

# Finding $k$ : Asynchronous Reinforcement Learning for $k$ -Induction Bounded Model Checking<sup>★</sup>

Ben Lloyd-Roberts<sup>1</sup>, Phil James<sup>1</sup>, Michael Edwards<sup>1</sup>, and Tom Werner<sup>2</sup>

<sup>1</sup> Swansea University, Swansea, UK

<sup>2</sup> Siemens Rail Automation UK, Chippenham, UK

{ben.lloyd-roberts, p.d.james, michael.edwards}@swansea.ac.uk  
WT.Werner@siemens.com

**Abstract.** Solid State interlockings in the railway domain are computerised safety critical systems responsible for operating the signalling components directing railway traffic. The formal verification of such systems using inductive methods, such as bounded model-checking (BMC), often face limited guarantees of correctness given they operate within a bounded region of  $k$  depth. In this work we devise a reinforcement learning (RL) strategy to support SAT-based bounded model checking by improving measures of completeness. We outline a mapping from labelled transition systems traditionally used as BMC abstractions, to Markov decision processes (MDP) used as interactive RL environments. Thereafter software agents are trained to explore this environment for the longest acyclic sequences of execution to an upper bound  $k$ . We test the feasibility of our approach against a series of artificially generated ladder logic programs with incrementally larger state spaces. A number of metrics suggested in RL literature are then gathered to evaluate the robustness and scalability of the approach, particularly in industrial contexts. Finally we produce an *exploration graph* representing the structure of the traversed state space. These visualisations are then fed back to engineers responsible for designing such programs.

**Keywords:** First keyword · Second keyword · Another keyword.

## 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. As early

---

<sup>★</sup> Supported by Siemens Mobility UK & EPSRC

as 1995, Groote et al. [5] applied formal methods to verify an interlocking for controlling the Hoorn-Kersenboogher railway station. They conjecture the feasibility of verification techniques as a mean of ensuring correctness criteria on larger railway yards. In 1998, Fokkink and Hollingshead [6] suggested a systematic translation of Ladder Logic into Boolean formulae. Newer approaches to interlocking verification have also been proposed in recent years [7, 8, 9, 10, 11]. This includes work by Linh et al. which explores the verification of interlockings written in a similar language to Ladder Logic using SAT-based model checking. After two decades of research, academic work [12, 13] has shown that verification approaches for Ladder Logic can indeed scale; in an industrial pilot, Duggan et al. [14] conclude: “Formal proof as a means to verify safety has matured to the point where it can be applied for any railway interlocking system.” In spite of this, such approaches still lack widespread use within the UK Rail industry. Industrial standards for the railway and related domains increasingly rely on the application of formal methods for system analysis in order to establish a design’s correctness and robustness. Recent examples include the 2011 version of the CENELEC standard on railway applications, the 2011 ISO 26262 automotive standard, and the 2012 Formal Methods Supplement to the DO-178C standard for airborne systems. However, the application of formal methods research within the UK rail industry has yet to make a substantial impact. Principally our work aims to address one of the issues hindering its uptake, by removing the need for manual analysis of false negative error traces produced during verification.

Bounded model-checking is an efficient means of verifying a system through refutation. Given an abstracted model of a target system  $M$  and some safety properties  $\phi$ , 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. 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. Additionally, the k-induction rule, which checks if the property  $\phi$  holds inductively for sequences of states, is constrained to acyclic graphs for guaranteeing completeness.

An invariance rule allows us to establish an invariant property  $\psi$  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 in an academic-industry partnership with Siemens Rail Automation UK to represent interlocking verification as a goal-orientated rein-

forcement learning task. Second, devising an appropriate strategy for learning invariants within those frameworks.

The aim of this paper is to address this first challenge. Reinforcement learning is a popular machine learning paradigm with a demonstrably impressive capacity for learning near optimal strategies for goal-orientated tasks. Such approaches are well suited to problems where one need systematically learn the behaviour of a deterministic system, record observations regarding different 'states' of behaviour and identify patterns across those states. Generalisation of learned policies across different 'environments' being one of the key challenges in RL, we first aim to learn conceptually simpler goals, such as finding the progressively larger upper bounds for BMC. Understanding how a verification problem can be formulated as one of where machine learning can be successfully applied suggests the promise of learning invariants this way.

