

Distributed Energy Management and Demand Response in Smart Grids: A Multi-Agent Deep Reinforcement Learning Framework

AMIN SHOJAEIGHADIKOLAEI, University of Kansas

ARMAN GHASEMI, University of Kansas

KAILANI JONES, University of Kansas

YOUSIF DAFALLA, University of Kansas

ALEXANDRU G. BARDAS, University of Kansas

REZA AHMADI, Amazon Web Services, Kirkland

MORTEZA HASHEMI, University of Kansas

This paper presents a multi-agent Deep Reinforcement Learning (DRL) framework for autonomous control and integration of renewable energy resources into smart power grid systems. In particular, the proposed framework jointly considers demand response (DR) and distributed energy management (DEM) for residential end-users. DR has a widely recognized potential for improving power grid stability and reliability, while at the same time reducing end-users' energy bills. However, the conventional DR techniques come with several shortcomings, such as the inability to handle operational uncertainties while incurring end-user disutility, which prevents widespread adoption in real-world applications. The proposed framework addresses these shortcomings by implementing DR and DEM based on real-time pricing strategy that is achieved using deep reinforcement learning. Furthermore, this framework enables the power grid service provider to leverage distributed energy resources (i.e., PV rooftop panels and battery storage) as dispatchable assets to support the smart grid during peak hours, thus achieving management of distributed energy resources. Simulation results based on the Deep Q-Network (DQN) demonstrate significant improvements of the 24-hour accumulative profit for both prosumers and the power grid service provider, as well as major reductions in the utilization of the power grid reserve generators.

CCS Concepts: • Theory of computation → Multi-agent reinforcement learning; • Information systems → Mobile information processing systems;

Additional Key Words and Phrases: Demand Response, Distributed Energy Management, Reinforcement Learning, Deep Q-Network (DQN)

1 INTRODUCTION

Recent proliferation in renewable energy sources (RESs) and distributed energy resources (DERs) along with the advances in information and communication technologies (ICT) has facilitated a paradigm shift from the traditional power consumers to the more resourceful energy *prosumers*. A prosumer is an active end-user with the ability to consume and produce energy [1]. The number of prosumers is increasing as more households are proceeding to install solar photovoltaic (PV) rooftop panels and energy storage system (ESS) [2]. The premise of such a distributed network of energy resources is to provide economic and ancillary benefits to both prosumers and power grid service provider (SP) [3]. The interaction between the prosumers and SP has been investigated in several prior studies. Many works utilize model-based control paradigms, such as model predictive control (MPC) [4], mixed-integer linear programming (MILP) [5], dynamic and stochastic programming [6], and alternating direction method of multiplier (ADMM) [7]. One key step for applying these techniques is to establish an accurate model that captures the dynamics of the microgrid and the interactions between various components with all the operational constraints.

The massive deployment of DERs and RESs dramatically increases the complexity and uncertainty of the system due to the more complex cross-area power balancing between SP and prosumers. This makes the accurate system model hard to obtain. The model-based nature of the aforementioned methods also exhibits limited generalization capabilities [8]. Difficulties in handling nonlinear behavior, volatile renewable generations, and heterogeneity of the end-users represent other major shortcomings of these methods. To overcome these challenges, model-free Reinforcement Learning (RL) techniques are proven beneficial for demand-side energy management since they do not require an explicit model of the environment. In general, RL frameworks are emerging as the pre-eminent tool for sequential decision-making problems within unknown environments. Similar to many other scientific and engineering domains, the RL-based solutions are receiving more attention from the power system society. For instance, voltage and frequency control [9], market bidding [10], microgrid energy management [11], and demand response (DR) [12] are just a few examples of power-related problems that can be solved using RL.

In this paper, and to bridge the gap between the power system and RL communities, we investigate the mutual interplay between participants and the different tasks in a residential microgrid. We propose a framework based on Deep Reinforcement Learning (DRL) for both SP and prosumers to enable dynamic decision-making. A service provider agent (SPA) is deployed for solving the distributed energy management (DEM) problem in the microgrid by dynamically determining the electricity buy price as a control parameter for optimization of the energy management across distributed prosumers. For instance, a higher electricity buy price during peak hours (e.g., evening hours) incentivizes the prosumers with surplus energy to discharge their battery. As a result, the prosumer receives economic benefit in terms of electricity bill reduction. This process is called Demand Response (DR). The purpose of DR is to encourage the prosumers to actively participate in a program and contribute to the optimal energy distribution in the electricity retail market. Moreover, buying electricity from distributed prosumers would enable the power grid service provider to support higher demands during peak hours. In our envisioned system model, the prosumer agent (PA) solves the prosumer energy management problem by determining the charge/discharge of the battery installation. Thus, each prosumer independently maximizes its own profit.

In summary, our main contributions are as follows:

- We formally define the interaction between SP and the prosumers as a Markov Decision Process (MDP), and develop a multi-agent DRL framework that interweaves the real-time energy management over the microgrid with the prosumer side real-time DR, using demand-dependent dynamic pricing for incentivizing prosumers' participation in DR.
- We show that the proposed DRL framework enables the service provider to leverage energy storage as dispatchable assets, while incentivizing the prosumers to actively participate in the program to support the grid activities.
- Our numerical results based on Deep Q-Network (DQN) demonstrate that the proposed framework reduces the average daily bill for prosumers, while at the same time it provides higher profits for the grid service provider (SP) by leveraging distributed energy resources instead of tapping into traditional reserve generation facilities with higher costs.

The remainder of this paper is organized as follows. Section 2 covers related work, followed by the system model and problem formulation in Section 3. Next, Section 4 formulates the DRL framework and Section 5 presents numerical results. Finally, Section 6 concludes the paper.

