

Reinforcement Learning State Space Properties for Model Checking of Interlockings^{*}

Ben Lloyd-Roberts¹, Phil James¹, and Michael Edwards¹

Swansea University, Swansea, UK

`{ben.lloyd-roberts, p.d.james, michael.edwards}@swansea.ac.uk`

Abstract. Railway interlockings serve as safety layer within railway signalling systems by ensuring proposed signalling requests are safe given the current state of the railway. As a vital part of any railway signalling system, interlockings are critical systems developed with the highest safety integrity level (SIL4) according to the CENELEC 50128 standard. The application of formal methods, in particular model-checking, in order to verify interlockings operate correctly is well established within academia and is beginning to see real applications in industry. However, the uptake of formal methods research within the UK rail industry has yet to make a substantial impact due to current approaches often producing false positives that require manual analysis during verification. Here, it is accepted that so-called invariants can be added to reduce the number of such false positives, however automatically computing these invariants remains a challenge. In this work we present a first approach illustrating how reinforcement learning can be used to learn properties of a state space for a given railway interlocking. Understanding how a verification problem can be formulated as one where machine learning can be successfully applied is the first step towards automated learning of invariants.

Keywords: Reinforcement Learning · Interlocking · Model Checking.

1 Introduction

Interlockings serve as a filter or ‘safety layer’ between inputs from operators, such as route setting requests, ensuring proposed changes to the current railway state avoid safety conflicts. As a vital part of any railway signalling system, interlockings are critical systems regarded with the highest safety integrity level (SIL4) according to the CENELEC 50128 standard. The application of model-checking to Ladder Logic programs in order to verify interlockings is well established within academia and is beginning to see real applications in industry. Bounded model-checking (BMC) is an efficient means of verifying a system through refutation. Given an abstracted model of a target system M and some safety properties ϕ , BMC searches for counterexamples up to some precomputed bound k . Search terminates when either an error trace is produced or the bound k is reached.

^{*} Supported by Siemens Mobility UK & EPSRC

Determining this completeness threshold to sufficiently cover all states is often computationally intractable given it's true value will depend on model dynamics and size of the search space. Additionally, the k-induction rule, which checks if the property ϕ holds inductively for k sequences of states, is constrained to acyclic graphs for guarantees of completeness.

An invariance rule allows us to establish an invariant property ψ which holds for all initial states and transitions up to a given bound. Invariants may hold for sub-regions of the state space, meaning complete coverage isn't necessary to learn them. Supporting k-induction with strengthening invariants helps reduce the overall search space for bounded model checking, proving that invariant aspects of the system need not be considered. This can help filter cases where false negative counter examples are triggered by unreachable states. Generating sufficiently strong invariants is a non-trivial process given their construction is heavily program dependent. Usually such invariants require domain knowledge, typically devised by engineers responsible for the program implementation.

We infer two principle challenges to address. First, formulating theoretical and practical frameworks to represent invariant finding as a reinforcement learning task. Second, devising an appropriate probabilistic strategy for identifying learnable state space properties within those frameworks. With the ultimate aim of adapting our approach to identify ladder logic invariants, we present a set of preliminary results illustrating the potential for RL schemes to maximise state space coverage via learning progressively larger values for the recurrence reachability diameter.

2 BMC & Program Invariants

- Not sure if this section is needed?

3 Reinforcement Learning

RL is a popular ML paradigm with demonstrably impressive capacity for modelling sequential decision making problems as the optimal control of some incompletely-known Markov Decision Process (MDP), known as the environment [16]. Given a permitted set of actions performable over a series of discrete time steps t , software agent(s) interact with and observe changes in the environment, receiving intermittent reward signals. Over accelerated trial-and-error, a function or policy π , is learned mapping states to optimal actions likely to return the greatest cumulative future reward. Interactions with the environment can be characterised as continuous or episodic, depending on constraints of the learned task.

Tabular reinforcement learning has existed for decades through dynamic programming [18], Monte-Carlo methods, and temporal difference learning. Resurgence in RL research over the last decade, applied to games [13, 15, 17], robotics [1, 3] and operations research [9], can be attributed to the fusion of a existing

value/policy iteration methods with powerful approximate learning techniques following the advent of deep learning [7].

Subsequent popularity in deep reinforcement learning (DRL) after breakthrough performance of deep Q-Learning on a set of gaming benchmarks [11] gave rise to improvements of historic methods, now leveraging deep learning [14]. Actor-Critic methods, combining policy learning and value function approximation saw particular successes.