## 2 Preliminaries

### 2.1 Ladder Logic

Ladder logic programs used to implement interlocking systems allow engineers to define system behaviour in terms of conditions (rungs) and outputs (coils) expressed in a derivation of boolean algebra. Existing SAT-based approaches have already demonstrated translations from ladder logic programs to propositional formulae. The valuation of variables comprising those expressions at any given time is an instance of a program state. The set of reachable states then constitutes the abstractions notion of complete system behaviour. We can further abstract the program's behaviour by considering any one unique valuation of all variables constitutes a unique program 'state'.

### 2.2 Labelled Transition Systems

Each rung computes a new coil value to be used in the next execution cycle hence the significance of rung ordering in ladder logic programs. Labelled transition systems are commonly used to model this change in valuation over execution cycles.

**Definition 1 (Labelled Transition System).** *A labelled transition system  $\mathcal{L}$ , given a ladder logic formula  $\psi L$  is denoted  $\mathcal{L}(\psi L)$  and represented by the 4-tuple  $\langle S, T, R, S_0 \rangle$ , where  $S$  is a finite set of states,  $T$  is a finite set of transition labels,  $R$  is a labelled transition relation and  $S_0$  is the set of initial states.*

### 2.3 Bounded Model-Checking

Model checking is a formal verification technique stemming from the need to systematically check whether certain properties hold for different configurations (states) of a given system. Given a finite transition system  $T$  and a formula  $F$ , model checking attempts to verify that  $s \vdash F$  for every system state  $s \in T$ , such

that  $T \vdash F$ . Model checking is typically performed over three distinct phases as described by Baier [20]. Initially the model, that is the mathematical representation of the system, must be constructed. This is accomplished using a model description language, of which there are several. Next the model must be tested through simulation to ensure the abstracted system functions as expected. Provided the constructed model behaves adequately the final stage of modelling requires formulating the properties to be verified according to some specification language. The constructed model and property definitions then enable the model checker to perform runs of the system for state verification. The model checking process culminates in the analysis phase wherein verification results are generated. Properties which hold for all tested states produce a ‘safe’ output. In the event a state is found to violate any specified properties, that is  $\neg \models F$ , a counter example trace is provided by the model checker indicating which state(s) caused the violation. Results may also indicate that the model, property formulation or simulation process are insufficient for verification and therefore require further refinement. If so, the adjustments are made and the pipeline is revised and repeated. Verification means abstraction - "...existing works have successfully mapped ladder logic semantics to propositional logic. Here expressions can be reordered in CNF ready for a SAT-solver."

Abstraction means translating logic to a model - "...any unique valuation of the variables comprising the target program is viewed as an individual system state. This process is expanded in sec 2"

Modelling means we can apply verification - "...decomposing the target system to an abstracted set of states allows model-checking to systematically check each state for safety conflicts"

Verification has limitations - "... while a systematic process there are no guarantees of complete coverage. The longest loop-free path would be useful knowledge" Bounded model checking (BMC) provides a subset of the state space to avoid combinatorial explosion given the number of unreachable states. A bound is defined as a finite number of transitions between states. However unless all states are reached within that number of transitions then the exploration process is incomplete. This trade-off ensures that unreachable states are not considered but does not guarantee completeness within the bounded model

Primary limitations of the solution presented by Kanso are address in [16, 12, 17, 18]. Due to the inductive step of verification the tool checks to see if a given state satisfies some condition but does not consider the possibility of unreachable states which violate the same safety condition. Consequently if such benign violations are detected the tool risks generating false negative counter examples which requires manual inspection by an experienced engineer. This concept is illustrated in Fig. 3.5. A potential solution to this problem was introduced in subsequent work by James in 2010 [17]. In their exploration of SAT-based bounded model checking, James addresses the limitations of Kanso’s approach in terms of state reachability by introducing invariants to suppress false negatives.