2 RELATED WORK

The interaction between service provider and end-users have been investigated in many prior works. Recently, there has been growing interests in adopting RL to address the concerns of this interaction. These attempts can be categorized as follows:

Prior Works Focused on Service Provider: A multitude of prior works use RL in the context of power systems, mainly to address service provider concerns. For example, a hierarchical electricity market with bidding and pricing over wholesale and retailer using DRL was proposed in [13]. In [14], a Deep Deterministic Policy Gradient (DDPG) algorithm was used to solve the bidding problem of several generation companies. Moreover, RL has been widely employed for energy management optimization [15, 16]. In [15], the authors proposed a multi-agent distributed energy management using Q-learning framework with time-of-use pricing for a microgrid energy market with a centralized battery pack and renewable generation. The authors in [16] designed a novel RL method that uses classical recurrent neural networks instead of SAR (State-Action-Reward) method to solve the microgrid energy management in the Industrial Internet of Things (IIoT). In contrast to these works, we propose a dynamic pricing scheme with DQN at the service provider side, combined with distributed battery and PV installations at the demand side.

Prior Works Focused on End-users: The authors in [17] developed a DDQN-based load scheduling algorithm with a time-of-use scheme to reduce the peak load demand and the operation cost for the distribution system, however without considering prosumers. In [18], the interaction between households and the service provider was modeled to jointly solve the energy scheduling consumption and preserving privacy for the households by using policy gradient RL method. Furthermore, [19] developed a home energy management system (HEMS) using Actor-Critic technique based on adaptive dynamic programming. Likewise, in [20], an energy management DR using Q-learning adjusts the consumption plan for thermostatically controlled loads. Along the same lines, [21] proposed a framework for home energy management based on multi-agent Q-learning to minimize the electricity bill as well as DR-induced dissatisfaction costs for end-users. Their method considered proliferation of rooftop PV systems, but did not consider home energy storage systems. Moreover, the works in [22, 23], focused on reducing the high peak load in large-scale HEM using Entropy-based multi-agent deep reinforcement learning, and the work in [24] applied prioritized DDPG in multi-carrier energy system to solve the real-time energy management problem. It is also worth noting that these works [17–23] do not consider the grid side and only propose optimization algorithms for end-users.

Prior Works Focused on Both Service Provider and End-users: Other research works have addressed the interactions between retailer/SP and end-users. For example, [25] investigated a gradient-based method to minimize the aggregate load demand by assuming the presence of consumption scheduling devices on the users side. Furthermore, works such as [26–28] leveraged the Stackelberg game to model and maximize the retailer profit and minimize the payment bill of its customers. The authors included renewable energy generation on the retailer side, which is different from our model with distributed PV installations across prosumers. The studies in [29] and [30] proposed a hierarchical agent-based framework to maximize both retailer and customers profits. A Q-learning DR is proposed in [31] for the energy consumption scheduling problem that can work without prior information and leads to reduced system costs. Also, in [32] both service provider profit and consumers' cost are considered to find the optimal retail price. However, unlike our effort, the works presented in [29–32] only consider regular electricity consumers, rather than prosumers with generation and storage capabilities. In contrast to these efforts, we model a microgrid that consists of both consumers and prosumers that are equipped with battery storage. As a result, the proposed multi-agent DRL framework achieves DR with dynamic pricing and distributed energy

Table 1. Taxonomy of Related Work on Management of Distributed Energy Resources. (✓: Considered, –: Not Considered)

Reference	Energy management DR	Retailer/end-user profit optimization	Renewable penetration consideration	Energy storage consideration	RL algorithm
[29],[31],[32]	✓	✓	–	–	✓
[27],[28],[18],[30]	✓	✓	–	–	–
[21]	–	✓	✓	–	✓
[22], [23]	–	✓	–	–	✓
[26]	✓	–	–	✓	✓
Our proposed framework	✓	✓	✓	✓	✓

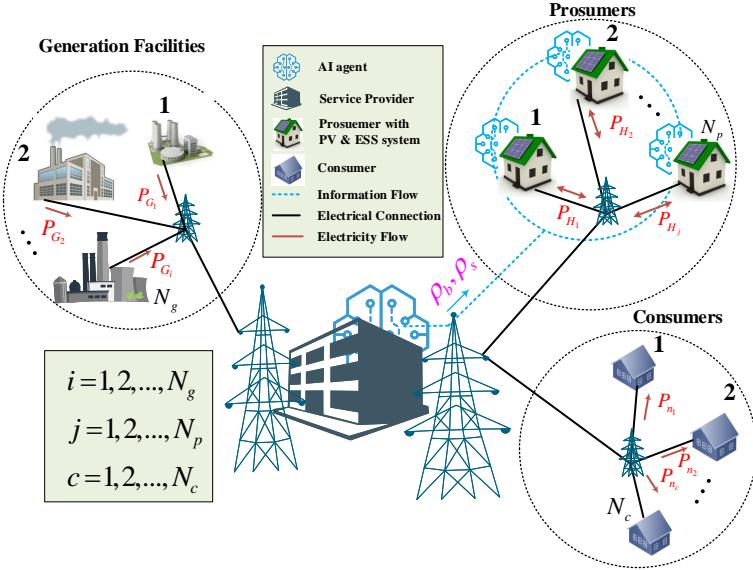


Fig. 1. Proposed microgrid system architecture that consists of generation facilities, traditional consumers, and prosumers that are equipped with solar rooftop panels and energy storage and deep reinforcement learning agents. Grid and prosumers are equipped with deep reinforcement learning agents to dynamically adjust their policies in terms of buy/sell prices and power injection.

management for a microgrid with large-scale rooftop PV panels and energy storage system. Table 1 compares our work with the previous related works.

