRL.

3 DEEP REINFORCEMENT LEARNING

Reinforcement Learning (RL) learns optimal actions for specific environment conditions by trial-and-error required to maximize the long term cumulative reward [9].

At each time step t , an agent perceives the environment's state s_t from state space S . The agent, then selects an action a_t for the available action set A . The agent receives a reward $r_t = R(s_t, a_t)$ when it transitions to the state s_{t+1} . **The discount factor $\gamma \in [0, 1]$ is applied to the future rewards. One of most widely used RL techniques is Q-Learning [16]**, in which an RL agents learn Q-values, i.e., quality of taking an action in a given state: $Q: S \times A \rightarrow \mathbb{R}$. Thus, the learning policy is represented as:

$$\pi(s) = \arg \max_{a \in A} Q(s, a) \quad (1)$$

$$Q(s, a) = (1 - \alpha)Q(s, a) + \alpha[R(s, a) + \max_{a' \in A} Q(s', a')] \quad (2)$$

where $0 < \alpha \leq 1$ is a learning rate.

In situations with a large state and action spaces it is unfeasible to learn Q-value estimates for each state-action pair independently as in standard tabular Q-Learning. Mnih et al. [6] propose Deep Q-Networks (DQN) as a technique to combine Q-Learning with CNN. DQN parameterizes an approximate value function $Q(s, a; \theta_i)$ where θ_i are the weights of the network at iteration i . The experience replay stores the agent's experiences $e_t = (s_t, a_t, r_t, s_{t+1})$ at each time step t in a dataset $D_t = \{e_1, \dots, e_t\}$ pooled over many episodes into a replay memory. Mini batches of experience are drawn at random from the dataset $(s, a, r, s) \sim U(D)$, and applied as

Q-updates during the training. The Q-learning update at iteration i uses the following loss function:

$$L_i(\theta_i) = E_{(s, a, r, s') \sim U(D)} \left[\left(r + \gamma \max_{a'} Q(s', a'; \theta_i^-) - Q(s, a; \theta_i) \right)^2 \right]$$

where the θ_i are the Q-network parameters at iteration i and θ_i^- are the target network parameters.

4 EMPIRICAL STUDY DESIGN

The *goal* of this preliminary study is to analyze the neuron coverage in a DeepRL system from different perspectives than those explored in existing DNN testing. We are interested in analyzing the evolution of neuron coverage at the (1) DeepRL system level and (2) neural network layer levels, and their correlation with cumulative environment reward (as reward is a feature that distinguishes DeepRL from DNN systems). This study uses two different models of DeepRL system implemented for solving a widely used benchmark problem in the DeepRL community (i.e., mountain car). In particular, this study addresses the following research questions (RQs):

RQ₁ *Is there a difference in the evolution patterns of neuron coverage for the different layers in a DeepRL system?* The goal of this **RQ** is to understand whether neuron coverage evolves differently for each of the layers in a DeepRL system.

RQ₂ *Is there a correlation between neuron coverage and cumulative rewards in DeepRL systems?* The purpose with this **RQ** is to investigate whether there is a relationship between higher cumulative rewards and higher coverage. Moreover, it could indicate how to select parameters during training, so we can manipulate coverage and obtain better rewards.

4.1 The Mountain Car Benchmark

The mountain car (MC) is a classic problem in RL [9]. In this problem, a car is positioned at the bottom of a valley between two mountains, and the engine of the car is not powerful enough to climb the mountain at the right and reach the goal at the mountain top. Hence, it is necessary to go back and forth between the two mountains to gain enough energy. The car's state is defined by its x-coordinate position and velocity.

At each time step, MC can take one of the three available actions: accelerate towards left, neutral/no acceleration, accelerate towards right. Its position is bounded by the coordinates $[-1.2, 0.6]$ and velocity by $[-0.07, 0.07]$.