Clearly, introducing appropriate invariants to avoid unattainable configurations requires a comprehensive understanding of both the abstracted and phys-

ical system. Additionally, as system complexity increases so do the number of unreachable states. This could potentially necessitate thousands of unique invariants to suppress triggered false negatives. Currently this process requires manual analysis by specialist engineers making automation highly desirable. Generating sufficiently strong invariants automatically is an extremely complex task, one which has received considerable attention in academic literature. From software engineering techniques [21, 22, 23, 24, 25, 26, 27] to hybrid methods incorporating machine learning [28, 29, 30, 31, 32], researchers have proposed various approaches to invariant finding with varying degrees of success.

## 2.4 Reinforcement Learning

Reinforcement learning (RL) is a principle ML paradigm which leverages human ability to formalise sequential decision making problems with respect to some goal-orientated tasks as optimal control of some incompletely-known MDP [1]. Given a user defined environment  $\mathcal{E}$  comprising a set of states  $S$ , a set of actions  $\mathcal{A}$  and a state transition function  $f : S \times \mathcal{A} \rightarrow S^+$ , a software agent is trained to learn the optimal action  $a$  to take from state  $s$  at time step  $t$  given a positive or negative reward signal  $r \in \mathcal{R}$  [45, 46, 47].

Agents then attempt to converge toward an optimal deterministic or stochastic policy  $P(s', r|s, a)$  which maximises cumulative future reward. Rewards signals are often sparse, meaning agents may traverse several states before receiving any feedback. Thus determining rewards over future time steps, the *expected return*, becomes desirable  $G_t = r_{t+1} + r_{t+2} + \dots + r_{\mathcal{T}}$ . In episodic tasks  $\mathcal{T}$  refers to the time step where a terminal condition is triggered. Subsequently the environment is reset to its initial state  $I_s \in$ . In continuous or sparsely rewarded episodic tasks, we apply a scalar discount factor  $\gamma \in [0, 1]$  at each time step to help enumerate future cumulative reward over a potentially infinite horizon.

**Policies** A policy is a rule set learned by an agent which determines the next best action. This is determined based on samples from the current probability distribution over any given action yielding a positive reward given the state.

**Definition 2 (Policy).** *A probability distribution  $\pi(a|s)$  describing the likelihood of action  $a$  returning a positive reward given environment state  $s$ .*

**Definition 3 (Trajectory).** *A sequence  $\tau$ , comprised of the states  $s_0, \dots, s_n$ , actions  $a_0, \dots, a_n$ , and rewards  $r_0, \dots, r_n$  experienced by an agent following a policy  $\pi$*

**Theorem 1 (Optimal Policy).** *: Given a finite MDP, there exists at least one optimal policy  $J$ , over a set of parameters  $\theta$ , such that:  $J(\theta) = E_{\pi}[r(\tau)]$  Finding the optimal policy requires some measure of decision quality. These are referred to as value functions.*

**Value Functions** Value functions return an expected reward, which determines the overall benefit of a given policy. They can be differentiated by the observations they use to determine the value. So called state-value functions, denoted  $v_\pi$ , provide the expected future reward given a state  $s_t$  and adhering to policy  $\pi$ .

$$v_\pi = E[G_t \mid S_t = s] = E \left[ \sum_{i=0}^{\infty} \gamma^i r_{t+i} \mid S_t = s \right] \quad (1)$$

**Theorem 2 (Optimal state-value function).** : *An optimal policy  $J(\theta)$  has an optimal state-value function for all  $s \in S$ , defined:*

$$v^*(s) = \max_{\pi} v_{\pi}(s) \quad (2)$$

Informally the optimal state-value function returns the maximum possible expected return of any policy for each state. Additionally action-value functions, denoted  $q_\pi$ , provides the value assigned to an action under policy  $\pi$ . Action-value functions, also known as Q-functions, produce a quality valuation assigned to the current action which determines the expected reward from subsequently adhering to the same policy  $\pi$ . We define the value function  $q_\pi(s)$ :

$$q_\pi(s, a) = E[G_t \mid S_t = s, A_t = a] = E \left[ \sum_{i=0}^{\infty} \gamma^i r_{t+i} \mid S_t = s, A_t = a \right] \quad (3)$$

The output of our Q-function is known as the Q-value.