3 MICROGRID MODEL AND PROBLEM FORMULATION

The envisioned microgrid model is shown in Figure 1. On the end-user side, the microgrid consists of regular consumers as well as prosumers that are equipped with PV panels, battery, and an intelligent agent. The presence of consumers ensure that the microgrid always has energy deficiency. The conventional generation units are connected to the microgrid to meet the deficiency. The base generation units have low cost for the service provider but if their capacity is not enough, the reserve and the more expensive generation facilities will be dispatched to meet the real-time microgrid demand [33]. Without loss of generality, this paper forgoes the ramp rate constraints of the generation facilities. The agent decides on whether to consume the excess energy generated by the PV system, store it in the battery, or sell it to the grid. On the other side, the SP is responsible

for dispatching power to the loads from various generation facilities, including distributed energy resources of prosumers.

In this paper, distributed energy management and demand response optimization problems are defined for the SP and prosumers, respectively. On the grid side, the SP is engaged in energy management with the objective of improving its own profit. On the household side, prosumers engage in DR with the goal of reducing their daily electricity bill. The complex interactions between the two sides denote the joint DEM-DR approach proposed in this paper.

Service Provider Energy Distribution Model: In the presence of prosumers, the DEM problem of the microgrid evolves into a much more complex problem since the prosumers also appear as generation facilities in certain time periods due to injection of power into the grid. Hence, in this paper, the DEM formulation is defined as follows ,

$$\max_{\rho_b} \Psi_G(T) = R_G(T) - \left\{ \sum_{i=1}^{N_g} F_{G_i}(T) + \sum_{j=1}^{N_p} F_{H_j}(T) \right\} \quad (1)$$

Subject to:

$$\sum_{i=1}^{N_g} F_{G_i}(T) = \sum_{i=1}^{N_b} F_{B_i}(T) + \sum_{i=1}^{N_r} F_{R_i}(T), \quad (2)$$

$$R_G(T) = \int_0^T P_D(t) \cdot \rho_s(t) dt, \quad (3)$$

$$F_{H_j}(T) = \int_0^T P_{H_j}(t) \cdot \rho_b(t) dt \quad \text{for } P_{H_j}(t) \geq 0, \quad (4)$$

$$P_D(t) = \sum_{i=1}^{N_g} (P_{G_i}(t) - P_{loss_i}(t)) + \sum_{j=1}^{N_p} P_{H_j}(t), \quad (5)$$

$$P_{loss_i} = \beta_i \times P_{G_i}^2(t), \quad (6)$$

$$P_{G_i}^{\min} \leq P_{G_i}(t) \leq P_{G_i}^{\max} \quad \text{for } i = 1, 2, \dots, N_g, \quad (7)$$

where $\Psi_G(T)$ denotes the SP profit over a time horizon of T , and $R_G(T)$ is SP revenue as a result of selling electricity to the loads, which is calculated by Equation (3) where $P_D(t)$ and $\rho_s(t)$ denote the total system demand and electricity sell price at any given time, respectively. $F_{G_i}(T)$ is the cost of generation unit i defined as the quadratic function of P_{G_i} , which means $F_{G_i}(P_{G_i}) = a_i P_{G_i}^2 + b_i P_{G_i} + c_i$ where a_i , b_i , and c_i are fitting parameters and P_{G_i} is the amount of power purchased from generation unit i [34]. This cost consists of the cost of base generation units $F_{B_i}(T)$ and cost of reserve generation units $F_{R_i}(T)$ in Equation (2) where N_b and N_r are the number of base and reserve generation facilities, respectively. N_g is the number of all generation facilities. If the power demand exceeds the amount of base generation, the overall cost will be significantly higher because procuring reserve power is more costly. $F_{H_j}(T)$ is the cost of buying electricity from the j^{th} prosumer and calculated by Equation (4) where $P_{H_j}(t)$ and $\rho_b(t)$ denote the power injected by the j^{th} prosumer and electricity buy price, respectively. Equation (5) illustrates the total generation and total demand power balance requirement, which needs to be maintained at any given time slot t . N_p and N_c denote the number of prosumers and consumers, respectively. P_{loss_i} denotes the losses induced by generation unit i , where β_i is the loss coefficient [35]. In addition, the power output of each generation unit must not exceed its operation limits, which are described in Equation (7).

The SP's goal is to maximize its profit. To do so, the SP dynamically changes the buy price to incentivize the prosumers to sell their excess energy to the grid, especially during the peak demand hours. Hence, the control variable in DEM problem is the buy price $\rho_b(t)$.

Prosumer Side Demand Response Model: The objective of the prosumers is to minimize their daily electricity bill by engaging in DR as a response to the dynamic buy price controlled by the SP. Therefore, the prosumer side profit maximization is formulated as follows,

$$\max_{P_{b_j}} \Psi_{H_j}(T) = F_{H_j}(T) + \mathbb{1}_{\{P_{H_j}(t)<0\}} \times C_{H_j}(T) \quad (8)$$

Subject to:

$$C_{H_j}(T) = \int_0^T P_{H_j}(t) \cdot \rho_s(t) dt \quad \text{for } P_{H_j}(t) < 0, \quad (9)$$

$$F_{H_j}(T) = \int_0^T P_{H_j}(t) \cdot \rho_b(t) dt \quad \text{for } P_{H_j}(t) \geq 0, \quad (10)$$

$$P_{H_j}(t) = P_{PV_j}(t) - P_{b_j}(t) - P_{C_j}(t), \quad (11)$$

$$SoC_{b_j}^{\min} \leq \frac{1}{C_{b_j}} \int_0^t P_{b_j}(\tau) d\tau + SoC_{b_j}(0) \leq SoC_{b_j}^{\max}, \quad (12)$$

