

# Defending Smart Electrical Power Grids against Cyberattacks with Deep $Q$ -Learning

Mohammadamin Moradi,<sup>1</sup> Yang Weng<sup>1</sup>,<sup>1</sup> and Ying-Cheng Lai<sup>1,2,\*</sup>

<sup>1</sup>*School of Electrical, Computer and Energy Engineering, Arizona State University, Tempe, Arizona 85287, USA*

<sup>2</sup>*Department of Physics, Arizona State University, Tempe, Arizona 85287, USA*

 (Received 6 June 2022; revised 23 September 2022; accepted 10 October 2022; published XX XX 2022)

A key to ensuring the security of smart electrical power grids is to devise and deploy effective defense strategies against cyberattacks. To achieve this goal, an essential task is to simulate and understand the dynamic interplay between the attacker and defender, for which stochastic game theory and reinforcement learning stand out as a powerful mathematical and computational framework. Existing works are based on conventional  $Q$ -learning to find the critical sections of a power grid to choose an effective defense strategy, but the methodology is only applicable to small systems. Additional issues with  $Q$ -learning are the difficulty in considering the timings of cascading failures in the reward function and deterministic modeling of the game, while attack success depends on various parameters and typically has a stochastic nature. Our solution for overcoming these difficulties is to develop a deep  $Q$ -learning-based stochastic zero-sum Nash strategy solution. We demonstrate the workings of our deep  $Q$ -learning solution using the benchmark Wood and Wollenberg 6-bus and the IEEE 30-bus systems; the latter is a relatively large-scale power-grid system that defies the conventional  $Q$ -learning approach. Comparison with alternative reinforcement learning methods provides further support for the general applicability of our deep  $Q$ -learning framework in ensuring secure operation of modern power-grid systems.

DOI: [10.1103/PRXEnergy.0.XXXXXX](https://doi.org/10.1103/PRXEnergy.0.XXXXXX)

## I. INTRODUCTION

Electric power grids, a critical infrastructure, are vulnerable to random failures and, more alarmingly, to hostile physical and/or cyberattacks that can often trigger large-scale cascading types of breakdowns. The US-Canadian blackout in 2003 affected approximately 50 million people in eight US states and two Canadian provinces. In the same year, there were two other significant blackouts in Europe [1]. The gigantic impacted geophysical area of these events and the economic consequences highlight the need for developing effective defense strategies against attacks on the power grids. In the past two decades, research on cybersecurity systems has attracted increasing attention. An important requirement is to make these systems automated and “intelligent,” as many power grids are unmanned and located in isolated, remote, rural, or mountainous areas [2]. In the field of cyberphysical systems and security, the year 2010 was a turning point, when the first ever cyberwarfare

weapon, known as Stuxnet [3], was created. Documented significant events of cyberattacks include a synchronized and coordinated attack in December 2015, which compromised three Ukrainian regional electric power distribution companies and resulted in power outages affecting approximately 225 000 customers for several hours [4]. Due to the extraordinarily large scale and complexity of the power-grid networks, developing effective defense strategies against attacks to prevent breakdown of the networks has become one of the most challenging problems of interdisciplinary research in science and engineering in the present time. In this regard, a pioneering approach is to use state estimation to detect the attack modes to power systems [5,6], assuming that the topology and parameters are known to both the attacker and defender in the transmission grid. Recently, this approach was extended to the distribution grid [7,8]. It is also recognized that attacks are possible, even if the attackers do not know the topology and parameters of the distribution grid [9].

From a general and mathematical point of view, cybersecurity is determined by the dynamic interplay between the attacker and the defender, where the former seeks to maximize, while the latter strives to minimize, damage to the power grid. Game theory [10], a well-established branch of mathematics for analyzing strategic interactions among rational players, thus represents a powerful

\*Ying-Cheng.Lai@asu.edu

*Published by the American Physical Society under the terms of the Creative Commons Attribution 4.0 International license. Further distribution of this work must maintain attribution to the author(s) and the published article's title, journal citation, and DOI.*

66 tool to probe the dynamics of cybersecurity, where the  
 67 attacker-defender interactions can be modeled as a nonco-  
 68 operative game. There are two categories of such games:  
 69 static and dynamic. In a static game, time and informa-  
 70 tion do not affect the action choice of the players, so the  
 71 game can be regarded as a one-shot process, in which  
 72 the players take their actions only once. In contrast, in a  
 73 dynamic game [11], the players have some information  
 74 about each other’s choices and can act more than once,  
 75 where time plays a central role in the decision-making.  
 76 Different game-theoretic techniques have been devised to  
 77 study the security of smart grids, such as the network  
 78 formation game technique used in smart grid communica-  
 79 tions systems, the Nash game and auction game methods  
 80 in demand-side management applications, and coalition  
 81 games used in microgrid distribution networks [12].

82 Recently, machine learning has been introduced to study  
 83 the security of smart power grids. For example, in Ref.  
 84 [13], the most vulnerable areas in a power grid are iden-  
 85 tified using unsupervised learning. Several state-of-the-art  
 86 machine-learning techniques have been devised to gener-  
 87 ate, detect, and mitigate cyberattacks in smart grids [14].  
 88 As one of the most developed machine-learning frame-  
 89 works, reinforcement learning (RL) has proven to be par-  
 90 ticularly useful for cybersecurity systems. Specifically, RL  
 91 is employed to derive false data injection attack policies  
 92 against automatic voltage control systems in power grids  
 93 [15]. In Ref. [16], a RL-based strategy was introduced  
 94 that aimed to choose the appropriate detection interval and  
 95 the number of CPUs allocated based on the defense pref-  
 96 erences through implementation inside the control center  
 97 of the power grid. Moreover,  $Q$ -learning [17] is used to  
 98 analyze the vulnerability of smart grids against sequen-  
 99 tial topological attacks, where the attacker can use  $Q$ -  
 100 learning to worsen the damage of sequential topology  
 101 attacks toward system failures with the least effort [18].  
 102 A fundamental difficulty with  $Q$ -learning is that it can  
 103 become extremely inefficient in the case of increasing  
 104 numbers of state-action pairs, as in a larger power grid. To  
 105 overcome this difficulty, deep RL has been employed in  
 106 large-scale power grids for topology attacks [19]; cyber-  
 107 attack mitigation [20]; and, more recently, to solve the  
 108 latency cyberattack detection problem [21]. In general,  
 109 deep  $Q$ -learning [22] uses neural networks to approximate  
 110 the  $Q$  function using only the state as the input and gener-  
 111 ate the  $Q$  values of all actions as the output. As a result, deep  
 112  $Q$ -learning is suited to problems with a large state-action  
 113 space, since it leverages the extent of deep neural net-  
 114 works to deal with complex cyberphysical systems, such  
 115 as the IEEE 30-bus system. Figure 1 provides a schematic  
 116 comparison of  $Q$ -learning and deep  $Q$ -learning.

117 Here, we develop a deep  $Q$ -learning-based defense strat-  
 118 egy for smart power-grid systems using transmission line  
 119 outages and generation loss as the concrete failure set-  
 120 tings. Broadly, we conceive the scenario in which the

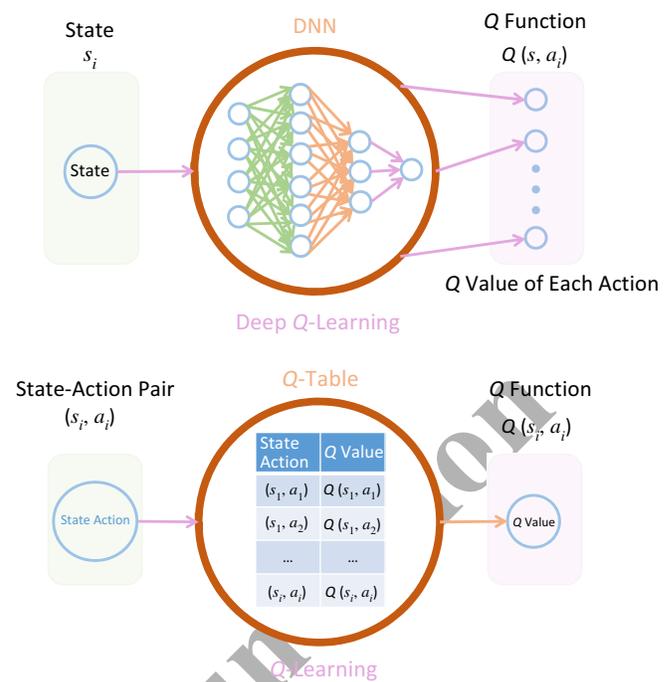


FIG. 1.  $Q$ -Learning versus deep  $Q$ -learning. Implementation of the  $Q$  table is the main difference between  $Q$ -learning and deep  $Q$ -learning. Instead of mapping a state-action pair to a  $Q$  value using the  $Q$  table, as is done in  $Q$ -learning, deep  $Q$ -learning uses neural networks to map the states to the action- $Q$  value pairs—the core reason that deep  $Q$ -learning can be used to solve large-scale problems.

121 defense management of a given large power grid performs  
 122 stochastic game playing to simulate the dynamic inter-  
 123 play between the attacker and the defender. The goal is  
 124 to uncover the “best” attack strategies that can result in  
 125 the maximal damage to the grid. Accordingly, protect-  
 126 ing the components in the grid that such attack strategies  
 127 entail provides the optimal defense tactics. We model the  
 128 attacker-defender interaction as a zero-sum game and solve  
 129 it by using deep  $Q$ -learning, where solving a game entails  
 130 finding its Nash equilibria (see Sec. II B for details). We  
 131 introduce a customized reward function for achieving the  
 132 desired objectives as directly as possible. Importantly, we  
 133 demonstrate that our deep  $Q$ -learning framework can be  
 134 used to address problems of cascading failures and tim-  
 135 ing delays, which, to the best of our knowledge, have  
 136 not been studied previously in the context of machine-  
 137 learning-enhanced or guaranteed security of power grids.  
 138 Our defense algorithm leads to the best protection sets  
 139 based on the defined objectives, taking into considera-  
 140 tion the defender’s policy. To demonstrate the workings  
 141 and advantages of our deep  $Q$ -learning scheme, we com-  
 142 pare its performance not only with the conventional  $Q$ -  
 143 learning method but also with other state-of-the-art algo-  
 144 rithms, such as actor-critic (AC), policy gradient (PG),  
 145 and proximal policy optimization (PPO). Overall, our deep

146  $Q$ -learning approach opens the door to applying RL to  
 147 large-scale smart grid cybersecurity problems to signifi-  
 148 cantly enhance the security of the system in an automated  
 149 manner.