**Theorem 3 (Optimal Q-function).** *An optimal policy,  $J\theta$  has an optimal Q-function for all  $s \in S$  and all  $a \in A$ :*

$$Q^*(s, a) = \max_{\pi} Q_{\pi}(s, a) \quad (4)$$

Informally, the optimal Q-function returns the maximum possible expected return of any policy for all state-action pairs. The Bellman principle of optimality [48], states for any state-action pair at the current time step, the expected return from an initial state  $s$ , taking action  $a$  and following the optimal policy  $J(\theta)$  thereafter is equal to the expected reward from taking action  $a$  in state  $s$ , plus the maximum achievable expected discounted return from any subsequent state-action pairs.

$$q^*(a, s) = E \left[ r_{t+1} + \gamma \max_{a'} q^*(s', a') \right] \quad (5)$$

The Bellman optimality equation is an integral metric used to learn the optimal Q-function which in turn is used to learn the optimal policy. Given an optimal Q-function action  $a^{prime}$ , a Q-learning algorithm will find the best action  $a'$  which maximises the Q-value for  $s'$ .

## 2.5 Reward/Goal Shaping

Designing appropriate reward systems is arguably one of the greatest challenges in RL. We determine the cumulative reward for a trajectory

$$\mathcal{R} = \sum_{t=0}^{T-1} R(s_t, a_t) \quad (6)$$

. We introduce a scalar discount factor  $\gamma$  to adapt modify returns such that rewards are weighted according to their distance from the current time step

$$G_t = \sum_{i=0}^{\infty} \gamma^i r_{t+i} \quad (7)$$

. While we wish to both encourage the agent to make optimal decisions at each time step, the discount factor helps prioritise rewards for immediate actions yet will determine the future reward of subsequent actions for several potential routes.

## 2.6 Exploration Strategy

Following Q-value initialisation, agents have share no estimates regarding optimal state-action pairs. Therefore we set an initial exploration rate  $\epsilon \in [0, 1]$ , and scalar decay value applied at each time step, to balance action sampling according to the learned policy or by some random probability distribution over  $A_t$  [53, 54]. Pure exploitation according to  $\pi$  likely results in sub-optimal solution where insufficient observations have been made regarding the global environment. Conversely, sustained exploration reduces the likelihood of an agent leveraging any previous observations to maximise reward. The  $\epsilon$ -greedy strategy is a popular approach to achieving this balance [49, 50, 51, 52]. Values nearing 0 represent a greedy strategy where the agent solely adheres to  $\pi$ . Values nearing 1 encourage random sampling of the  $A_t$ . epsilon decay value to decrease the exploration rate for each episode. Including an epsilon decay value provides some beneficial learning properties. As the agent learns more about the environment it relies less on exploration. Instead we expect to have observed a sufficiently comprehensive understanding of the environment thus settle on a purely greedy strategy. We also influence exploration and exploitation at every time step. First we randomly generate a number,  $0 \leq i \leq 1$  to compare with  $\epsilon$ . If  $i > \epsilon$  the agent selects an action according to  $\pi$ . Where  $i < \epsilon$ , the agent randomly samples an action to explore the environment.

**Replay Memory** To support the training of our policy network we generate a dataset of maximum  $N$  experiences representing an agent’s replay memory, sampled at each time step  $t$ . We characterise experiences as the 4-tuple  $(S, A, R, S')$ , associating each sample with an importance weight to avoid learning linear correlations over short sequences of steps.