$$0 \leq P_{PV_j}(t) \leq P_{PV_j}^{\max}, \quad (13)$$

$$|P_{H_j}(t)| \leq P_{H_j}^{\max}, \quad (14)$$

$$SoC_{b_j}(0), \quad SoC_{b_j}^{\min}, \quad P_{H_j}^{\max} \geq 0,$$

where $\Psi_{H_j}(T)$ is the profit of the j^{th} prosumer and $C_{H_j}(T)$ is the cost of buying electricity from the grid over the time horizon T , which is described by Equation (9). $F_{H_j}(T)$ describes the prosumer j revenue as a result of selling power to the grid and is calculated by Equation (4). Each prosumer has an internal household-scale power system with a power balance equation of its own that needs to be satisfied at all times. This is formulated by Equation (10) where P_{PV_j} , P_{b_j} and P_{C_j} are the amount of PV generation, battery charging/discharging power, and power consumption of the j^{th} prosumer, respectively. According to Equation (10), P_{H_j} can take a positive or negative value depending on the amount of PV generation and battery charge or discharge action. A positive (negative) value means the prosumer j injects (purchases) power to (from) the grid. Hence, $\mathbb{1}_{\{P_{H_j}(t)<0\}}$ is one/zero when the injected power from prosumer j is negative/positive. The state of charge of each prosumer battery should be maintained within a safe operational range as illustrated in Equation (11) where C_{b_j} is the nominal battery capacity and $SoC_{b_j}(0)$ represents the initial state of charge of the battery. $SoC_{b_j}^{\min}$ and $SoC_{b_j}^{\max}$ are the lowest and highest allowable state of charge for the battery, respectively. Furthermore, the charging and discharging power of the battery is limited as Equation (12), where $P_{b_j}^{\max, \text{charge}}$ and $P_{b_j}^{\max, \text{dis}}$ are the maximum charge/discharge power ratings of the battery. The constraint in Equation (13) is imposed by the physical limitations of the size of the PV system and shows that the amount of PV generation is limited by the $P_{PV_j}^{\max}$. Dynamic power injection into the grid might cause instability in some cases. Thus, the maximum allowable power injection is limited by Equation (14).

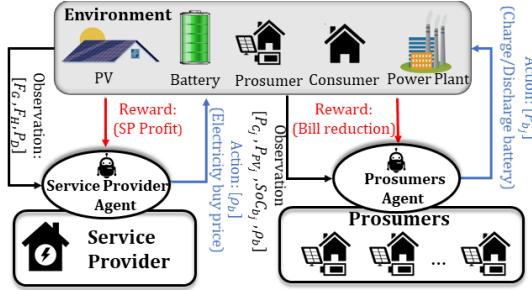


Fig. 2. Implemented Reinforcement Learning framework

In this paper, the formulated DR method is different from the conventional time-based DR, such as Real-Time Pricing (RTP) methods [25]. Conventionally, RTP approaches rely on dynamically changing the electricity sell price to motivate the customers to alter their energy use profile, while the proposed DEM-DR algorithm relies on dynamically changing the buy price instead of the sell price. This potentially incentivizes prosumers to participate in the program without any customer dissatisfaction. At any given time slot t , both aforementioned optimization problems must be solved at the same time. Dynamic SP pricing scheme in DEM-DR problem reinforces the need to implement the optimization problems in short intervals (e.g., 15 minutes).

4 REINFORCEMENT LEARNING FOR DEM AND DR

Reinforcement Learning (RL) is a model-free machine learning technique that trains an agent to take optimal actions through repeated interactions with an environment. The general RL framework consists of a set of states \mathcal{S} , a set of actions \mathcal{A} , a reward function R , and probability of transition between the states. The policy π is a mapping from the states of the environment (s^t) to the actions (a^t), i.e., $\pi : \mathcal{S} \rightarrow \mathcal{A}$. The agent's ultimate goal is to learn which action $a^t \in \mathcal{A}$ to take at each time instance t to maximize its cumulative reward over time. The action a^t results in reward r^t , and the environment transitions from the state s^t to $s^{t+1} \in \mathcal{S}$.

In this section, we use t as superscript to denote discrete time slots. It should be noted that there is no interaction between the prosumers, such that their goals are to maximize their local profit regardless of other prosumers' performance.

4.1 Decision-Making Problem Formulation

In order to tackle the optimization problems in Section 3, we leverage a multi-agent DRL approach such that we define agent for the grid side SP and agents for the prosumer side with PV and battery storage capabilities. The SPA is representative of the service provider agent, while PA_j denotes the j^{th} prosumer agent. Each agent observes the environment and subsequently takes an action, receiving a reward commensurate with the merit of the action. The agents will receive a reward through their action selection, which is profit maximization for SPA and electricity bill reduction for PA. Figure 2 graphically demonstrates the interaction between the SPA and PA agent with the environment, their observation, and the action for each agent.

4.1.1 Service Provider Agent. The SPA observes the cost of buying electricity (base and reserve generation) from N_g generation facilities at time slot t , which is denoted by $F_G^t = [F_{G_1}^t, F_{G_2}^t, \dots, F_{G_{N_g}}^t]$. In addition, the SPA observes the cost of buying electricity from N_p prosumers denoted by $F_H^t = [F_{H_1}^t, F_{H_2}^t, \dots, F_{H_{N_p}}^t]$, and the total power demand of the loads P_D^t at the time slot t . Hence, environment states $s_{SPA}^t = \{F_G^t, F_H^t, P_D^t\} \in \mathcal{S}_{SPA}$ are observable by the SPA. Subsequently, the SPA

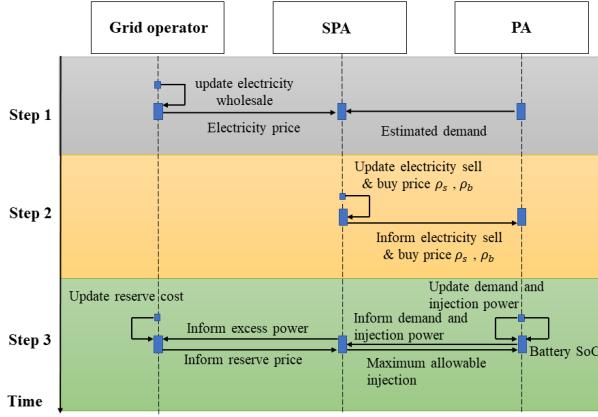


Fig. 3. Timeline of interaction between the SPA and PA. At each time slot t , the SPA dynamically determines the buy price in Step 2 (based on collected information in Step 1), and the PA decides on injected power and changing the battery state of charge in Step 3.

takes an action by adjusting the electricity buy price denoted by $a_{SPA}^t = \rho_b^t \in \mathcal{A}_{SPA}$, where \mathcal{A}_{SPA} is the finite set of available actions for SPA, i.e., all the possible buy prices.