150 The rest of this paper is organized as follows. The RL  
 151 formulation of the attacker-defender stochastic zero-sum  
 152 game, problem description, reward function definition,  
 153 and an illustration of why  $Q$ -learning is not viable for  
 154 large-scale problems are given in Sec. II. In Sec. III, we  
 155 formulate our deep  $Q$ -learning method and present the  
 156 optimal defense strategy. Simulation scenarios and com-  
 157 parative results are detailed in Sec. IV. Section V presents  
 158 a discussion.

## 159 II. REINFORCEMENT-LEARNING-BASED 160 FORMULATION OF ATTACKER-DEFENDER 161 GAME

162 We describe the essential quantities needed for modeling  
 163 the attacker-defender interactions using a stochastic zero-  
 164 sum game and  $Q$ -learning algorithm. We then define the  
 165 reward function based on the objectives of the attack sce-  
 166 narios. The efficiencies of  $Q$ -learning and deep  $Q$ -learning  
 167 are compared using an illustrative example. In the formu-  
 168 lation below, player 1 is the attacker, while player 2 is the  
 169 defender.

### 170 A. Attacker-defender stochastic zero-sum game and 171 Nash equilibrium

172 A game is closely related to a Markov decision process  
 173 that can be viewed as a single-player decision problem, so  
 174 its extension to two players results in a stochastic game  
 175 [23]. Mathematically, a *two-player stochastic zero-sum*  
 176 *game* is a sextuple  $\langle S, A^1, A^2, r^1, r^2, p \rangle$ , where  $S$  is the dis-  
 177 crete state space,  $A^i$  is the discrete action space of player  $i$   
 178 (for  $i = 1, 2$ ),  $r^i: S \times A^1 \times A^2 \rightarrow \mathbb{R}$  is the payoff function  
 179 for player  $i$ , whereas  $r^1(s, a^1, a^2) = -r^2(s, a^1, a^2)$  for all  
 180  $s \in S, a^1 \in A^1, a^2 \in A^2$ . For the cases studied in this work,  
 181 intuitively, rewards are the game payoffs that are either the  
 182 generation loss caused by the attacks or a function of the  
 183 transmission line outages [cf., Eq. (10) below]. Moreover,  
 184  $p: S \times A^1 \times A^2 \rightarrow \Delta(S)$  is the transition probability map-  
 185 ping, with  $\Delta(S)$  being the set of probability distributions  
 186 over the state space,  $S$ . During a game, player 1 aims to  
 187 maximize, but player 2 strives to minimize, the sum of the  
 188 discounted rewards. Given an initial state  $s$ , discount fac-  
 189 tor  $\gamma$ , and  $\pi^1$  and  $\pi^2$  (the strategies of players 1 and 2,  
 190 respectively), the values of the game for the two players  
 191 are

$$192 \quad v^1(s, \pi^1, \pi^2) = \sum_{t=0}^{\infty} \gamma^t \mathbb{E}\{r_t^1 | \pi^1, \pi^2, s_0 = s\}, \quad (1)$$

$$193 \quad v^2(s, \pi^1, \pi^2) = \sum_{t=0}^{\infty} \gamma^t \mathbb{E}\{r_t^2 | \pi^1, \pi^2, s_0 = s\}, \quad (2)$$

194 where  $\pi^{1,2} = (\pi_0^{1,2}, \dots, \pi_t^{1,2}, \dots)$ , with  $\pi_t^{1,2}$  denoting the  
 195 decision rules of players 1 and 2 at time  $t$  and  $\mathbb{E}\{\cdot\}$  is the  
 196 conditional expectation. For instance,  $\mathbb{E}\{r_t^i | \pi^1, \pi^2, s_0 = s\}$   
 197 is the expectation of the player  $i$ 's instant reward at time  
 198  $t$ , following the decision rules  $\pi^{1,2}$  with  $s$  as the initial  
 199 state. These strategies are “stationary,” in the sense that the  
 200 decision rules are fixed over time, in contrast to the “behav-  
 201 ior” strategies often used in economics, where the decision  
 202 rules depend on the history of states and the actions up  
 203 to the present time. Assuming each player has complete  
 204 information about the reward function of the other player, a  
 205 Nash equilibrium can emerge. Specifically, the *Nash equi-*  
 206 *librium for a two-player stochastic zero-sum game* is a pair  
 207 of strategies,  $(\pi_*^1, \pi_*^2)$ , such that for all  $s \in S$ , the following  
 208 hold:

$$209 \quad v^1(s, \pi_*^1, \pi_*^2) \geq v^1(s, \pi^1, \pi_*^2) \quad \forall \pi^1 \in \Pi^1, \quad (3)$$

$$210 \quad v^2(s, \pi_*^1, \pi_*^2) \geq v^2(s, \pi_*^1, \pi^2) \quad \forall \pi^2 \in \Pi^2, \quad (4)$$

211 where  $\Pi^i$  is the set of all possible policies for player  $i$ .  
 212 Intuitively, a Nash equilibrium means that each player's  
 213 strategy is the best response to the other player's strategy:  
 214 neither one has anything to gain by changing only their  
 215 own strategy.  
 216

217 In general, based on the structure of the information that  
 218 the players possess, attacker-defender stochastic zero-sum  
 219 games can be classified into four categories, depending on  
 220 whether the information is complete or incomplete, per-  
 221 fect or imperfect. In particular, in a complete information  
 222 game, the players know the structure of the game being  
 223 played, such as the number of players and their payoff  
 224 functions. Any missing information will lead to an incom-  
 225 plete information game. In addition, a game is regarded  
 226 as being of the perfect information type if all the players  
 227 know the historical actions of each other at the time of their  
 228 move; otherwise, the game is of the imperfect information  
 229 type [24]. For simplicity, in our work, we assume both the  
 230 attacker and defender can observe each other's immedi-  
 231 ate reward and have access to their actions throughout the  
 232 learning process. This assumption, while ideal and offering  
 233 mathematical convenience, is based on the consideration  
 234 that the goal of our work is to solve the attacker-defender  
 235 stochastic zero-sum game for defensive planning. In fact,  
 236 our aim is to find the best scenario for the attacker, so  
 237 that the defender can be prepared for the worst, and thus,  
 238 assuming the availability of complete information may not  
 239 be unreasonable. Possible scenarios to obtain the required  
 240 information include the observation of the state of the  
 241 transmission lines by the defender, the defender's access  
 242 to the resulting generation loss when an attack happens,  
 243 and some insider information about the defender obtained  
 244 by the attacker.

## B. Q-Learning-based solution to attacker-defender stochastic zero-sum game

Reinforcement learning belongs to the field of decision-making, where the “agent” explores the “environment,” interacts with it, and observes its reactions to find an optimal behavior to maximize a long-term “reward.” Contrary to supervised learning, in RL, the agent must act independently to find an optimal sequence of actions that maximizes a given reward function in an unknown environment.

While RL is capable of directly solving certain cybersecurity problems, it can also serve as a powerful vehicle to gain insights into the attacker-defender interactions modeled as a game. In general, solving a game means finding its Nash equilibria. Especially, an appealing feature of RL is that it can yield solutions (Nash equilibria) of both the attacker-defender interplay and cybersecurity in a knowledge-free manner, i.e., based solely on data. For example, the Nash equilibrium for the two-player zero-sum game can be determined online based on RL [25]. RL has also been employed to solve a zero-sum stochastic game [26]. The min-max solutions of a dynamic Markov zero-sum game are derived using  $Q$ -learning [27], yielding optimal risk management strategies to meet the performance criteria with the parameters of the Markov game model completely unknown. A distributed RL algorithm is proposed to solve a non-zero-sum stochastic game in which each player needs only minimal information about the other player [28]. RL is also used in a stochastic adversarial game coupled with an expert advice framework to analyze the optimal attack strategies against predictors [29]. While game theory has been applied to many problems that require rational decision-making, there are some limitations in applying such methods to security games.  $Q$ -Learning was introduced to secure the system by devising proper actions against the adversarial behavior of a suspicious user [30].  $Q$ -Learning has also been employed in solving security games, as studied in Refs. [31,32].

In  $Q$ -learning, the  $Q$  function is a mapping of all possible state-action pairs (where actions refer to action profiles of each player) to a scalar value and represents the total discounted reward that a player is expected to obtain, starting from a determined state taking a specified action. For a two-player stochastic game, the optimal  $Q$  function for each player can be defined as

$$Q_*^1(s, a^1, a^2) = r^1(s, a^1, a^2) + \gamma \sum_{s'=1}^N p(s'|s, a^1, a^2) v^1(s', \pi^1, \pi^2), \quad (5)$$

$$Q_*^2(s, a^1, a^2) = r^2(s, a^1, a^2) + \gamma \sum_{s'=1}^N p(s'|s, a^1, a^2) v^2(s', \pi^1, \pi^2), \quad (6)$$

where  $s'$  is the next state evolving from state  $s$  taking actions  $a^1$  and  $a^2$ . Equations (5) and (6) define  $Q_*$ , the optimal value of the  $Q$  function associated with state  $s$  and action pair  $(a^1, a^2)$ . For each player, the optimal value is equal to the total discounted reward received by the player, when both the attacker and defender perform actions  $(a^1, a^2)$  in state  $s$  and subsequently follow their Nash equilibrium strategies  $(\pi^1, \pi^2)$ . For each player, the value of  $Q_*$  can be solved [Eq. (8)]. A player then generates a policy by following the action with the largest  $Q$  value in each state.