**Theorem 4 (Experience Replay).** *An agent's experience at time step  $t$  is defined:*

$$e_t = (s_t, a_t, r_{t+1}, s_{t+1}) \quad (8)$$

**DQNs** Deep reinforcement learning (DRL) addresses the limited capacity of Q-functions in determining an optimal policy for complex problems. Instead a policy network, or deep Q-Network, is implemented to approximate the optimal Q-function [46]. The DQN receives a number of states from the environment as input to approximate a function generating the optimal Q-value such that it satisfies the Bellman equation. Q-values output by the network are then compared to an approximated 'optimal' target value in order to compute the loss. Q-learning refers to a policy learning method which uses Q-functions to calculate maximum expected future reward. Learning is formulated as an iterative process of parameter adjustment known as value iteration

**Network Updates and Future Estimates** Computing model loss for  $(s_t, a_t)$ , we determine the difference between the estimates. Without a discrete comparable Q-value for  $s'$  we must approximate this value by passing  $s'$  to the policy network using different samples from replay memory. Thus we can determine the greatest Q-value output by the network over all actions  $a'_1, \dots, a'_n$ , in  $s'$  [math for maximum q-value given action]. Next we perform gradient descent, typically stochastic gradient descent (SGD), to minimise the loss and update the weights of our DQN [math for weight update given loss]. Alternatively we incorporate a second DQN, called the target network to mitigate over-approximations of maximum expected future reward. Under a single policy network both  $q(s, a)$  and  $q(s', a')$  are approximated using the same weights. This introduces undesirable properties for learning as both sets of network parameters are coupled. As the network weights adjust,  $q(s, a)$  approaches the target value while our approximation of  $q(s', a')$  continues to change. A second DQN allows us to adjust a separate set of weights at arbitrary time steps to match those of the policy network.

We implement our environment as the subset of all reachable states produced by executions of some arbitrary ladder logic program. Thus our action set comprises all propositional variables in  $V$  where a single action entails negating a single boolean.

### 3 Mapping Formal Methods to RL

Machine learning means translating the model to another model - "...Which requires some partial translation, process is expanded in sect 2".

labelled transition systems to MDPs

### 4 Implementation

[action space is one hot encoded binary vector of possible boolean valuations]



### 4.1 LL program generator

### 4.2 Learning Environment

Constructing an environment is arguable the longest phase of training an RL agent. To ensure our agent performs well in real-world applications, the learning environment must be sufficiently representative of the problem domain. Through what is essentially the process of gamification, we design an environment that records agent actions and issues a reward based on the behaviour we wish to reinforce. This process presents a number of distinct challenges in our context. First, given a simple ladder logic program, is it possible to encode its functionality and constraints using an imperative programming language. Second, an efficient method of indexing visited states during exploration to determine a k-number of steps before the agent is forced to revisit states. Finally, we must identify the optimal reward scheme and learning parameters to enforce the desired behaviour. We design an openai-gym [102] like environment to train our agent (Appendix A). Consider the environment as a game. The agent makes successive transitions between states from some arbitrary start position sampled from  $I_s$ . Recall the set  $I$  defined our physical inputs received by the ladder logic program shown in sect. 3.2. At each step the agent is presented with a binary decision to make, selected from our action set  $A$ , where:  $A \equiv I$

Our single input variable pressed is valued either 0 or 1. Once an action is selected, each transition must then be computed as a single execution of the ladder logic program. We define a transition function, Alg. 1, which receives an action  $a_i$ , and the current state as input. Depending on whether the agent selects 0 or 1, the function returns a new state with an updated valuation. These are the transitions  $\sigma : q \rightarrow q'$  defined in sect 3.3 w.r.t finite state automata. Under the assignment rules of our transition function, the agent will eventually discover all six reachable states. Fig. 8.1 represents the finished environment with all paths the agent can explore. Note this is essentially a more complete version of our automaton in sect. 3.3. The formulation of our environment means the agent benefits from having to traverse a directed graph. This can often reduce the number of decisions to consider at each step particularly for large graphs with many edges. Additionally our agent samples from a small action set, meaning fewer potential transitions from a single node. With our transition function generating a total of six states, we continue to implement a collection of traces to record agent exploration. With each action we store the ensuing transition, a 3-tuple  $(q, a_i, q')$ , and the new state  $q'$  in separate hash tables. Our traces let us determine whether the current state or latest transition is unique given we have a record of past steps.

### 4.3 Agent Training

Given we wish to find the maximal depth of the state space without repeating transitions, a positive reward is issued for new discoveries. Contrarily we issue a negative reward for adding previously recorded states with a harsher reward for

repeating transitions. In regard to the exploration-exploitation trade-off we issue a larger reward for discovering new transitions over new states. This is primarily due to agent’s prioritising actions that are guaranteed to issue positive reward while attempting to avoid negative scores. To avoid infinite loops we introduce a terminal condition which resets the environment when all states and transitions have been discovered. To encourage the agent to satisfy the terminal condition as soon as possible we issue a large positive reward once it has been reached.