Probabilistic learning is unlikely to provide guarantees of completeness, but can be used to supplement more suitable methods of doing so. For purposes of learning heuristics, we would like to maximise the information gathered, so that there is more data from which to identify patterns. With the aim of learning invariant properties that hold across states, state space coverage would be ideally maximised. Fortunately existing works have illustrated the efficacy of RL methods in learning such heuristics over large graph structures [8]. Distributing the task among concurrent workers for accelerated performance [5] and prioritising exploration in unfamiliar environments [2, 4, 12].

We select asynchronous advantage actor-critic (A3C) [10] to estimate both the value function and behaviour policy when exploring an environment. The asynchronous nature of the algorithm facilitates distributed exploration of state-action pairs via separate workers with a shared global policy network.

4 Learning Framework

Contacts, the binary variables constituting ladder logic program inputs, comprise the *action space* \mathcal{A} . Transitions between program states are simulated via interactions with this action space, initially uniformly at random. Each state, formally the *observation space* \mathcal{S} , is formed by the union of contacts and coils under their present valuation. Subsequently the environment unfolds as a set of reachable states for the respective ladder logic program. As workers improve their value estimates according to the reward function, a balance is maintained between stochastic action sampling and predictions from the policy network.

Reinforcing this behaviour requires motivating agents via a sparse reward signal with potentially strong temporal dependence. In light of this, we influence agents to pursue the longest loop free path where a reward scheme positively rewards sequences of novel observations. Inversely, negative rewards are issued for repeated observations within a training episodes. Workers are initialised with separate environment instances to accumulate experience independently. A global set of observations is asynchronously updated by each worker at set intervals to compile shared experiences.

For practicality, each environment has an associated max number of episodes T_{\max} to bound training runs. Additionally we utilise two forms of early termination to avoid superfluous interaction with the environment. First, if the worker performance curves converges to some local minima. Second, if all reachable states have been observed at least once. This update frequency depends on both the number of workers and precomputed size of the environment.

5 Environment Generation

Given exhaustive search of large state spaces is often computationally intractable, we generate a set of ladder logic programs where the number of reachable states and recurrence reachability diameter are known. Using existing models of ladder logic structures as a base template [6], we derive progressively larger programs by sequentially introducing additional rungs. This way a constrained yet predictable pattern of growth is devised. If $|S(p_i)|$ represents the number of reachable states for a program p_i , a subsequently generated program p_{i+1} with one additional rung, has $2|S(p_i)| + 1$ reachable states. Through a series of training runs on each environment we record the number of states observed by workers to gauge the overall state space coverage.

6 Results

Preliminary results from applying our approach to a number of generated programs are outlined in Table 1. Coverage metrics are expectedly maximised for environments with a small number of reachable states. Still, acceptable levels of coverage are observed for programs with more than $1e5$ states.

Disparate correlations between the coverage metrics and the total episodes may stem from the stochasticity of environment and network initialisation, leading to differing local minima.

We empirically observe three scenarios in the event environment coverage is not maximised: (i) Performance curves converge to a reward bound by the longest loop free path; (ii) Performance curves converge to a suboptimal reward; (iii) Performance improves linearly before collapsing to some suboptimal reward.

The on-policy nature of actor critic meant previous experiences from old policy iterations were forgotten, thus biasing behaviour to the most recent model updates and introducing sample inefficiency.

Given the A3C algorithm requires workers to asynchronously update their shared network every T_{\max} steps or on episode termination, larger values for T_{\max} consolidate more information regarding worker trajectories before applying gradient updates to their local network. Therefore we observe a strong dependence on this parameter both in terms of increasing the k bound when introducing workers to larger environments. Prior to random episode initialisation and fine-tuning network update frequencies, workers seldom covered 100% of the smaller environments.

7 Conclusion & Future Work

A3C shows some promising preliminary results but is clearly limited in its capacity to scale.

The algorithm shows proclivity to collapse following small network adjustments in previously unexplored regions of large environments.

Table 1. Initial results applying A3C learning over progressively larger environments.

Environment Metrics			Training Metrics			
States	Reachable States	Actions	Observed K	States Observed	Coverage	Total Episodes
2^{12}	7	2	6	7	100	1.00E+04
2^{14}	15	4	14	15	100	1.00E+04
2^{16}	31	8	28	31	100	1.00E+04
2^{18}	63	16	48	63	100	1.00E+05
2^{20}	127	32	33	127	100	1.00E+05
2^{22}	255	64	76	255	100	1.00E+04
2^{24}	511	128	49	377	100	1.00E+05
2^{26}	1023	256	306	1023	100	3.00E+05
2^{28}	2047	512	538	2047	100	3.00E+05
2^{30}	4095	1024	1418	4084	99.731	3.00E+05
2^{32}	8191	2048	1712	7907	96.532	3.00E+05
2^{34}	16383	4096	1498	15654	95.550	3.00E+05
2^{36}	32767	8192	2879	27752	84.694	3.00E+05
2^{38}	65535	16384	1969	58373	89.071	3.00E+05
2^{40}	131071	32768	2692	108638	82.884	2.00E+05
2^{42}	262143	65536	1406	5199317	76.033	3.00E+05
2^{44}	524287	131072	1782	325781	62.137	2.00E+05
2^{46}	1048575	262144	1593	671645	64.053	3.00E+05
2^{48}	2097151	524288	1598	1206867	57.547	3.00E+05
2^{50}	4194303	1048576	2527	1739939	41.48	1.50E+05