We remark that, in the reinforcement learning literature, the notation  $r$  is usually reserved for “instant reward” or “instant payoff,” whereas  $v$  is the “value function.” In Eq. (5), the term  $r^1(s, a^1, a^2)$  means the instant payoff that player 1 gets when the game is in state  $s$  and player 1 chooses action  $a^1$  while player 2 selects action  $a^2$ . The quantity  $v^1(s', \pi^1, \pi^2)$  denotes the total discounted payoff starting from the next state  $s'$  while the players follow the policies  $\pi^1$  and  $\pi^2$ . Thus,  $Q_*^1(s, a^1, a^2)$  in Eq. (5) represents the instant reward added to the best possible future rewards for player 1. Intuitively, this means the best reward player 1 can achieve starting from state  $s$  with the two players taking actions  $a^1$  and  $a^2$ , respectively.

Because of the zero-sum nature of the game,  $Q_*^1(s, a^1, a^2) + Q_*^2(s, a^1, a^2) = 0$ , or

$$Q_*^1(s, a^1, a^2) = -Q_*^2(s, a^1, a^2), \quad (7)$$

the learning agent needs to learn (or approximate) only one  $Q$  function. This should be contrasted with a general sum game characterized by  $Q_*^1(s, a^1, a^2) + Q_*^2(s, a^1, a^2) \neq 0$ , where two  $Q$  functions need to be learned, increasing substantially the computation complexity. To solve Eqs. (5) and (6), we use the following algorithm [23]:

$$Q_{t+1}(s, a^1, a^2) = (1 - \alpha_t) Q_t(s, a^1, a^2) + \alpha_t \left[ r_t + \gamma \max_{\pi^1(s') \in \sigma(A^1)} \min_{\pi^2(s') \in \sigma(A^2)} \pi^1(s') Q_t(s') \pi^2(s') \right], \quad (8)$$

where  $Q_{t+1}(s, a^1, a^2) = Q_{t+1}^1(s, a^1, a^2)$ . Convergence requires that all state-action pairs be visited infinitely often, which is practically infeasible. To obtain a reasonable functional approximation, a sufficiently large state-action space needs to be explored. This is the main reason that prevents  $Q$ -learning from being applicable to large-scale smart grids.

## C. Transmission line outage, generation loss, and reward functions

We focus on two representative attack scenarios on smart power grids [33–35]. The first is the switching line

340 problem, where the attacker attempts to cause a predeter-  
 341 mined percentage of the transmission lines to go down. In  
 342 the second scenario, the attacker attempts to maximize the  
 343 generation loss in the power system through a sequence of  
 344 attacks. In both cases, the defender strives to mitigate the  
 345 attack consequences, regardless of whether they are due to  
 346 transmission line outages or are caused by generation loss.  
 347 [We use a dc load flow simulator of cascading (separation)  
 348 in power systems, named DCSIMSEP [33,34], to calculate  
 349 the generation loss.] The state space for both attacks is  
 350 the state of transmission lines denoted as a  $l \times 1$  binary-  
 351 valued vector, where  $l$  is the number of transmission lines;  
 352 this value for each transmission line is 0 if the respective  
 353 line is down and is 1 otherwise. The attacker's actions for  
 354 both attacks are chosen from the set  $A = \{1, 2, 3, \dots, l\}$ ,  
 355 where action  $i$  means attacking transmission line  $i$ . The  
 356 defender's action for both attacks is considered to be a  
 357 set consisting of  $n$  transmission lines, denoted as the pro-  
 358 tection set. The attacker's reward for the line switching  
 359 attack is given by Eq. (10) and for the generation loss  
 360 attack is the average generation loss [Eq. (9)] caused by  
 361 the attack. Since the game is considered to be zero sum,  
 362 for the defender, the payoff is the negative of the attacker's  
 363 reward for both attacks. The transition probability distri-  
 364 bution is represented with power-grid transitions simulated  
 365 with the DCSIMSEP tool.

366 We incorporate the cascading failure timing into the  
 367 reward function. We assume that the attacker's next attack  
 368 will be launched at time  $T = 1.2t_{\text{cas}}$ , where  $t_{\text{cas}}$  is the  
 369 cascading failure length caused by the attacks. The propor-  
 370 tional constant 1.2 is chosen somewhat arbitrarily, insofar  
 371 as it is greater than 1, so that the system settles into a  
 372 steady state after an attack on the transmission lines. The  
 373 choice of the value  $T$  does not have a significant effect  
 374 because the generation loss is relative among different  
 375 attacks and our goal is to minimize the total loss. To take  
 376 into account the timing delays of the cascading failures, we  
 377 use a weighted average of generation loss during a reason-  
 378 able time interval. Specifically, the average generation loss  
 379  $G_{\text{loss}}^-$  is

$$380 \quad G_{\text{loss}}^- = G_{\text{loss}}^{\text{init}} \frac{t_{\text{cas}}}{T} + G_{\text{loss}}^{\text{stead}} \frac{T - t_{\text{cas}}}{T}, \quad (9)$$

381 where  $G_{\text{loss}}^{\text{init}}$  is the generation loss caused initially by the  
 382 attack, while  $G_{\text{loss}}^{\text{stead}}$  represents the generation loss during  
 383 the steady state of the system after a transient phase caused  
 384 by the attack. The reason is that, after an attack, the power  
 385 grid will enter into a transient state, during which cascad-  
 386 ing failures occur. We assume that the defender's policy  
 387 is passive while the attacker's policy evolves according to  
 388 deep  $Q$ -learning (as described in Sec. II D). The defender's  
 389 protection set is updated at the end of each run, mean-  
 390 ing that the attacker must learn the optimal sequences in  
 391 a constantly updated environment. In general, the defender

is not able to protect all lines simultaneously because of  
 limited resources. This highlights the need for  $Q$ -learning  
 because the defender should wisely select the set of lines  
 to protect.

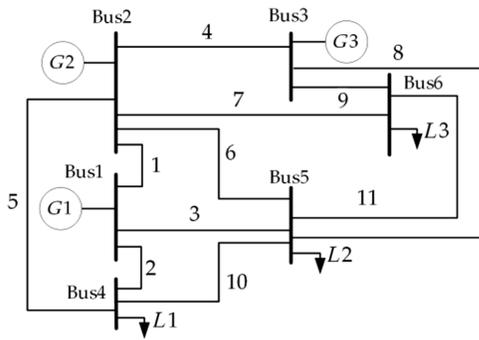
For the first attack scenario, the reward function is given  
 by

$$\begin{aligned} r &= r_1, & \text{for } \text{IO} > \text{AO}, \\ r &= r_2, & \text{if attack is final,} \\ r &= \text{IO}/\text{AO}, & \text{otherwise,} \end{aligned} \quad (10)$$

where IO is the instant number of transmission line out-  
 ages caused by the attack, AO is the attack objective,  
 and  $r_1 > r_2$ . For example, in the Wood and Wollenberg  
 (W&W) 6-bus system shown in Fig. 2, when the protec-  
 tion set consists of lines 1 and 2, attacking line 5 will cause  
 an instant outage of five lines ( $\text{IO} = 5$ ), which is more than  
 the attack objective ( $\text{AO} = 4$ ). In this case, the reward of  
 attacking line 5 is equal to  $r_1$ . This is the best scenario, and  
 therefore,  $r_1$  is chosen to be large enough to persuade the  
 agent to learn this action, if possible. This will also lead to  
 $G_{\text{loss}}^{\text{init}} = 210$  MW and  $G_{\text{loss}}^{\text{stead}} = 83.5$  MW, and the cascad-  
 ing failure length is  $t_{\text{cas}} = 331.61$  s. The cascading failure  
 timing delays caused by attacking line 5 in the W&W 6-  
 bus system are illustrated in Fig. 3. Equation (9) provides  
 the average generation loss, taking into account the timing  
 delay of cascading failures as  $G_{\text{loss}}^- = 167.83$  MW. Like-  
 wise, attacking line 3 will cause lines 1, 2, and 3 to go  
 down, leading to the reward  $r = 3/4$ . Eventually, if the  
 number of currently downed transmission lines is less than  
 AO, but an attack causes the number of downed lines to be  
 equal to or larger than AO, the attacker will have achieved  
 the objective in this specific step, executing the chosen  
 action. In this case, the attack is called final and the reward  
 is  $r_2$ , as the attacking agent is motivated to take the final  
 blow when an opportunity rises.

#### D. Necessity of deep $Q$ -learning

A standard way to implement  $Q$ -learning is through the  
 sample base variant called "tabular  $Q$ -learning." In a  $Q$   
 table, the rows list the states of the underlying system,  
 and the columns are indexed by the action set. Training  
 the table is helpful in finding an optimal action for each  
 state with the goal of maximizing the long-term reward.  
 This is a straightforward yet powerful approach to the  
 security of small cyberphysical systems. For example, a  
 one-shot game with a multiline switching attack between  
 the attacker and defender in a smart grid was studied  
 [36]. In another work [37], the dynamics of the electric  
 power grid were taken into account and the attacks were  
 modeled as a multistage game, where the percentage of  
 visited states with respect to the total number of states  
 was 1.81% for the W&W 6-bus system (37 states out of



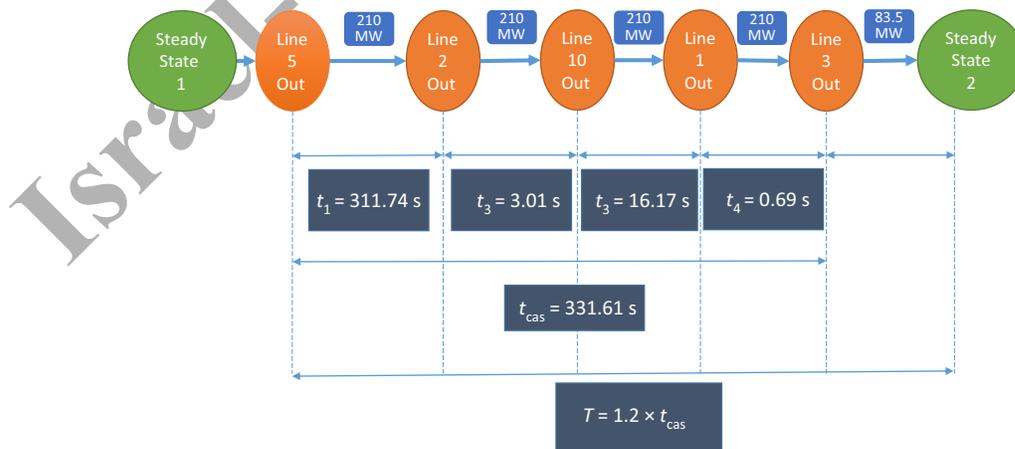
F2:1 FIG. 2. Wood and Wollenberg 6-bus system. It has 6 buses, 3  
 F2:2 generators (denoted by  $G$ ), 3 loads (denoted by  $L$ ), and 11 trans-  
 F2:3 mission lines. IEEE 30-bus system simulated in this paper has  
 F2:4 a similar topological structure but at a much larger scale: it has  
 F2:5 6 generators, 30 buses, and 41 transmission lines. Simulation of  
 F2:6 the smart power grids (they are “smart” because they support  
 F2:7 renewable sources) is performed using the DCSIMSEP package, a  
 F2:8 simulator of cascading failures in power systems. DCSIMSEP does  
 F2:9 not use any specific stress-mitigating controls under the assump-  
 F2:10 tion that the cascades are propagating too fast for the operators  
 F2:11 to react, so it is suitable for cyberattack problems.