We can use DQN(s) to approximate the optimal Q-function and, consequently the optimal policy. To this end we implement two DQNs, the first policy network to train our agent and the second target network to approximate our target values. DQNs train differently to Q-tables in their use of replay memory. We implement a basic Replay class to store sequences of 5-tuples  $(s, a, s', r)$ , where  $s$  is the current state,  $a$  denotes the action taken,  $s'$  representing the next state and  $r$  being the reward. We then define a max capacity  $N$ , which dictates the number of experiences our agent remembers. Once the max capacity is reached, the agent pops the oldest entry and adds the latest experience to memory. During the training process our agent samples experiences according to some arbitrary batch size. If the replay memory is less than this batch size, we cannot sample. Therefore we introduce some basic condition checking to determine whether sampling is feasible at any given time step. In regard to the selected exploration strategy,  $\epsilon$ , we employ the same implementation used for Q-tables. Here we set an initial value for  $\epsilon$ , a minimum value and our  $\epsilon$ -decay rate. Concerning our DQN architecture, we construct an input layer of 12 nodes, one for each  $v_i \in V$ , two fully connected hidden layers and a binary output layer to reflect the action set  $A$ . To avoid the problem of moving targets, we set an update rate of 10 epochs for our target network. That is to say we update the weights of the target network for every 10 passes through our policy network. Our chosen optimiser is Adam [103], an SGD variant based on AdaGrad [104] and RMSProp [105], due to its speed and accuracy.

## Network Architecture

## 5 Experiments

### 5.1 Extended state-spaces

The agent learns over time to avoid repeated actions until an optimal trace is output. This trace can then be used to find the maximal depth of the state space without looping for all states in  $I_s$ , if such a path exists. We discover through a series of experiments that there are two sets of unique transitions for all start states such that no transition is repeated twice. Concerning the discovery of loop free paths, as discussed in [17], it may be possible to either explicitly state which path should be verified by the SAT-solver or leverage discoveries made by the agent in identifying a loop free path. By increasing the max step threshold within our environment allows the agent sufficient time to explore the state space and

understand what satisfies the terminal condition. Consequently, provided a long enough training time, the agent will always find an optimal path. In fact it was discovered through a series of experiments, specifically ones with a low  $\epsilon$ -decay rate, that there exist two optimal paths for every start state. These traces can then be analysed to determine the maximum k-steps before looping.

5.2 Interlocking examples

Table 1. Table captions should be placed above the tables.

Heading level	Example	Font size and style
Title (centered)	<b>Lecture Notes</b>	14 point, bold
1st-level heading	<b>1 Introduction</b>	12 point, bold
2nd-level heading	<b>2.1 Printing Area</b>	10 point, bold
3rd-level heading	<b>Run-in Heading in Bold.</b> Text follows	10 point, bold
4th-level heading	<i>Lowest Level Heading.</i> Text follows	10 point, italic

Displayed equations are centered and set on a separate line.

$$x + y = z \tag{9}$$

Please try to avoid rasterized images for line-art diagrams and schemas. Whenever possible, use vector graphics instead (see Fig. ??).

**Theorem 5.** *This is a sample theorem. The run-in heading is set in bold, while the following text appears in italics. Definitions, lemmas, propositions, and corollaries are styled the same way.*

*Proof.* Proofs, examples, and remarks have the initial word in italics, while the following text appears in normal font.

For citations of references, we prefer the use of square brackets and consecutive numbers. Citations using labels or the author/year convention are also acceptable. The following bibliography provides a sample reference list with entries for journal articles [?], an LNCS chapter [?], a book [?], proceedings without editors [?], and a homepage [?]. Multiple citations are grouped [?,?,?], [?,?,?,?].

References

1. Sutton, R.S., Barto, A.G.: Reinforcement learning: An introduction. MIT press (2018)