The reward is received at each time step for MC being in a state (x, v) and is defined as expressed in Equation (3).

$$R(x, v) = \begin{cases} 10, & \text{if } x \geq 0.5 \\ -1, & \text{otherwise} \end{cases} \quad (3)$$

We used in this study DeepRL-based models for MC because they (1) have been widely studied by the DeepRL community, and (2) there are numerous open source models that allow us to focus on the models' execution and evaluation. Thus, we selected two of such DeepRL-based models. Here in after we will refer to the models

as MODEL A¹ and MODEL B.² Both of these models use DQNs, but with different neural network architectures. MODEL A, shown in Figure 1a, takes as input the state before and after performing an action. The initial state is used as previous and next state, in the first input. MODEL A uses four hidden layers, where layer two is flatten (*i.e.*, it maps exactly layer one). MODEL B (Figure 1b), takes as input only a state, and uses two dense hidden layers. Note that both models use a dense output layer with three neurons.

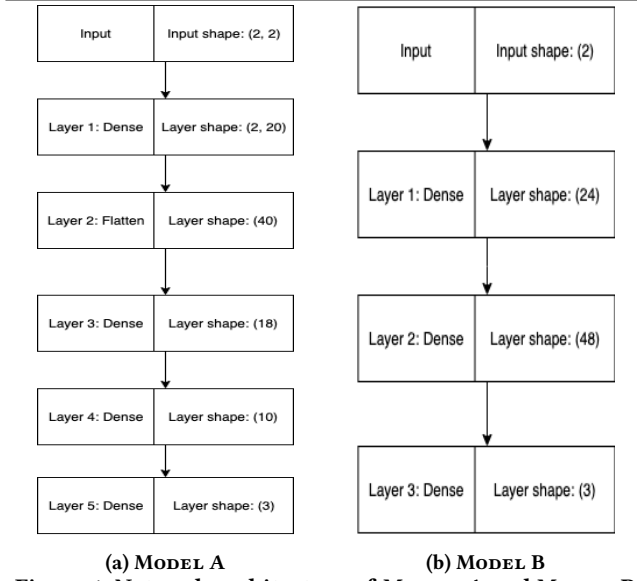


Figure 1: Network architecture of MODEL A and MODEL B

Model training and testing. Both models share the main structure for training and testing. MODEL A originally did not adjust the environment reward when the car reached the goal, which resulted in MC never learning how to reach the goal (*i.e.*, climb the mountain in maximum 200 steps). Therefore, we modified the original source code of MODEL A to use the same reward model as MODEL B as outlined above. Modifying MODEL A allow us to run the experiments with same criteria, and analyze across different architectures.

The models were trained and tested in sequences of iterations, episodes, and stages. An iteration i is a step or an action performed on the environment, and an episode e is a set of at most 200 iterations. A stage is a training - testing fold with 1000 training episodes and 20 testing episodes for each training episode in which the goal was achieved; *i.e.*, if all the training episodes of a stage achieve the problem goal, then, this stage will have 20k testing episodes. For this study, we used 5 stages.

A training iteration includes feeding the model with a batch of 32 inputs with the target output adjusted following Equation (2).

The models use $\alpha = 0.001$ (learning rate), and $\gamma = 0.99$ (discount factor) as parameters. To balance exploration and exploitation, both models use a decaying ϵ -greedy method in which ϵ starts at 1 and decays 0.05 after each episode. The minimum value of ϵ is 0.01.

4.2 Analysis Method

During training and testing of DeepRL systems, we record the activated neurons for each iteration and episode. Thus, given a system s , a layer $l \in s$, and a neuron $n \in N_l^s$, the set of neurons, the activation of n for a given iteration i and episode e is denoted by the indicator function $\lambda(n, i, e)$, returning 1 if the neuron is activated, and 0 otherwise. We denote $|s|$ as the total number of neurons in the system s , and $|l|$ as the total number of neurons in layer l .