441 a possible  $2^{11}$  states) and  $1.87 \times 10^{-8}\%$  for the IEEE 39-  
 442 bus system (13 130 states out of a possible  $2^{46}$  states).  
 443 The tabular  $Q$ -learning method is thus incapable of suffi-  
 444 cient state-space exploration, leading to suboptimal poli-  
 445 cies for the given reward functions. In general, for larger  
 446 power-grid systems, such as the benchmark IEEE 30-bus  
 447 system that has 41 transmission lines, tabular  $Q$ -learning is  
 448 impractical. This is because each line has two states, opera-  
 449 tional or out of service, so there are  $2^{41}$  number of states for  
 450 all the transmission lines. If only a single line is attacked,  
 451 the total number of actions is 41. Because there are  $2^{41}$

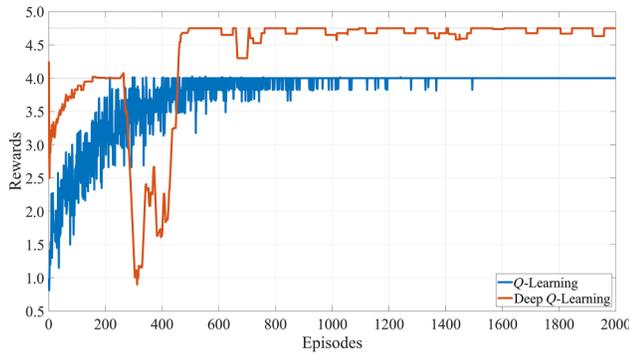
452 states for each action, the table will have  $2^{41} \times 41$  cells, 452  
 453 rendering infeasible any computation based on the table. 453

454 To appreciate the necessity of adopting deep  $Q$ -learning 454  
 455 in tackling the cybersecurity problem of smart power- 455  
 456 grid systems in a concrete way, we use the switching line 456  
 457 problem as a prototypical example. For the W&W 6-bus 457  
 458 system, consider the specific formulation in which AO is 458  
 459 4, the protection set is  $[1, 2]$ , the maximum number of 459  
 460 attacks is 4, and the reward function is given by Eq. (10) 460  
 461 with  $r_1 = 4$  and  $r_2 = 1$ . The optimal attacking sequence 461  
 462 derived using  $Q$ -learning after 20 independent runs (each 462  
 463 with 2000 episodes) is to attack line 5, which will lead to 463  
 464 a maximum reward of 4. However, the optimal attacking 464  
 465 sequence derived using deep  $Q$ -learning is to attack line 465  
 466 9, then line 8, and finally line 6. In particular, the outage 466  
 467 of line 9 will lead to reward  $r = 0.25$ ; attacking line 8 will 467  
 468 bring down lines 8 and 4 together, so the reward is  $r = 0.5$ ; 468  
 469 and disabling line 6 will cause lines 1, 2, 3, 6, 10, and 11 to 469  
 470 go down, leading to the reward  $r = 4$ . As a result, the deep 470  
 471  $Q$ -learning strategy will result in a total reward of 4.75. A 471  
 472 detailed comparison of the rewards achieved as a function 472  
 473 of time from executing the optimal attack strategies from 473  
 474  $Q$ -learning and deep  $Q$ -learning is shown in Fig. 4. It can 474  
 475 be seen that, while there is a brief time period (between 200 475  
 476 and 500 episodes of the game) in which the reward of  $Q$ - 476  
 477 learning is greater than that of deep  $Q$ -learning, after 500 477  
 478 episodes, deep  $Q$ -learning leads to a persistently higher 478  
 479 reward than  $Q$ -learning. 479

480 The main reason that the tabular  $Q$ -learning results in 480  
 481 lower reward in the long run lies in insufficient state-space 481  
 482 exploration, generating a suboptimal policy for the defined 482  
 483 reward function. In a larger power grid, such as the IEEE 483  
 484 30-bus system that has 41 transmission lines, there are  $2^{41}$  484  
 485 distinct states. Practically, a state space of this large size 485  
 486 cannot be solved using conventional tabular  $Q$ -learning 486



F3:1 FIG. 3. Cascading failure timing delays caused by attacking line 5 in the W&W 6-bus system derived using DCSIMSEP package.  
 F3:2 Average generation loss ( $G_{loss}$ ) caused by this attack can be calculated using these timings in Eq. (9).



F4:1 FIG. 4. Comparison of the performance of deep  $Q$ -learning  
 F4:2 and conventional tabular  $Q$ -learning using a concrete exam-  
 F4:3 ple. Setting is the switching line problem in the W&W 6-bus  
 F4:4 system. Shown are the values of reward function [Eq. (10)]  
 F4:5 with  $r_1 = 4$  and  $r_2 = 1$ ] from deep  $Q$ -learning and conventional  
 F4:6  $Q$ -learning with similar simulation parameter values. Deep  $Q$ -  
 F4:7 learning algorithm manages to find an optimal attack sequence,  
 F4:8 which results in the reward of  $r = 4.75$ , while conventional  $Q$ -  
 F4:9 learning is unable to find a sequence with a reward of larger than  
 F4:10  $r = 4$ .

487 [38]. This difficulty with  $Q$ -learning is fundamental. As  
 488 the system becomes larger, the deficiency of  $Q$ -learning  
 489 will become more apparent and pronounced. To address  
 490 the cyberattack and defense problem for large-scale power  
 491 grids, invoking deep  $Q$ -learning is necessary.

### 492 III. DEEP $Q$ -LEARNING-BASED FORMULATION 493 OF ATTACKER-DEFENDER GAME

494 We introduce the deep  $Q$ -learning algorithm and exploit  
 495 it to formulate and solve the attacker-defender stochastic  
 496 zero-sum game problem. We also analyze the proposed  
 497 defense strategy for smart power grids against cyberat-  
 498 tacks. The zero-sum nature of the game dynamics stip-  
 499 ulates that the deep  $Q$ -learning agent needs to learn (or  
 500 approximate) only one  $Q$  function. It should be noted  
 501 that, mathematically, convergence to a Nash equilibrium  
 502 requires that all state-action pairs be visited infinitely often,  
 503 which is practically infeasible. To obtain a reasonable  
 504 functional approximation, a sufficiently large state-action  
 505 space needs to be explored, which can be accomplished by  
 506 deep  $Q$ -learning.

#### 507 A. Deep $Q$ -learning solution to attacker-defender 508 stochastic zero-sum game

509 The core of deep  $Q$ -learning is an online multilayered  
 510 neural network [39] that for any given state  $s$  outputs a  
 511 vector of action values  $Q(s, \cdot, \cdot; \theta)$ , where  $\theta$  denotes the  
 512 set of parameters of the online network. Two foundations  
 513 of the deep  $Q$ -learning algorithm are the target network  
 514 and the use of experience replay. The target network, with  
 515 parameter set  $\theta^*$ , is the same as the online network, except

that, for every  $c$  episodes, its parameters are copied from 516  
 the online network,  $\theta_t^* = \theta_t$ , which are kept fixed during 517  
 the  $c$  episodes. The target used by deep  $Q$ -learning can be 518  
 described as 519

$$Q_t^* = r_{t+1} + \gamma \max_a Q_t(s_{t+1}, a^1, a^2; \theta_t^*). \quad (11) \quad 520$$

The deep  $Q$ -learning agent gets the initial state and com- 521  
 puts the  $Q$ -function values for all possible actions, which 522  
 in our problem is the transmission lines of the power 523  
 grid. We use the epsilon greedy method [40] to select 524  
 a proper action, where the action with the largest  $Q$ - 525  
 function value is chosen with the probability of  $1 - \epsilon$ , and 526  
 a random action is performed with the probability of  $\epsilon$ . 527  
 The state, attacker, and defender's actions; the next state 528  
 derived from the stochastic transition function; and the 529  
 gained reward are stored for some time. These data are 530  
 then sampled uniformly from this memory bank to update 531  
 the network, which is called experience replay, as some 532  
 random batches of transition are sampled. The difference 533  
 between the target  $Q$  function and the predicted  $Q$  function 534  
 is calculated as 535