The SPA reward function is defined based on the grid profit,

$$r_{SPA}^t = P_D^t \times \rho_s^t - \sum_{i=1}^{N_g} F_{G_i}^t - \sum_{j=1}^{N_p} \mathbb{1}_{\{P_{H_j}(t) \geq 0\}} \times P_{H_j}^t \times \rho_b^t, \quad (15)$$

where r_{SPA}^t is the SPA reward at time slot t . The value of $F_{G_i}^t$ is obtained using the incremental cost curve of the i^{th} generation facility. Given the definition of the immediate reward, the ultimate goal for the SPA is to maximize the cumulative reward $R_{SPA}^t = \sum_{k=0}^{\infty} \gamma^k r_{SPA}^{t+k+1}$ over an infinite time horizon known as expected return, where $0 \leq \gamma \leq 1$ is the discount factor.

4.1.2 Prosumer Agent. PA_j observes its own local parameters at each time slot t including power consumption $P_{C_j}^t$, state of charge of the battery SoC_j^t , and the PV generation $P_{PV_j}^t$, as well as the electricity buy price ρ_b^t which is the result of SPA action. Hence $s_{PA_j}^t = \{P_{C_j}^t, SoC_j^t, P_{PV_j}^t, \rho_b^t\} \in \mathcal{S}_{PA_j}$. Subsequently, the charge/discharge command to the energy storage in prosumer j is the action determined by PA_j , which is shown by $a_{PA_j}^t = P_{b_j}^t \in \mathcal{A}_{PA_j}$. In this case, \mathcal{A}_{PA_j} is the finite set of available actions to PA_j . The reward function for the prosumer agents is defined as,

$$r_{PA_j}^t = \mathbb{1}_{\{P_{H_j}(t) \geq 0\}} \times P_{H_j}^t \times \rho_b^t + \mathbb{1}_{\{P_{H_j}(t) < 0\}} \times P_{H_j}^t \times \rho_s^t. \quad (16)$$

Similar to the SPA, the PA_j tries to maximize its infinite-horizon accumulative reward $R_{PA_j}^t = \sum_{k=0}^{\infty} \tilde{\gamma}_j^k r_{PA_j}^{t+k+1}$, where $0 \leq \tilde{\gamma}_j \leq 1$ is the discount rate for PA_j .

The timeline of interactions between the SPA and PA agents is illustrated in Figure 3. At a given time slot t , the SPA makes a decision first by determining the market electricity buy price (a_{SPA}^t) based on the complete information of the prices of electricity generation facilities and the electricity demand from the prosumers' agents. The SPA then receives a reward (r_{SPA}^t) for taking the action. Next, the PAs observe the sell and buy price, and make the decision ($a_{PA_j}^t$) in terms of changing the battery state of charge and receive a reward ($r_{PA_j}^t$) for that decision. Finally, the PAs inform the SPA about the demand and the injected power.

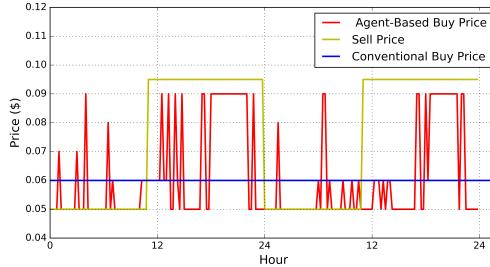


Fig. 4. Retail prices for the last two days of the simulation.

4.2 Deep Q-Network Learning

A systematic review of Deep Learning (DL) methods applied to the smart grid is provided in [36]. Deep Reinforcement Learning is the combination of Deep Neural Networks (DNNs) and RL. In general, RL solutions can be divided into value- and policy-based algorithms. The value-based solutions are straight-forward to implement. The main value-based method is Q-learning. Due to the high dimensionality nature of the problem Deep Q-Network (DQN) agents, which can handle a large state space, are used to solve their respective Markov Decision Processes (MDPs) and maximize their accumulative rewards. In addition, since the action spaces of SPA and PAs are discrete, the DQN can be used for solving this problem. The transition formula of Q-learning is,

$$Q(s^t, a^t) \leftarrow Q(s^t, a^t) + \alpha[r^{t+1} + \gamma \max_{a^{t+1}} Q(s^{t+1}, a^{t+1}) - Q(s^t, a^t)], \quad (17)$$

where α is the learning rate. Additionally, γ is the discount factor, with a range of 0 to 1. If γ is closer to zero, the agent prioritizes the immediate reward over long-term rewards. On the other hand, if γ is closer to one, the agent places greater weights over the future rewards compared with the immediate reward. The estimated Q-values are used to find the optimal policy that maximizes the accumulative rewards. An ϵ -greedy strategy is used to balance between exploration and exploitation, i.e.,

$$a^t = \begin{cases} \arg \max_{a^t} E[Q(s^t, a^t)] & \text{with probability } 1 - \epsilon, \\ \text{random action} & \text{with probability } \epsilon, \end{cases}$$

where the probability of random action is ϵ that decays from 1 to 0.01 over the training episodes in our simulations.