By collecting λ for each layer, episode, and iteration, at any moment during training or testing we can compute metrics at different levels (*i.e.*, systems, and layers) and showing cumulative or snapshot-based (*i.e.*, at a given moment) results:

Neuron Coverage (NC): total number of neurons activated at a given iteration i and episode e divided the total number of neurons of the system s . Then, or a system s with L layers

$$NC(s, i, e) = \frac{\sum_{l \in L} \sum_{n \in N_l^s} \lambda(n, i, e)}{|s|}$$

Neuron Layered Coverage (NLC): total number of neurons of system s activated at layer l given iteration i and episode e divided the total number of neurons of the layer l , computed as

$$NLC(s, l, i, e) = \frac{\sum_{n \in N_l^s} \lambda(n, i, e)}{|l|}$$

Cumulative Neuron Coverage (CNC): total number of unique neurons activated until a given iteration i^* and episode e^* divided the total number of neurons of the system s . To do this, we require to compute the set of neurons that have not been activated until the iteration i^* and episode e^* ; we denote this set by $\hat{N}(s, e^*, i^*)$. CNC is defined as

$$CNC(s, e^*, i^*) = \frac{\sum_{e \leq e^*} \sum_{i \leq i^*} \sum_{n \in \hat{N}(s, e^*, i^*)} \lambda(n, i, e)}{|s|}$$

Cumulative Neuron Layered Coverage (CNLC): total number of unique neurons of a layer l^* activated until a given iteration i^* and episode e^* divided the total number of neurons of the layer l^* . To do this, we require to compute the set of neurons that have not been activated until the iteration i^* and episode e^* in a given layer l^* ; we denote this set by $\hat{N}(s, e^*, i^*, l^*)$. CNLC is defined as

$$CNLC(s, e^*, i^*, l^*) = \frac{\sum_{e \leq e^*} \sum_{i \leq i^*} \sum_{n \in \hat{N}(s, e^*, i^*, l^*)} \lambda(n, i, e)}{|l^*|}$$

To answer **RQ₁**, we use CNC and $CNLC$, as these metrics show all the neurons that are activated over the time (*i.e.*, count neurons activated on a past iteration or episode). Therefore, it is a simple way to observe evolution in coverage. We calculate these metrics in the models during training and testing phases. Then, we built time series with the values for each metric (*i.e.*, NLC and $CNLC$), and computed the descriptive statistics for each layer.

To answer **RQ₂**, we computed the Pearson and Spearman correlation coefficients between the coverage metrics (*i.e.*, NC and NLC) and the cumulative rewards values for both models, with the purpose of identifying whether there are relationships (linear or monotonic) between coverage and cumulative rewards.

¹<https://github.com/branavg/Deep-Q-learning>

²<https://github.com/pylser/Deep-Reinforcement-learning-Mountain-Car>

To record the neuron activation of hidden layers we used the functional API of Keras which allows us to define a model with multiple outputs and obtain all values from layers.³

5 RESULTS

To answer **RQ₁** and **RQ₂** for both models we executed 5 stages, with 1000 training episodes each, with at most 200 iterations per episode. The experiments were executed on a macOS Catalina V10.15.2 with 16 GB of RAM and a six core 2.2 GHz Intel core i7 processor. We used the libraries Keras (v2.2.4), Tensorflow (v1.15.0), Gym⁴ (v0.15.4), and NumPy (v1.17.4) to obtain and manipulate the data.

5.1 Evolution Patterns of Neuron Coverage

Figures 2 and 3 depict the distribution of CNC and CNLC for both models over several executions; average values are reported with a triangle shape. CNC is reported as network-coverage and individual CNLCs are reported as layer-x-coverage.

Training phase. On average, 95% and 86% of neurons were activated across the training episodes for MODEL A and MODEL B, correspondingly. In MODEL A, all the layers were able to execute all of their neurons during training. Across the shown executions, on average, the layers with the top cumulative coverage were layers 1-2 (97%) and layer 4 (90%). In MODEL B, layers 1 and 3 were able to execute all of their neurons during training. Layer 1, on average, achieved the top cumulative coverage (99%) followed by layer 3 (88%). Layer 5 in MODEL A reported the lower average cumulative coverage (80%) from the two models despite having only 3 neurons, showing that having less neurons does not necessarily result in higher cumulative coverage.