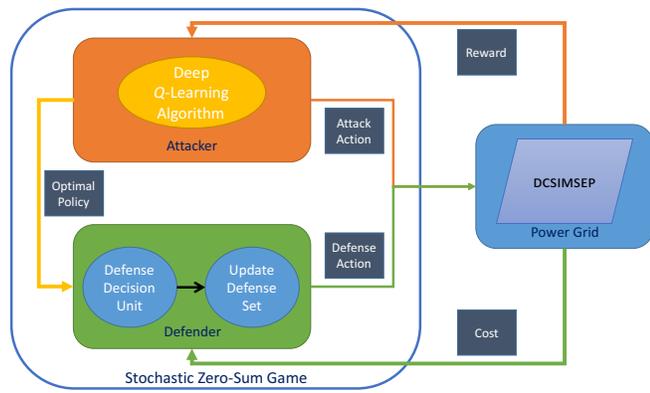
$$\text{error} = Q_t^* - Q_t(s_{t+1}, a^1, a^2; \theta_t), \quad (12) \quad 536$$

where a small error indicates a well-trained algorithm. 537  
 Typically, a gradient descent algorithm can be used to opti- 538  
 mize the online network parameter values to minimize the 539  
 error. The target network's parameters are updated peri- 540  
 odically to match the ones of the online network. Both 541  
 the target network and experience replay can dramatically 542  
 improve the performance of the algorithm [38]. Using the 543  
 $Q$  functions defined in Eqs. (5) and (6) for the stochas- 544  
 tic zero-sum game, we determine the optimal attacking 545  
 sequence so that the defender can choose the best defense 546  
 strategy. 547

The main difference between  $Q$ -learning and deep  $Q$ - 548  
 learning lies in the implementation of the  $Q$  table. In a 549  
 problem with a large number of state-action pairs, the  $Q$  550  
 table becomes unmanageably large and impractical. This is 551  
 because the greater the number of rows and columns, the 552  
 more time it requires for the agents to explore the states 553  
 and to update their values. In deep  $Q$ -learning, the idea is 554  
 that, rather than mapping a state-action pair to a  $Q$  value 555  
 using the  $Q$  table, neural networks can be exploited to 556  
 map the states to the action- $Q$ -value pairs. That is, instead 557  
 of visiting different state-action pairs and filling in the  $Q$  558  
 table, a deep neural network is trained to approximate the 559  
 $Q$  function. 560

#### 561 B. Defensive strategy algorithm using deep $Q$ -learning

Figure 5 presents the proposed algorithm for articulat- 562  
 ing a defense strategy to protect a smart power grid from 563  
 cyberattacks. The attacker and defender play a stochastic 564



F5:1 FIG. 5. Defensive strategy algorithm based on deep  $Q$ -  
 F5:2 learning in a stochastic zero-sum game. Attacker and defender  
 F5:3 are the two players of this game. Attacker uses the deep  $Q$ -  
 F5:4 learning algorithm to find an optimal attack sequence to maxi-  
 F5:5 mize the generation loss or transmission line outage, while the  
 F5:6 defender updates its defense set based on the attacker's previous  
 F5:7 policy. Chosen actions of both players are given to the DCSIM-  
 F5:8 SEP power flow simulator and the reward (cost) is then calculated  
 F5:9 and returned to the players. Process continues until the defender's  
 F5:10 protection set remains unchanged for a number of cycles.

565 zero-sum game with the defined objective of disabling a  
 566 fixed number of transmission lines or maximizing (mini-  
 567 mizing) the generation loss. The attacker attacks the power  
 568 system while the defender protects some transmission  
 569 lines. The payoff, which is either the generation loss or  
 570 the number of downed transmission lines, is determined  
 571 using DCSIMSEP based on the players' actions. Both players  
 572 receive the reward for (cost of) their actions. The attacker  
 573 uses deep  $Q$ -learning to optimize the attack sequence.  
 574 Once an optimal attacking strategy is reached, it is trans-  
 575 mitted to the defender. The defense decision management  
 576 unit will decide whether or not to update the protection set.  
 577 More specifically, the decision unit will simply update the  
 578 protection set with the sweet targets of the previous learn-  
 579 ing process, which are the transmission lines that have the  
 580 largest  $Q$ -function value for the current state. The defense  
 581 decision unit will not update the protection set in the case  
 582 of periodic alternation of sweet targets, which is the indi-  
 583 cator of convergence of the algorithm. This procedure  
 584 continues until a Nash equilibrium (equilibria) is reached.

#### 585 IV. RESULTS

586 To demonstrate the workings and power of our deep  
 587  $Q$ -learning algorithm in generating optimal defense strate-  
 588 gies against attacks, we use the benchmark W&W 6-bus  
 589 and IEEE 30-bus systems. Specifically, for the relatively  
 590 small W&W 6-bus system, the generation loss problem is  
 591 studied in more detail with physical insights. For the larger  
 592 IEEE 30-bus system, we focus on both the switching line  
 593 (transmission line outage) and the maximum generation

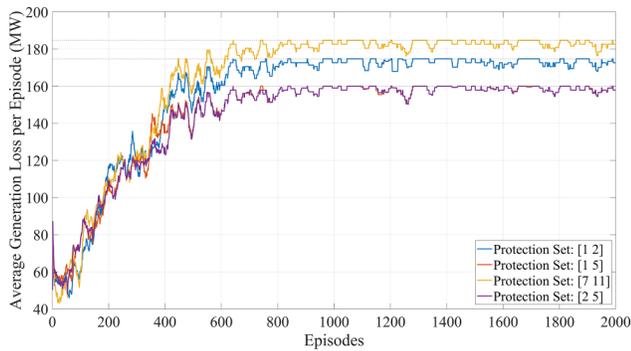
TABLE I. Simulation parameters for W&W 6-bus system generation loss and IEEE 30-bus system generation loss and switching line problems.

Parameters	W&W6 gen	IEEE30 switch	IEEE30 gen
Trans. lines	11	41	41
Episodes	2e3	2e3	1e4
Attack length	5	4	5
Epsilon	1	1	1
Eps. decay	0.005	0.0008	0.005
Eps. min	0.01	0.001	0.01
Learn. rate	0.001	0.001	0.001
Disc. factor	0.7	0.7	0.8
Minibatch size	256	1024	256
FF. neurons	100	200	200
Attack succ. prob.	0.8	0.9	0.9

594 loss problems. All the simulations are carried out using  
 595 the MATLAB R2021b reinforcement learning toolbox on a  
 596 desktop PC with an Intel Core i7-6850K CPU and 128  
 597 GB of RAM. Table I lists the simulation parameter val-  
 598 ues for each problem. In our simulations, we assume that  
 599 an attack on a specific line is successful with a preassigned  
 600 probability that depends on the defender's protection set,  
 601 which is updated after the attacker's learning process. For  
 602 example, in the W&W 6-bus system, suppose the defender  
 603 protects line 5. If the attacker attacks any line other than  
 604 5, the probability of that line's outage will be  $p$ . How-  
 605 ever, if the attacker attacks line 5, it will not go down,  
 606 since the defender protects it, but failures can occur with  
 607 the same probability  $p$ . The value of  $p$  may depend on  
 608 the available resources allocated to the defender or the  
 609 attacker at each time step. During the dynamic interplay  
 610 between the attacker and defender, the value of  $p$  is treated  
 611 as a constant. The reason lies in the tacit assumption that  
 612 both sides have equal access to the resources, so assigning  
 613 extra resources to any specific transmission line is disal-  
 614 lowed. It is worth noting that deep  $Q$ -learning generally  
 615 runs much faster than the equivalent  $Q$ -learning algorithm  
 616 on a per episode basis, because the computation complex-  
 617 ity of deep  $Q$ -learning can be significantly reduced when  
 618 neural networks are used instead of a table, as in con-  
 619 ventional  $Q$ -learning. In all cases, the core of our deep  
 620  $Q$ -learning system is a neural network consisting of two  
 621 fully connected and two ReLu layers.

#### 622 A. Optimal defense strategy for W&W 6-bus system 623 against generation loss

624 We study the maximum generation loss problem, a  
 625 stochastic zero-sum game in which the attacker aims to  
 626 maximize, but the defender aims to minimize, the  
 627 generation loss caused by the attacks, with probabilistic  
 628 state transitions. The attacker's reward at each step is equal  
 629 to  $G_{\text{loss}}^-$  defined in Eq. (9). The zero-sum nature of the  
 629



F6:1 FIG. 6. Effect of choosing an effective protection set in the  
 F6:2 worst-case scenario of generation loss in the W&W 6-bus sys-  
 F6:3 tem. Attacker uses deep  $Q$ -learning to find an optimal attack  
 F6:4 sequence, while the defender updates its protection set accord-  
 F6:5 ing to the attacker's policy. Starting from a random protection  
 F6:6 set {7, 11}, the defender finds the optimal defense set to be {2, 5},  
 F6:7 which causes the worst-case scenario of the generation loss to be  
 F6:8 reduced by %13.41.

630 game dynamics stipulates that the defender's reward must  
 631 be  $-G_{\text{loss}}$ . To be concrete, we assume that the defender is  
 632 able to defend two lines at a time, while the attacker can  
 633 attack up to five lines in a sequential manner. The spec-  
 634 ific numbers can be chosen arbitrarily. Figure 6 depicts  
 635  $G_{\text{loss}}$  per episode for different protection sets. First, for a  
 636 random protection set {7, 11}, we apply deep  $Q$ -learning  
 637 to find the attacker's sweet targets, the transmission lines  
 638 that have the largest  $Q$ -function value for the initial state.  
 639 From the specific random protection set, the sweet targets  
 640 are determined to be lines 1 and 2, so the protection set is  
 641 updated to lines {1, 2}. We apply deep  $Q$ -learning again,  
 642 resulting in lines 1 and 5 becoming the updated sweet tar-  
 643 gets. For the protection set {1, 5}, the new sweet targets  
 644 are lines 2 and 5. Further steps of the game plan will result in  
 645 a Nash equilibrium of 159.93 MW generation loss, alter-  
 646 nating between the protection sets {1, 5} and {2, 5}, which  
 647 represent the solution of the optimal defense sets to this  
 648 problem. Intuitively, the derived sequence of the attacker's  
 649 actions and the protection set constituting a Nash equi-  
 650 librium can be interpreted as pairs of actions from which  
 651 neither the attacker nor the defender is inclined to deviate  
 652 unilaterally. As shown in Fig. 6, this optimal choice of the  
 653 protection set results in a 13.41% decrease in the worst-  
 654 case scenario of generation loss where the attacker plays  
 655 the optimal sequence strategy.