5 NUMERICAL RESULTS

In this study, we leverage the grid and prosumers' agents to implement the proposed DEM-DR scheme referred to as *Agent-Based* scheme. To demonstrate the efficacy of the proposed approach, we compare it with the *Conventional* approach, which simply injects the excess power to the grid when the battery is fully charged. The simulation parameters for DQN agents are tabulated in Table 2. As previously mentioned, the set of buy prices is defined as \mathcal{A}_{SPA} , representing the action space of the service provider agent. We designate the action space for prosumers as follows,

$$P_{b_j}^t = \begin{cases} P_{b_j}^{\max, \text{charge}} & \text{Charge action,} \\ 0 & \text{No charge or discharge action,} \\ P_{b_j}^{\max, \text{dis}} & \text{Discharge action.} \end{cases} \quad (18)$$

We compare the Agent-Based with the Conventional scheme through two scenarios: (1) A small-scale setting with 5 prosumers in which we demonstrate the agents' effectiveness in interacting

Table 2. DQN Hyperparameters and Simulation Parameters.

Hyperparameters	Value for SPA	Value for PA_j
Batch size	64	64
Discount factor	$\gamma=[0.95-0.99]$	$\tilde{\gamma}_j=[0.95-0.99]$
Learning rate	$\alpha=1e-3$	$\tilde{\alpha}_j=1e-3$
Soft update interpolation	1e-5	1e-5
Hidden Layer-nodes	1-[1000]	2-[1000,1000]
Activation	Tanh	Tanh
Optimizer	Adam	Adam
Simulation Parameter	Description	Value
$P_{PV_j}^{\max}$	Max. PV Generation	[2-6] kW
$P_{b_j}^{\max, \text{charge}} / P_{b_j}^{\max, \text{dis}}$	Max. allowable charge/discharge	2/-2.5 kW
$P_{H_j}^{\max}$	Max. allowable power injection	10 kW
$SoC_{b_j}^{\max}$	Max. state of charge	$0.9 \times C_{b_j}$
$SoC_{b_j}^{\min}$	Min. state of charge	$0.1 \times C_{b_j}$
C_{b_j}	Energy storage capacity	[8-15] kWh
$\phi_j(0)$	Initial state of charge	[1-4] kWh
ρ_s	Sell price [before 11am, after 11am]	[0.05, 0.095] \$/kWh
\mathcal{A}_{SPA}	Buy price set for Agent-Based scenario	{0.05, 0.06, 0.07, 0.08, 0.09, 0.1} \$/kWh
ρ_b^t	Buy price for Conventional scenario	0.06 \$/kWh
$[P_{G_1}^{\min}, P_{G_1}^{\max}]$	Limitation of base generation	[5, 45] kW
$[P_{G_2}^{\min}, P_{G_2}^{\max}]$	Limitation of reserve generation	[0, 100] kW
$[\beta_1, \beta_2]$	Transmission loss coefficient of two generators	[0.0002, 0.0002]

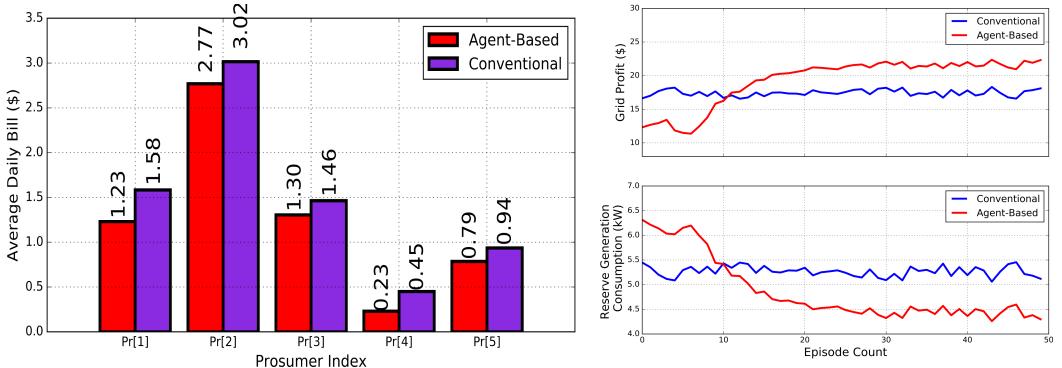


Fig. 5. (a) Average daily bill comparison of the Conventional vs. Agent-Based schemes in a small-scale microgrid with 5 prosumers. (b) Performance comparison of the Conventional vs. Agent-Based scenarios in terms of grid profit and reserve unit utilization.

with the energy marketplace environment, and improvements made to both SP and prosumers' economic benefits. (2) A medium-scale setting with 50 prosumers wherein we investigate the system's reliability and scalability by simulating a larger-scale microgrid. We implemented our DQN agents using Python3 with PyTorch 1.5.0. All simulations were performed via episodic updating across 10,000 episodes, each of which represents a 24-h cycle. A cycle consists of 96 iterations (sampling time is 15 minutes).

5.1 Scenario I: Small-Scale Simulation with 5 Prosumers

The first scenario studies the microgrid in Figure 1 consisting of two generation facilities (a base plant and a spinning reserve) managed by one distribution SP, five prosumers, and one non-generational consumer. The DQN agents are implemented for the SP and the prosumers. The operational parameters of the five prosumers and the single consumer are tabulated in Table 2. The

utilized consumption and generation profiles for the prosumers mimic real-world trends reported by California ISO [37] to exemplify real-world operation.

5.1.1 Macroscopic Evaluation of Prosumers and Grid Benefits. As reported in Table 2, the sell price for both Conventional and Agent-Based scenarios takes two values depending on the time of a day. On the other hand, the buy price is constant during a day for the Conventional scenario, while in the Agent-Based scenario the buy price is dynamically determined by the SPA. Figure 4 compares the dynamic pricing achieved by the Agent-Based approach vs. the fixed-pricing of the Conventional scheme. The results shown in Figure 5 (a) compares the average daily bill of prosumers when utilizing the two approaches. As pictured, the average daily bills have decreased by 22.1%, 8.2%, 10.9%, 48.8%, and 15.9% for prosumer 1 to prosumer 5, respectively. Figure 5 (b) on the other hand, illustrates the moving average of SP profit and reserve generation consumption. According to these results, the Agent-Based approach is offering an 25.7% increase in the SP profit, made possible by a 16% reduction in reserve power consumption, compared to the conventional approach. This further demonstrates that the proposed Agent-Based approach is capable of providing a higher profit for the SP while offering greater economic benefit for the prosumers, warranting a win-win scenario for the grid and prosumers.