When looking into the details of the CNLC series, the series are less stable in MODEL A, *i.e.*, the behavior is less consistent in terms of the number of iterations required to activate all of the neurons in a layer and range of values. All the neurons from Input and Output layers in MODEL B are activated quicker than in MODEL A; however, this is the opposite for the hidden layers in MODEL B, where the progression is (in general) slower than in MODEL A.

Testing phase. In both models, we can see a reduction in the cumulative coverage at the layer and network levels. In the case of MODEL A, there is a reduction of 45% on the average cumulative coverage of the network; at the layers level the reduction is 42% on average. For MODEL B, the reduction at the network level is 30%, and 25% at the layers level on average. In both cases, only the output layers were able to activate all their neurons. While MODEL A executes 37,260 episodes with a success rate (*i.e.*, the car climbed the mountain) of 42%, MODEL B executes 74,240 episodes with a success rate of 77%. Note that there are 20 testing episodes for each training episode in which the problem goal was achieved.

Our first insight here is that MODEL A, despite having better cumulative coverage during training, has a lower success rate and lower cumulative coverage during testing –that is, better coverage in training does not necessarily mean better coverage in testing. Additionally, the evolution patterns of cumulative coverage during training are not necessarily the same during testing.

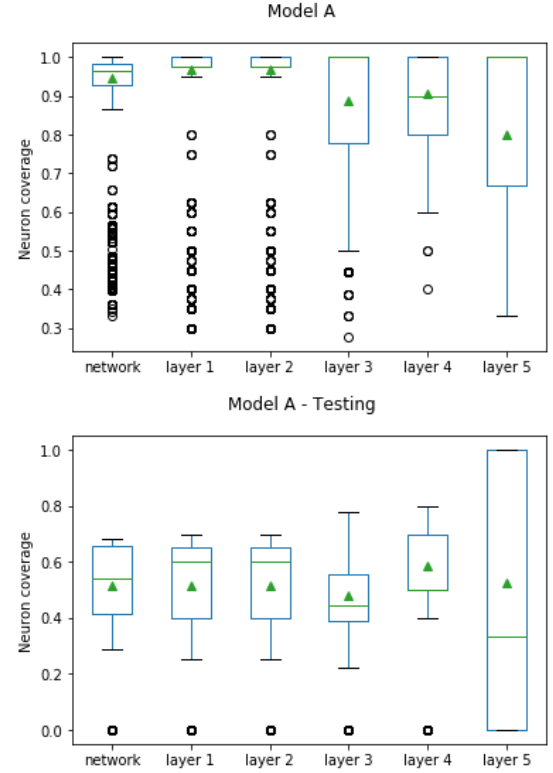


Figure 2: Distribution of cumulative coverage for Model A during the training and testing phases

5.2 Coverage-Reward Correlation

Figure 4 shows network (NC) and layered (NLC) coverage for each of the models (due to space restrictions, we depict three representative stages out of the total of five). We can observe that neuron coverage behavior for the last network layer is erratic across the different runs for both models. For the first two presented stages of MODEL A, the coverage in the first layer is higher than that of the following layers. The third stage, however, flips such behavior, showing a higher coverage in the layers towards the end of the net. In this case, layer_4, has the least coverage of all runs (under 25% except for a peak around episode 600). The behavior of flatten layer_2 corresponds exactly to layer_1 and is therefore not displayed in the figure. MODEL B, presents similar results. Earlier net layers have a higher coverage than later layers, across all episodes. Again, the behavior of the last layer is erratic. Nonetheless, in this model we observe that towards the last episodes, the coverage of the last layer converges to the network coverage.

In comparing the reward (Figure 5) obtained for each of the models with their corresponding coverage, we observe that the shape of the last layer's coverage is similar to the shape of the reward, as if the success of the mountain car would be associated to the coverage of the last layer of the model.

To evaluate this last observation, we calculate the correlation between coverage and reward. Both Pearson and Spearman correlation metrics show negative correlations. In the case of, NC the Pearson correlation is of -0.06 for MODEL A and -0.56 for MODEL B.