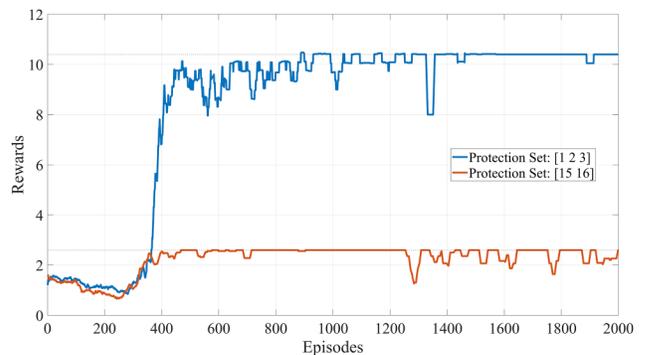
### 656 B. Optimal defense strategy for IEEE 30-bus system 657 against attacks on switching lines

658 In the switching line problem, the attacker has a fixed  
 659 objective of disabling a specific set of transmission lines.  
 660 Our concrete setting is that the defender is able to defend  
 661 up to three lines at a time, while the attacker can attack

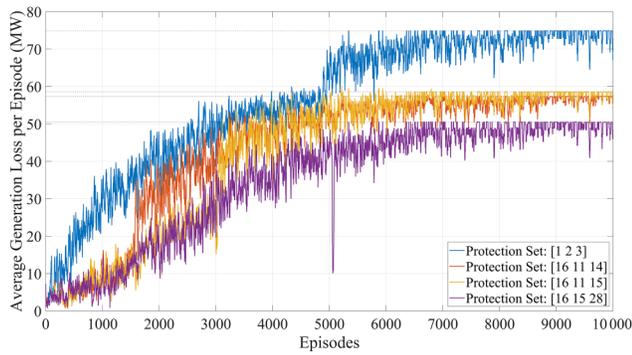
up to four lines sequentially with the AO set to five lines. 662  
 The reward function is given by Eq. (10) with  $r_1 = 10$  663  
 and  $r_2 = 1$ . Starting with a random protection set {1, 2, 3}, 664  
 we apply our deep  $Q$ -learning algorithm and identify the 665  
 sweet targets as lines 15 and 16. The protection set is then 666  
 updated to {15, 16}, and the worst-case scenario reward is 667  
 decreased significantly, as shown in Fig. 7. Further gam- 668  
 ing steps result in the protection set {15, 16} as the Nash 669  
 equilibrium. The intuitive reason is that, when protecting 670  
 lines {15, 16}, the attacker is not able to find a sequence 671  
 that will result in a large instantaneous outage. As a result, 672  
 the attack receives a much smaller reward compared to the 673  
 case when the defender defends a random protection set. 674  
 This phenomenon is helpful for the defender in the scen- 675  
 ario where the generation loss can be compensated for by 676  
 somewhere else for the demand, making the transmission 677  
 line outage a priority. 678

### 679 C. Optimal defense strategy for IEEE 30-bus system 680 against attack-induced generation loss

We demonstrate the power of our deep  $Q$ -learning 681  
 algorithm to solve the generation loss problem for the 682  
 IEEE 30-bus system, which otherwise is not solvable using 683  
 conventional tabular  $Q$ -learning. Figure 8 shows  $G_{\text{loss}}$  per 684  
 episode for different protection sets, where the simula- 685  
 tion setting is that the defender is able to defend up to 686  
 three lines at a time, while the attacker can attack up to 687  
 five lines sequentially. Starting from a random protection 688  
 set {1, 2, 3}, with the worst-case scenario generation loss 689  
 per episode of 74.87 MW, the protection set evolves from 690  
 {16, 11, 14} to {16, 11, 15} and finally to the optimal pro- 691  
 tection set {16, 15, 28} that results in 50.49 MW generation 692



F7:1 FIG. 7. Evolution of reward function values during the learn-  
 F7:2 ing phase in the switching line problem in the IEEE 30-bus  
 F7:3 system for a random and an optimal protection set. While the  
 F7:4 defender chooses a random protection set {1, 2, 3}, the attacker  
 F7:5 finds an optimal sequence to obtain the reward of  $r = 10.4$  [cal-  
 F7:6 culated by Eq. (10) with  $r_1 = 10$  and  $r_2 = 1$ ]. After a number of  
 F7:7 cycles, the defender chooses {15, 16} as its protection set. As a  
 F7:8 result, the attacker fails to find a sequence with a reward of more  
 F7:9 than  $r = 2.6$ .



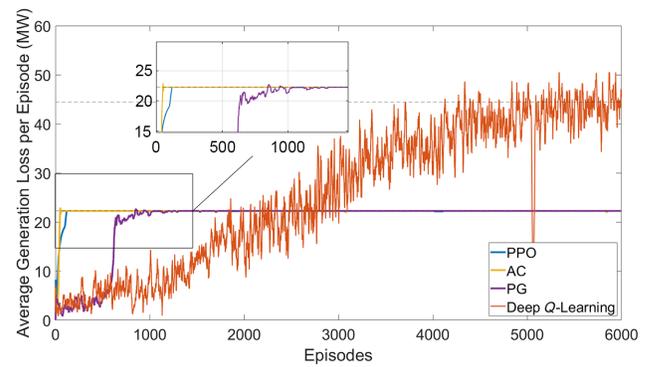
F8:1 FIG. 8. Optimal protection set against the worst-case scenario  
 F8:2 of generation loss in the IEEE 30-bus system. Defender chooses  
 F8:3 a random protection set  $\{1, 2, 3\}$ , whereas the attacker finds an  
 F8:4 optimal policy to maximize the generation loss. After a number  
 F8:5 of cycles, the defender chooses  $\{16, 15, 28\}$  as its protection set  
 F8:6 and, as a result, the worst-case scenario generation loss caused  
 F8:7 by the optimal attack sequence is reduced by 48.28%.

693 loss. Using the optimal protection set can result in 48.28%  
 694 mitigation of the worst-case generation loss, even if the  
 695 attacker chooses the optimal attacking sequence.

696 It is worth noting that the IEEE 30-bus system simulation  
 697 is used to demonstrate that conventional  $Q$ -learning  
 698 is unable to deal with this system, while our deep  $Q$ -  
 699 learning can. The system is only regarded as “large” in  
 700 a relative sense: it is much larger than the W&W 6-bus  
 701 benchmark system. Much larger systems are available,  
 702 e.g., the IEEE 300-bus or IEEE 3000-bus systems, which  
 703 can be simulated using specific power-grid software, such  
 704 as Gridlab-D. Deep RL methods are applicable to these  
 705 larger systems, but the required computations are beyond  
 706 our current capability.

#### 707 D. Comparison with alternative RL algorithms

708 We compare the performance of our deep  $Q$ -learning  
 709 algorithm with three widely used RL algorithms for dis-  
 710 crete state-action space systems: PG, AC, and PPO. The  
 711 PG algorithm [41] is a rudimentary policy-based model-  
 712 free online on-policy method, while the AC algorithm aims  
 713 to optimize the policy (actor) directly and train a critic  
 714 to estimate the return or future rewards [42]. PPO [43]  
 715 is an actor-critic model-free online on-policy algorithm  
 716 that alternates between data sampling by interacting with  
 717 the environment and optimization of a clipped objective  
 718 function, which leads to improved training stability by lim-  
 719 iting the size of the policy change at each step. We set  
 720 the learning rate, discount factor, and other applicable key  
 721 simulation parameters to the same values as in deep  $Q$ -  
 722 learning. The actor and critic networks in both the PPO  
 723 and AC algorithms have the same structure as the critic  
 724 network in our deep  $Q$ -learning algorithm and the actor



F9:1 FIG. 9. Comparison with representative existing RL algo-  
 F9:2 rithms. Shown is the performance comparison of the deep  $Q$ -  
 F9:3 learning with PG, AC, and PPO algorithms for the generation  
 F9:4 loss problem in the IEEE 30-bus system. Maximum genera-  
 F9:5 tion loss caused by the optimal attack sequences derived by  
 F9:6 the PPO, AC, and PG agents is 22.24 MW, while our deep  
 F9:7  $Q$ -learning agent is able to obtain 50.49 MW. While the deep  
 F9:8  $Q$ -learning algorithm takes a longer time to converge, reliability  
 F9:9 and efficiency are guaranteed.

725 network in the PG algorithm for fair comparison. The  
 726 protection set for all algorithms is set to  $\{16, 15, 28\}$ , which  
 727 is the Nash equilibrium in Sec. IV C. Figure 9 shows that  
 728 the maximum generation loss caused by the attacker in  
 729 the PPO, AC, and PG algorithms converges to 22.24 MW,  
 730 while that in our deep  $Q$ -learning algorithm converges to  
 731 50.49 MW. Generally, the deep  $Q$ -learning algorithm takes  
 732 a long time to converge, but the reliability and efficiency  
 733 compensate for the slow convergence since real-time com-  
 734 putation is not needed in strategy planning. Moreover, due  
 735 to the large size of action and state spaces, asymmetric and  
 736 stochastic state transitions, and insufficient exploration of  
 737 the state space intrinsic to the other algorithms, our deep  
 738  $Q$ -learning algorithm significantly outperforms the PPO,  
 739 AC, and PG algorithms.

#### 740 V. DISCUSSION

741 The problem of devising optimal defense strategies to  
 742 protect smart power grids from cyberattacks is of signifi-  
 743 cant current interest. Given a grid system, a general prin-  
 744 ciple is to simulate attacks to identify the scenario(s) that  
 745 can result in the most severe damage to define the best pos-  
 746 sible defense strategies. This attacker-defender interaction  
 747 problem can be modeled as a stochastic zero-sum game,  
 748 for which machine learning can provide effective solutions.  
 749 In recent years, conventional RL, in particular,  $Q$ -learning,  
 750 has been applied to the attacker-defender game problem,  
 751 but a fundamental shortcoming is the exponentially grow-  
 752 ing state space as the size of the system increases linearly.  
 753 We articulate a general deep  $Q$ -learning framework to  
 754 solve the game problem in arbitrarily large power-grid sys-  
 755 tems. We demonstrate that our deep  $Q$ -learning algorithm

756 typically leads to a Nash equilibrium, and the correspond- 805  
 757 ing strategy represents the optimal solution. We test the 806  
 758 proposed framework under different attack-defense scen- 807  
 759 arios for the W&W 6-bus system used in the current 808  
 760  $Q$ -learning literature and the relatively large IEEE 30-bus 809  
 761 system that cannot be solved with the conventional  $Q$ - 810  
 762 learning algorithm. We also compare the results of our 811  
 763 deep  $Q$ -learning algorithm to those from three alterna- 812  
 764 tive but state-of-the-art RL algorithms and demonstrate the 813  
 765 superiority of our method.