5.1.2 Microscopic Evaluation of Prosumers' Behavior. Figure 6 illustrates the temporal profiles of several internal states of all five prosumers over the course of the last two simulated days for the Agent-Based and Conventional approaches. These graphs provide an intuitive interpretation of how each demand-side participants respond to the DR scheme in real-time. As pictured, using the Conventional method, the households' PV systems can charge their batteries only when their generation is more than consumption, while only being able to sell their excess generation when their batteries are fully charged.

Observations on prosumers 1 and 2: As demonstrated in Figure 6 (a) and (b), the generated power by prosumers 1 and 2 is less than their consumption (except for a very short period of time for prosumer 2). Therefore, throughout simulation of the Conventional approach, their batteries are never charging, denying them potential economic benefits otherwise possible by incorporating the energy storage system, and prohibiting them from participating in grid support during peak demand hours. On the other hand, using the Agent-Based approach, the prosumer batteries are charging during off-peak-low-price hours (i.e., beginning of the day), and prosumers are engaging in grid support during peak demand hours. Although charging the battery at the beginning of the day incurs some cost for prosumers, it eventually provides higher economic benefits to them through selling the energy back to the grid at a higher price, while assisting with grid power balance during peak demand hours. In other words, the stored energy in prosumer 1 and 2 batteries is dispatched to the grid by properly incentivizing them to participate in DR.

Observations on prosumer 3: According to consumption and generation profiles for prosumer 3 illustrated in Figure 6 (c), only a small amount of excess PV generation during a few hours is discernible. In this case, the results demonstrate that deploying the Agent-Based algorithm yields a higher SoC compared with the Conventional method. This is because the RL agent learns to charge the battery at the beginning of the day when the selling price by the grid is low, and to support the grid in the afternoon when the buy price is high.

Observations on prosumers 4 and 5: According to Figure 6 (d) and (e), these prosumers have excess PV generation during peak sun hours (e.g., around noon). The Conventional scheme fully charges the batteries of these prosumers during the peak sun hours, while the Agent-Based scheme charges the batteries starting from the beginning of the day when the sell price is low. As a result,

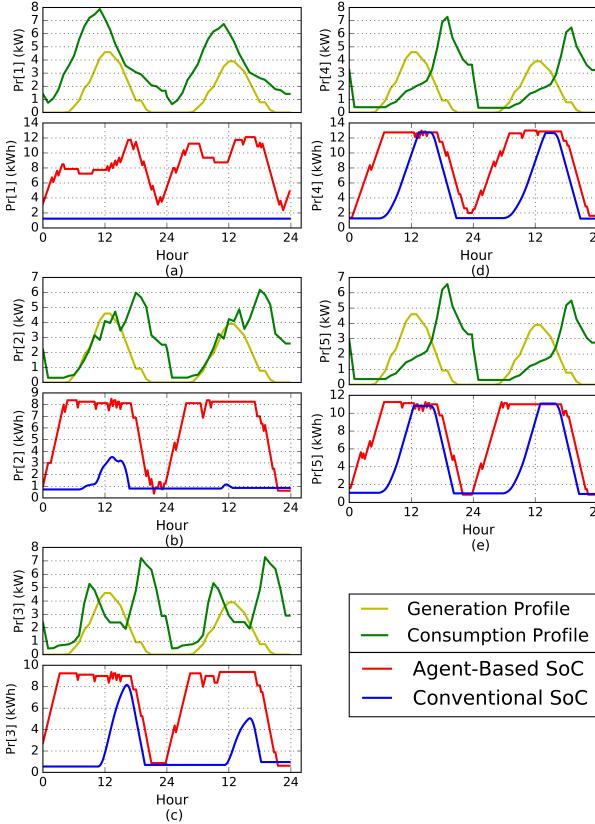


Fig. 6. Generation and consumption profiles along with the State-of-Charge (SoC) for 5 prosumers.

these prosumers can sell their excess power during sun hours, rather than storing it, at even a higher price rate than the Conventional scheme could, after fully charging the batteries.

5.1.3 Effect of Battery Size. Figure 7 demonstrates the prosumers' daily bill reduction and grid daily profit achieved by increasing the size of prosumer batteries from 2kWh to 25kWh. As pictured, the daily bill for all five prosumers decreases as the battery capacities are increased. Similarly, the grid daily profit also increases by increasing the battery capacities. Nevertheless, the improvements start to level down as the battery capacities approach to around 15kWh. Based on this observation, we use batteries with capacities in the range of 8kWh to 15kWh for all prosumers in the simulations.

Figure 8 demonstrates the impact of transmission losses on the total grid profit. In particular, the results show that as the value of transmission loss β increases, the service provider profit decreases since it has to purchase more power to meet the total demand at any given time.

5.2 Scenario II: Medium Scale Simulation with 50 Prosumers

The second scenario studies a medium scale microgrid with 50 prosumers and $N_c = 40$ non-generational consumers. In this scenario, the proposed method's scalability is the main criteria under investigation. The average daily electricity bills for 50 prosumers each with a distinct consumption profile, while using Conventional and Agent-Based approaches are illustrated in

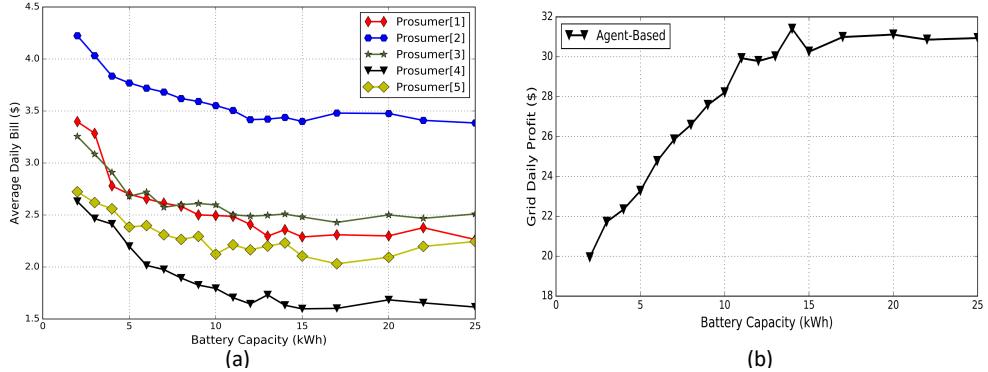


Fig. 7. (a) Effect of increasing battery capacity on the prosumers' average daily bill reduction in small-scale simulation. (b) Effect of increasing battery capacity on the grid profit improvement in small-scale simulation.