³<https://keras.io/getting-started/faq/#how-can-i-obtain-the-output-of-an-intermediate-layer>

⁴<https://gym.openai.com/>

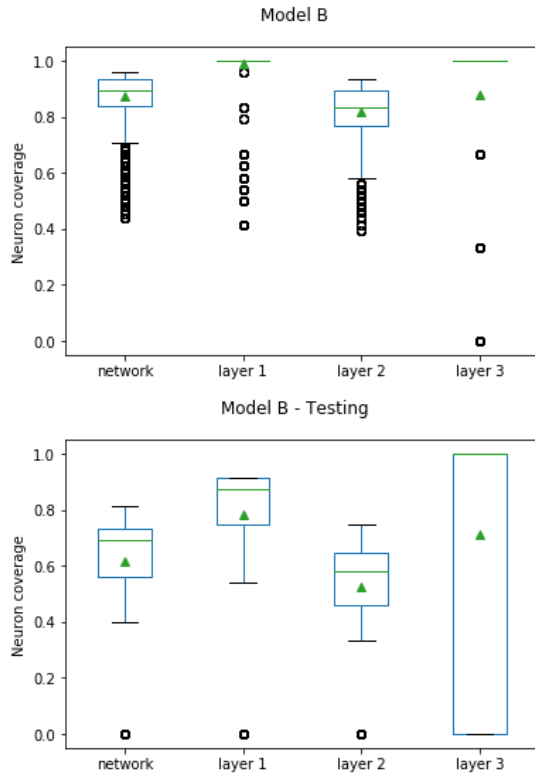


Figure 3: Distribution of cumulative coverage for Model B during the training and testing phases

The Spearman correlation is respectively of -0.07 and -0.56 ; as the coverage in MODEL B is higher than that of MODEL A. NLC's correlation presents similar results in a higher range. In the case of MODEL A, the Pearson correlation is in the range $[-0.36, 0.02]$, with a higher correlation for the first layer, and a lower correlation for the third one. The Spearman correlation is in the range $[-0.38, 0.06]$ presenting a similar behavior. MODEL B shows a similar correlation for both Pearson $[-0.62, -0.12]$ and Spearman $[-0.65, -0.04]$, where the highest correlation is in the first layer and the lowest in the middle layer.

6 CONCLUSION & FUTURE WORK

The advent of Deep Neural Network (DNN) has brought a special focus to Machine Learning (ML) algorithms and software development. Recently, testing of ML, and in particular DNNs, is being used as a means to understand and assure the quality of DNN-based software systems. This paper presents a first study in testing a particular type of DNN-based ML systems, namely, DeepRL. In particular, we study the use of neuron coverage as reliable white-box testing technique for DeepRL systems.

Our results, in contrast to neuron coverage in DNNs, show that **neuron coverage is not sufficient to reach substantial conclusions about the design or structure of DeepRL networks**. In particular, the negative correlation, or lack thereof, obtained in our results confirms the expected behavior of Deep Q-networks. The best possible coverage is achieved by extensive exploration, however,

this leads to exploring multiple actions that do not contribute in reaching the objective, and present a negative reward. Therefore, for the case of DeepRL, neuron coverage seems to assess state space exploration for the earlier layers of the network, as a means of exploring different actions.

As future work, the evaluation of neuron coverage as adequacy metric should be extended to other DeepRL scenarios to consolidate the validity of our results, in exploring for a metric that can correlate coverage and maximize reward. Additionally, a similar approach could be used to evaluate other DNN-based learning models, as means to evaluate the applicability of coverage testing to ML.

ACKNOWLEDGMENTS

We thank the anonymous reviewers for their comments to our work. The research in this paper has been partially supported by the Royal Irish Academy through the Charlemont Grant.