766 Immediate future work is expanding the deployment of 814  
 767 the deep RL algorithms to a general sum problem, in which 815  
 768 both the attacker and defender have limited resources that 816  
 769 they can use for their actions. The reward function would 817  
 770 also be different from the one used in this paper, where the 818  
 771 defender attempts to mitigate the consequences, whereas 819  
 772 the attacker has a set objective. The results in this paper 820  
 773 suggest that deep  $Q$ -learning can be effective at address- 821  
 774 ing the general sum game to devise the optimal resource 822  
 775 allocation policy. 823

## 776 ACKNOWLEDGMENTS

777 This work is mainly supported by the US-Israel Energy 829  
 778 Center managed by the Israel-US Binational Industrial 830  
 779 Research and Development (BIRD) Foundation. This work 831  
 780 is also supported by AFOSR under Grant No. FA9550-21- 832  
 781 1-0438. 833

## 782 APPENDIX: A DETAILED DESCRIPTION OF THE 783 DEEP $Q$ -LEARNING METHOD

784 Deep  $Q$ -learning is a model-free framework in which 838  
 785 the agent uses a neural network architecture to train a 839  
 786 critic to estimate the future cumulative rewards charac- 840  
 787 terizing how valuable one action is at each state. While 841  
 788 there are reinforcement learning methods for continuous 842  
 789 action spaces (e.g., deep deterministic policy gradient 843  
 790 [44] and twin-delayed deep deterministic policy gradient 844  
 791 [45]), deep  $Q$ -learning is only applicable to discrete action 845  
 792 spaces. 846

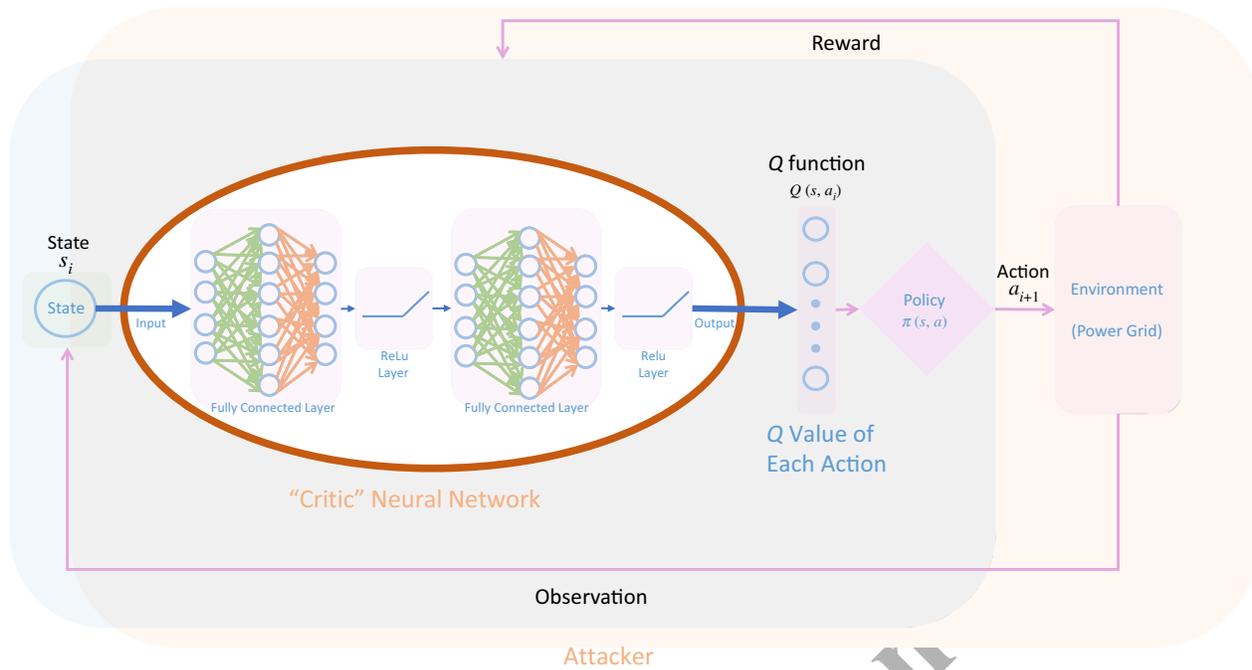
793 The structure of the deep  $Q$ -learning method in our work 847  
 794 is shown Fig. 10, which illustrates what happens inside the 848  
 795 attacker block in Fig. 5. Modeling the attacker-defender 849  
 796 interaction as a zero-sum game has the advantage of learn- 850  
 797 ing a single  $Q$  function (in a general sum game, learning 851  
 798 multiple  $Q$  functions would be necessary). For each state 852  
 799 input, the deep  $Q$ -learning structure returns an approxima- 853  
 800 tion of the  $Q$  function for that state and all possible actions. 854  
 801 In our problem, by “state” we mean the state of the trans- 855  
 802 mission lines in the power grid, which is denoted as a 856  
 803 binary-valued vector. The attacker’s action is chosen from 857  
 804 the set  $A = \{1, 2, 3, \dots\}$ , where action  $i$  means attacking 858  
 transmission line  $i$ . The defender’s action is a set consisting 859

of  $n$  transmission lines denoted as the protection set. The 805  
 environment block in Fig. 10 represents the power grids 806  
 studied in this paper. As described in the main text, we 807  
 employ DCSIMSEP, a dc load flow simulator of cascading 808  
 (separation) in power systems, to simulate the dynamics 809  
 of the power grid. Using our modified DCSIMSEP code, we 810  
 generate the observation and rewards for each attack (and 811  
 defense) actions and feed them to the algorithm in the next 812  
 step. 813

A deep  $Q$ -learning agent is represented by a critic 814  
 neural network. During the training phase, this critic is 815  
 trained to approximate the expectation of the cumulative 816  
 future rewards. The critic neural network is parameterized. 817  
 During training, the agent tunes the parameter values to 818  
 improve the accuracy of the estimation. The neural net- 819  
 work structure consists of two fully connected and two 820  
 ReLu layers (as detailed in Table I). In particular, a fully 821  
 connected layer multiplies the input by a weight vector 822  
 and adds a bias into it, which is similar to a nonlinear 823  
 principal component analysis for improving the estima- 824  
 tion accuracy. The ReLu layers set the negative values of 825  
 the input to zero and perform a threshold operation on the 826  
 input; these are nonlinear transformations to expedite the 827  
 training process. 828

Here, we model the attacker and defender interaction as 829  
 a zero-sum game, with the goal of disabling a fixed num- 830  
 ber of transmission lines or maximizing (minimizing) the 831  
 generation loss. Both players receive the reward for (or 832  
 cost of) their actions. The attacker uses deep  $Q$ -learning 833  
 to optimize the attack sequence. During the training process, 834  
 the agent explores the state space, i.e., the attacker attacks 835  
 different transmission lines to observe the results. This 836  
 exploration follows a standard greedy algorithm method, 837  
 where sometimes the attacker launches random attacks and 838  
 at other times the attack is based on what the attacker 839  
 has learned so far. The past experiences are stored using 840  
 an experience buffer. The critic neural network is updated 841  
 based on a pool of experiences randomly sampled from this 842  
 buffer. Once an optimal attacking strategy is reached, it is 843  
 transmitted to the defender, and the defender will update its 844  
 protection set to be better prepared against future attacks. 845  
 This process continues until the Nash equilibrium of the 846  
 game is reached. 847

We perform the simulation using MATLAB’s reinforce- 848  
 ment learning toolbox. For the deep  $Q$ -learning algorithm, 849  
 we use the rIDQNAgent object. The options set for rIDQ- 850  
 NAgentOptions are listed in Table I. The state space 851  
 is defined using rlNumericSpec, and the action space 852  
 type is selected as rlFiniteSetSpec. No external lower 853  
 or upper limits are applied to these spaces. The envi- 854  
 ronment (*env* object) is customized using the modified 855  
 DCSIMSEP. Eventually, the critic is a rlQValueRepresent- 856  
 ation object with the neural network layer depicted in 857  
 Fig. 10. The codes and simulation results are available at 858  
 Github [46]. 859



F10:1 FIG. 10. Structure of deep  $Q$ -learning algorithm used in this paper. Structure describes the processes inside the attacker block in Fig.  
 F10:2 5. Environment block contains the power grids simulated using our modified DCSIMSEP algorithm. DCSIMSEP generates the observation  
 F10:3 and rewards for each attack (and defense), which are fed to the algorithm in the next step. Through interacting with the environment,  
 F10:4 the critic returns an approximation of the  $Q$  function for the input state (the state of transmission lines) and all possible actions (attack  
 F10:5 actions or protection sets). This critic neural network is parameterized. During training, the agent tunes the parameter values to make  
 F10:6 the estimation more accurate. Critic consists of two fully connected and two ReLU layers, the specifications of which are listed in Table  
 F10:7 I. Attacker uses this algorithm to optimize the attack sequence. Once an optimal attacking strategy is reached, the defender will update  
 F10:8 its protection set (Fig. 5) to be better prepared against future attacks. This repeats until the optimal protection set has been found.