Introduce experience replay for distributed learning to improve on-policy bias and sample efficiency. Additional algorithms to implement for scalability - we can guarantee the increase of action spaces when transitioning to industry scale ladder logic programs where the number of contacts are known. It may be worth exploring DDPG to improve learning solely based on the observation space.

Intrinsic motivation for exploration.

IMPALA to improve both sample efficiency over A3C and robustness to network architectures and hyperparameters. Use of LSTM also improves performance given GPU acceleration is maximised on larger batch updates. A3C in our setting experiences variable episodic updates.

References

1. Bloesch, M., Humplik, J., Patraucean, V., Hafner, R., Haarnoja, T., Byravan, A., Siegel, N.Y., Tunyasuvunakool, S., Casarini, F., Batchelor, N., et al.: Towards real robot learning in the wild: A case study in bipedal locomotion. In: Conference on Robot Learning. pp. 1502–1511. PMLR (2022)
2. Gordillo, C., Bergdahl, J., Tollmar, K., Gisslén, L.: Improving playtesting coverage via curiosity driven reinforcement learning agents. arXiv preprint arXiv:2103.13798 (2021)

3. Gu, S., Holly, E., Lillicrap, T., Levine, S.: Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. In: 2017 IEEE international conference on robotics and automation (ICRA). pp. 3389–3396. IEEE (2017)
4. Haarnoja, T., Zhou, A., Abbeel, P., Levine, S.: Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor (2018)
5. Hoffman, M., Shahriari, B., Aslanides, J., Barth-Maron, G., Behbahani, F., Norman, T., Abdolmaleki, A., Cassirer, A., Yang, F., Baumli, K., et al.: Acme: A research framework for distributed reinforcement learning. arXiv preprint arXiv:2006.00979 (2020)
6. James, P., Lawrence, A., Moller, F., Roggenbach, M., Seisenberger, M., Setzer, A., Kanso, K., Chadwick, S.: Verification of solid state interlocking programs. In: International Conference on Software Engineering and Formal Methods. pp. 253–268. Springer (2013)
7. LeCun, Y., Bengio, Y., Hinton, G.: Deep learning. *nature* **521**(7553), 436–444 (2015)
8. Manchanda, S., Mittal, A., Dhawan, A., Medya, S., Ranu, S., Singh, A.: Learning heuristics over large graphs via deep reinforcement learning. arXiv preprint arXiv:1903.03332 (2019)
9. Mazyavkina, N., Sviridov, S., Ivanov, S., Burnaev, E.: Reinforcement learning for combinatorial optimization: A survey. *Computers & Operations Research* **134**, 105400 (2021)
10. Mnih, V., Badia, A.P., Mirza, M., Graves, A., Lillicrap, T.P., Harley, T., Silver, D., Kavukcuoglu, K.: Asynchronous methods for deep reinforcement learning (2016)
11. Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., Riedmiller, M.: Playing atari with deep reinforcement learning. arXiv preprint arXiv:1312.5602 (2013)
12. Ostrovski, G., Bellemare, M.G., van den Oord, A., Munos, R.: Count-based exploration with neural density models (2017)
13. Schaul, T., Quan, J., Antonoglou, I., Silver, D.: Prioritized experience replay. arXiv preprint arXiv:1511.05952 (2015)
14. Schulman, J., Levine, S., Moritz, P., Jordan, M.I., Abbeel, P.: Trust region policy optimization (2017)
15. Silver, D., Huang, A., Maddison, C.J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al.: Mastering the game of go with deep neural networks and tree search. *nature* **529**(7587), 484–489 (2016)
16. Sutton, R.S., Barto, A.G.: Reinforcement learning: An introduction. MIT press (2018)
17. Vinyals, O., Babuschkin, I., Czarnecki, W.M., Mathieu, M., Dudzik, A., Chung, J., Choi, D.H., Powell, R., Ewalds, T., Georgiev, P., et al.: Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature* **575**(7782), 350–354 (2019)
18. Watkins, C.J., Dayan, P.: Q-learning. *Machine learning* **8**(3), 279–292 (1992)