REFERENCES

- [1] A. Agarwal et al. *Automated test generation to detect individual discrimination in AI models*. arXiv preprint arXiv:1809.03260 (2018).
- [2] N. Cardozo, I. Dusparic, and M. Linares-Vásquez. *Perspectives in Testing Deep RL*. DeepTest'19. 2019.
- [3] D. Gopinath et al. *Symbolic execution for deep neural networks*. arXiv preprint arXiv:1807.10439 (2018).
- [4] L. Ma et al. *Combinatorial testing for deep learning systems*. arXiv preprint arXiv:1806.07723 (2018).
- [5] L. Ma et al. *DeepGauge: Multi-Granularity Testing Criteria for Deep Learning Systems*. Proc. of the 33rd ACM/IEEE Intl. Conf. on Automated Software Engineering. Association for Computing Machinery, 2018, pp. 120–131.
- [6] V. Mnih et al. *Human-level control through deep reinforcement learning*. Nature 518.7540 (2015), pp. 529–533.
- [7] K. Pei et al. *Deepxplore: Automated whitebox testing of deep learning systems*. Proc. of the 26th Symp. on Operating Systems Principles. ACM. 2017, pp. 1–18.
- [8] A. Ramanathan et al. *Integrating symbolic and statistical methods for testing intelligent systems: Applications to machine learning and computer vision*. Design, Automation & Test in Europe Conf. & Exhibition. IEEE. 2016, pp. 786–791.
- [9] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning, Second Edition An Introduction*. Second Edi. MIT Press, 2018, p. 550.
- [10] J. Schmidhuber. *Deep learning in neural networks: An overview*. Neural Networks 61 (2015), pp. 85–117.
- [11] Y. Sun, X. Huang, and D. Kroening. *Testing Deep Neural Networks*. CoRR abs/1803.04792 (2018).
- [12] Y. Sun et al. *Concolic testing for deep neural networks*. Proc. of the 33rd ACM/IEEE Intl. Conf. on Automated Software Engineering. ACM. 2018, pp. 109–119.
- [13] Y. Sun et al. *Structural Test Coverage Criteria for Deep Neural Networks*. IEEE/ACM 41st Intl. Conf. on Software Engineering: Companion Proc. (ICSE-Companion). 2019, pp. 320–321.
- [14] Y. Tian et al. *DeepTest: Automated testing of deep-neural-network-driven autonomous cars*. Proc. of the 40th Intl. Conf. on software engineering. ACM. 2018, pp. 303–314.