860 [1] P. Pourbeik, P. S. Kundur, and C. W. Taylor, The anatomy of  
 861 a power grid blackout—root causes and dynamics of recent  
 862 major blackouts, *IEEE Power Energy Mag.* **4**, 22 (2006).  
 863 [2] J. Xie, A. Stefanov, and C.-C. Liu, Physical and cyber  
 864 security in a smart grid environment, *Wiley Interdis. Rev.*  
 865 *Energy Envir.* **5**, 519 (2016).  
 866 [3] R. Langner, Stuxnet: Dissecting a cyberwarfare weapon,  
 867 *IEEE Secur. Priv.* **9**, 49 (2011).  
 868 [4] G. Liang, S. R. Weller, J. Zhao, F. Luo, and Z. Y. Dong, The  
 869 2015 Ukraine blackout: Implications for false data injection  
 870 attacks, *IEEE Trans. Power Sys.* **32**, 3317 (2017).  
 871 [5] Y. Liu, P. Ning, and M. K. Reiter, False data injection  
 872 attacks against state estimation in electric power grids,  
 873 *ACM Trans. Info. Sys. Secu.* **14**, 1 (2011).  
 874 [6] L. Xie, Y. Mo, and B. Sinopoli, Integrity data attacks in  
 875 power market operations, *IEEE Trans. Smart Grid* **2**, 659  
 876 (2011).  
 877 [7] M. Mohammadpourfard, Y. Weng, and M. Tajdinian,  
 878 Benchmark of machine learning algorithms on capturing  
 879 future distribution network anomalies, *IET Gene. Transmi.*  
 880 *Distri.* **13**, 1441 (2019).  
 881 [8] A. Shefaei, M. Mohammadpourfard, B. Mohammadi-  
 882 ivatloo, and Y. Weng, Revealing a new vulnerability of  
 883 distributed state estimation: A data integrity attack and an

unsupervised detection algorithm, *IEEE Trans. Cont. Net.* 884  
*Sys.* (2021),. 885  
 [9] N. Enriquez and Y. Weng, in *Asian Conference on Machine* 886  
*Learning (ACML), PMLR 157* (2021) p. 1333. 887  
 [10] J. F. Nash, Equilibrium points in  $n$ -person games, *Proc. Nat.* 888  
*Aca. Sci. (USA)* **36**, 48 (1950). 889  
 [11] T. Başar and G. J. Olsder, *Dynamic Non-Cooperative Game* 890  
*Theory* (SIAM, 1998), 2nd ed. 891 Q9  
 [12] W. Saad, Z. Han, H. V. Poor, and T. Basar, Game-theoretic 892  
 methods for the smart grid: An overview of microgrid 893  
 systems, demand-side management, and smart grid com- 894  
 munications, *IEEE Sig. Proc. Mag.* **29**, 86 (2012). 895  
 [13] S. Poudel, Z. Ni, X. Zhong, and H. He, in *2016 Interna-* 896  
*tional Joint Conference on Neural Networks (IJCNN)*, p. 897  
 2730. 898  
 [14] N. I. Haque, M. H. Shahriar, M. G. Dastgir, A. Deb- 899  
 nath, I. Parvez, A. Sarwat, and M. A. Rahman, Machine 900  
 learning in generation, detection, and mitigation of cyber- 901  
 attacks in smart grid: A survey, arXiv preprint (2020), 902  
 ArXiv:2010.00661. 903  
 [15] Y. Chen, S. Huang, F. Liu, Z. Wang, and X. Sun, Evaluation 904  
 of reinforcement learning-based false data injection attack 905  
 to automatic voltage control, *IEEE Trans. Smart Grid* **10**, 906  
 2158 (2019). 907  
 [16] B. Ning and L. Xiao, in *2021 40th Chinese Control Confer-* 908  
*ence (CCC)* (IEEE), p. 8598. 909

- 910 [17] C. J. C. H. Watkins and P. Dayan, *Q*-learning, *Mach. Learn.* **8**, 279 (1992). 959
- 911 960
- 912 [18] J. Yan, H. He, X. Zhong, and Y. Tang, *Q*-Learning-based 961
- 913 vulnerability analysis of smart grid against sequential topol- 962
- 914 ogy attacks, *IEEE Trans. Info. Foren. Secu.* **12**, 200 (2017). 963
- 915 [19] Z. Wang, H. He, Z. Wan, and Y. Sun, Coordinated topol- 964
- 916 ogy attacks in smart grid using deep reinforcement learning, 965
- 917 *IEEE Trans. Indust. Info.* **17**, 1407 (2020). 966
- 918 [20] C. Roberts, S.-T. Ngo, A. Milesi, S. Peisert, D. Arnold, 967
- 919 S. Saha, A. Scaglione, N. Johnson, A. Kocheturov, and 968
- 920 D. Fradkin, in *2020 IEEE International Conference on* 969
- 921 *Communications, Control, and Computing Technologies* 970
- 922 *for Smart Grids (SmartGridComm)* (IEEE), p. 1. 971
- 923 [21] Y. Li and J. Wu, Low latency cyberattack detection in smart 972
- 924 grids with deep reinforcement learning, (2022),. 973
- 925 [22] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. 974
- 926 Antonoglou, D. Wierstra, and M. A. Riedmiller, Play- 975
- 927 ing Atari with deep reinforcement learning, (2013), 976
- 928 [ArXiv:1312.5602](https://arxiv.org/abs/1312.5602). 977
- 929 [23] J. Hu and M. P. Wellman, in *ICML '98 Proceedings of the* 978
- 930 *Fifteenth International Conference on Machine Learning* 979
- 931 (1998), p. 242. 980
- 932 [24] S. R. Etesami and T. Basar, Dynamic games in cyber- 981
- 933 physical security: An overview, *Dyn. Games Appl.* **9**, 884 982
- 934 (2019). 983
- 935 [25] D. Vrabie and F. Lewis, in *The 2010 International Joint* 984
- 936 *Conference on Neural Networks (IJCNN)*, p. 1. 985
- 937 [26] Q. Zhu, H. Tembine, and T. Bbar, Heterogeneous learning 986
- 938 in zero-sum stochastic games with incomplete information, 987
- 939 49th IEEE Conference on Decision and Control (CDC), 219 988
- 940 (2010). 989
- 941 [27] P. Bommannavar, T. Alpcan, and N. Bambos, Security 990
- 942 risk management via dynamic games with learning, 2011 991
- 943 IEEE International Conference on Communications (ICC), 992
- 944 1 (2011). 993
- 945 [28] A. Truong, S. R. Etesami, J. Etesami, and N. Kiyavash, 994
- 946 Optimal attack strategies against predictors - learning from 995
- 947 expert advice, *IEEE Trans. Info. Foren. Secu.* **13**, 6 (2018). 996
- 948 [29] Q. Zhu, H. Tembine, and T. Başar, in *Reinforcement Learn-* 997
- 949 *ing and Approximate Dynamic Programming for Feedback* 998
- 950 *Control* (2013), p. 303. 999
- 951 [30] K.-W. Chung, C. A. Kamhoua, K. A. Kwiat, Z. T. Kalbar- 1000
- 952 czyk, and R. K. Iyer, in *2016 IEEE 17th International Sym-* 1001
- 953 *posium on High Assurance Systems Engineering (HASE)* 1002
- 954 (2016), p. 1. 1003
- 955 [31] X. He, H. Dai, and P. Ning, Improving learning and adapta- 1004
- 956 tion in security games by exploiting information asymme- 1005
- 957 try, 2015 IEEE Conference on Computer Communications 1006
- 958 (INFOCOM), 17872015. 1007
- [32] K. K. Trejo, J. B. Clempner, and A. S. Poznyak, in *2016* 959
- IEEE 55th Conference on Decision and Control (CDC)*, 960
- p. 5484. 961
- [33] M. J. Eppstein and P. D. H. Hines, A “random chemistry” 962
- algorithm for identifying collections of multiple contingen- 963
- cies that initiate cascading failure, *IEEE Trans. Power Sys.* 964
- 27**, 1698 (2012). 965
- [34] P. Rezaei, P. D. H. Hines, and M. J. Eppstein, Estimating 966
- cascading failure risk with random chemistry, *IEEE Trans.* 967
- Power Sys.* **30**, 2726 (2015). 968
- [35] G. Liang, S. R. Weller, J. Zhao, F. Luo, and Z. Y. Dong, 969
- A framework for cyber-topology attacks: Line-switching 970
- and new attack scenarios, *IEEE Trans. Smart Grid* **10**, 1704 971
- (2019). 972
- [36] S. Paul and Z. Ni, in *2018 International Joint Conference* 973
- on Neural Networks (IJCNN)*, p. 1. 974
- [37] Z. Ni and S. Paul, A multistage game in smart grid secu- 975
- rity: A reinforcement learning solution, *IEEE Trans. Neural* 976
- Netw. Learn. Syst.* **30**, 2684 (2019). 977
- [38] D. S. H. van Hasselt and A. Guez, in *Proceedings of the* 978
- thirtieth AAAI Conference on Artificial Intelligence* (2016), 979
- p. 1928. 980
- [39] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. 981
- Veness, M. G. Bellemare, A. Graves, M. Riedmiller, 982
- A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, 983
- A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. 984
- Wierstra, S. Legg, and D. Hassabis, Human-level control 985
- through deep reinforcement learning, *Nature* **518**, 529 986
- (2015). 987
- [40] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An* 988
- Introduction* (The MIT Press, 2018), 2nd ed. 989
- [41] R. J. Williams, Simple statistical gradient-following algo- 990
- rithms for connectionist reinforcement learning, *Mach.* 991
- Learn.* **8**, 229 (1992). 992
- [42] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, 993
- T. Harley, D. Silver, and K. Kavukcuoglu, in *Proceedings* 994
- of the 33rd International Conference on Machine Learning* 995
- Research 2016*, Vol. 48 (2016), p. 1928. 996
- [43] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. 997
- Klimov, Proximal policy optimization algorithms, (2017), 998
- [ArXiv:1707.06347](https://arxiv.org/abs/1707.06347). 999
- [44] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, 1000
- Y. Tassa, D. Silver, and D. Wierstra, Continuous control 1001
- with deep reinforcement learning, (2015), arXiv preprint 1002
- [ArXiv:1509.02971](https://arxiv.org/abs/1509.02971). 1003
- [45] S. Fujimoto, H. van Hoof, and D. Meger, Addressing func- 1004
- tion approximation error in actor-critic methods, (2018), 1005
- arXiv preprint [ArXiv:1802.09477](https://arxiv.org/abs/1802.09477). 1006
- [46] [https://github.com/AminMoradiXL/DQN\\_grid](https://github.com/AminMoradiXL/DQN_grid). 1007