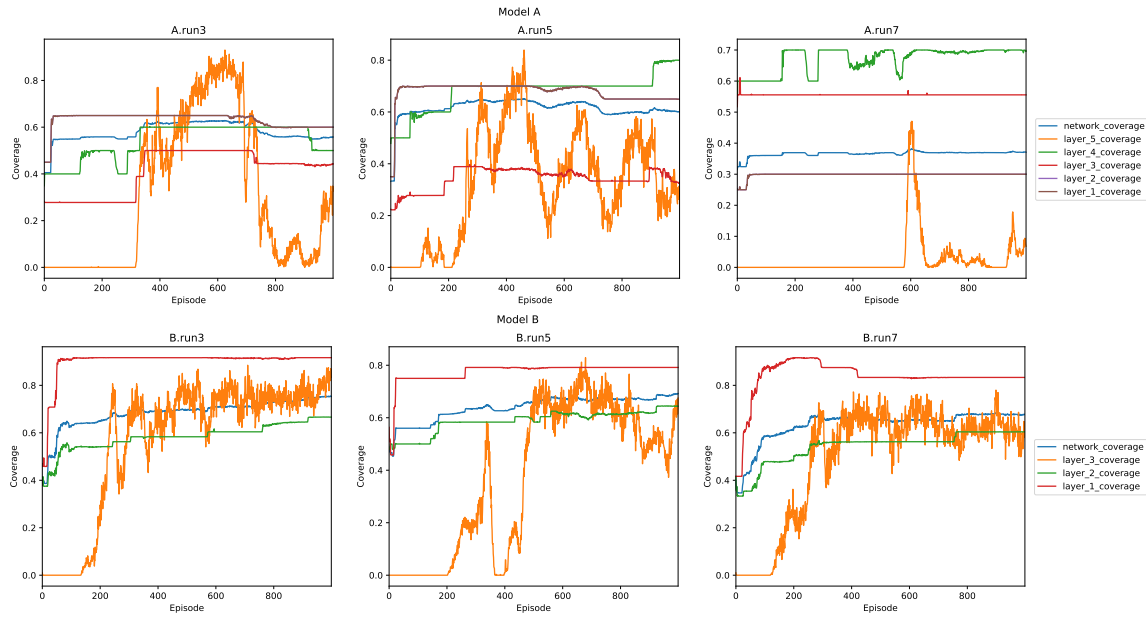


Figure 4: Neuron coverage per episode for the two models

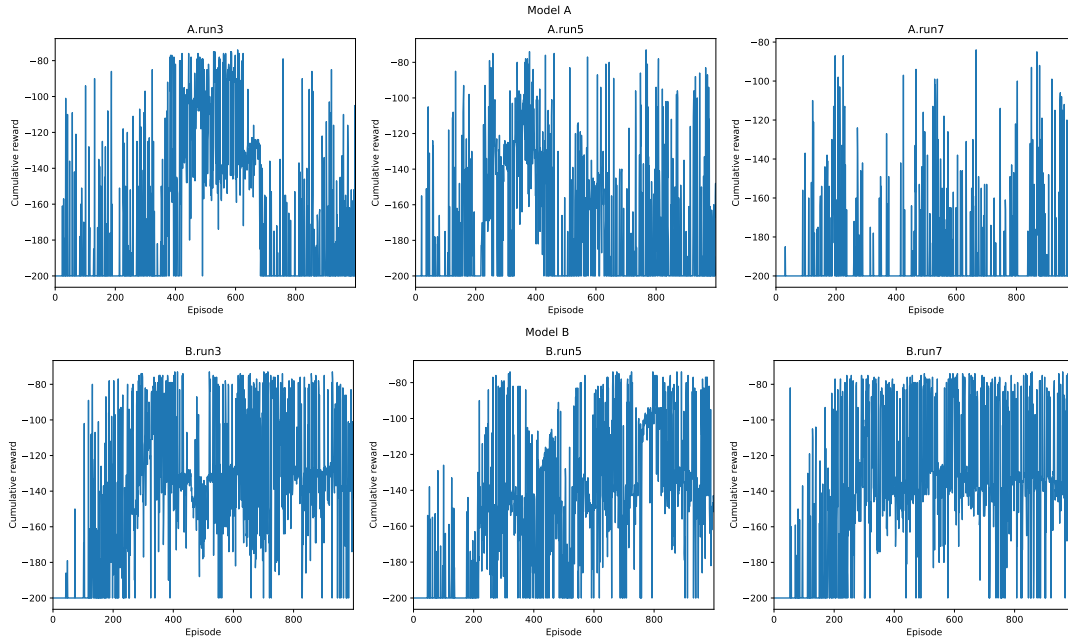


Figure 5: Reward per episode for the two models

- [15] H. Wang, D. Sun, and E. P. Xing. *What if we simply swap the two text fragments? A straightforward yet effective way to test the robustness of methods to confounding signals in nature language inference tasks*. Proc. of the AAAI Conf. on Artificial Intelligence. 2019, pp. 7136–7143.
- [16] C. J. C. H. Watkins and P. Dayan. *Technical Note: Q-Learning*. Machine Learning 8.3 (1992).
- [17] M. Wicker, X. Huang, and M. Kwiatkowska. *Feature-guided black-box safety testing of deep neural networks*. Intl. Conf. on Tools and Algorithms for the Construction and Analysis of Systems. Springer. 2018, pp. 408–426.
- [18] X. Xie et al. *Coverage-guided fuzzing for deep neural networks*. arXiv preprint arXiv:1809.01266 (2018).
- [19] J. M. Zhang et al. *Machine Learning Testing: Survey, Landscapes and Horizons*. CoRR abs/1906.10742 (2019).