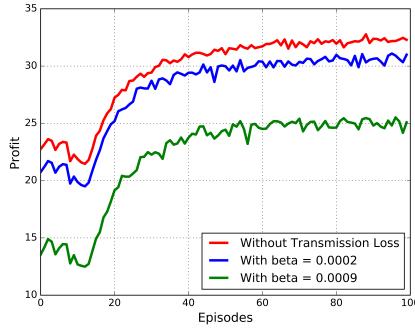


Fig. 8. Impact of transmission losses on total grid profit.

Figure 9(a). According to this figure, the average daily bill is reduced for all 50 prosumers in the Agent-Based scenario. Figure 9(b) represents the probability distribution function of the average daily bill for 50 prosumer data points. This figure illustrates an average bill of \$1.96 for Conventional method, vs. \$1.74 for Agent-Based method, amounting to around 12% bill reduction for a 24 hour billing cycle.

Figure 9 (c) compares the accumulative grid profit and reserve power utilization of the Conventional and Agent-Based scenarios in our medium scale microgrid. As pictured, the Agent-Based method provides around 4.3% higher profit over a 24-hour period. This improvement is due to the fact that the SPA and PAs learn how to dispatch the batteries' power to rely on the prosumers' PV generation, rather than using the costly reserve power. Moreover, properly incentivizing the prosumers to dispatch their stored energy has a positive impact on shaving the peak load demand of the grid, as compared in Figure 10 over a 24-hour period. As pictured, the Agent-Based approach shifts the demand at the peak hours between 18pm to 24pm to early morning hours before 6am.

6 CONCLUSION

In this paper, we propose a Deep Reinforcement Learning (DRL) framework for energy management and demand response in prosumer dominated microgrids. The service provider DRL agent dynamically changes the electricity buy price and determines the direction and amount of power flow according to the household's load demand. Further, the prosumers DRL agent controls the

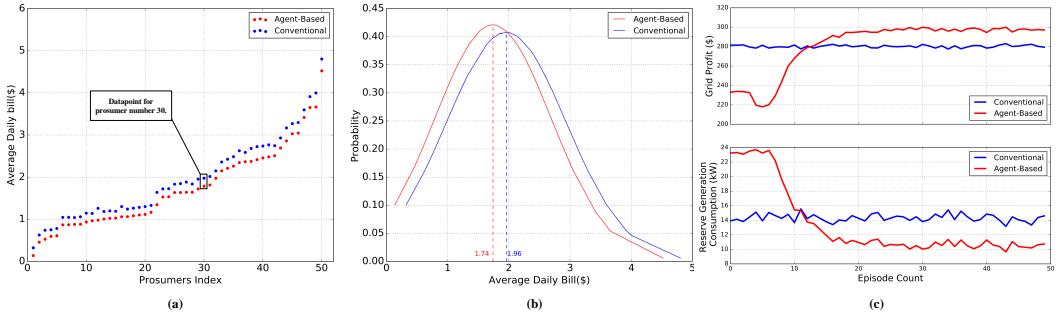


Fig. 9. Performance of the Conventional vs. Agent-Based scenarios in the medium scale microgrid with 50 prosumers. (a) Average daily bills of prosumers; (b) Distribution of average daily bills; (c) Performance comparison of the Conventional vs. Agent-Based scenarios in terms of grid profit and reserve power consumption with 50 prosumers.

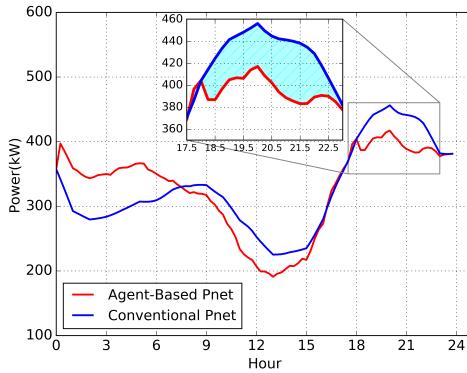


Fig. 10. Average net power comparison for the Conventional vs. Agent-Based scenarios during a 24-hour period.

battery charge and discharge rate and amount of power injection into the grid. These DRL agents collectively provide a dynamic decision-making framework. Our simulation results demonstrate that the proposed framework provides higher economic benefits for both power grid and prosumers. Specifically, properly incentivizing prosumers through dynamic pricing and leveraging the capacity of distributed battery resources result in: (i) reduced average daily bills for prosumers, (ii) enhanced profits for the grid by decreasing the reserve generation power demand, and (iii) reduced net power demands during peak hours. The proposed framework can be extended by including other external observations such as weather conditions, or investigating the performance of DRL agents under peer-to-peer energy sharing across prosumers.

REFERENCES

- [1] S. Grijalva and M. U. Tariq. Prosumer-based smart grid architecture enables a flat, sustainable electricity industry. In *ISGT 2011*, 2011. doi: 10.1109/ISGT.2011.5759167.
- [2] Nian Liu, Xinghuo Yu, Cheng Wang, and Jinjian Wang. Energy sharing management for microgrids with pv prosumers: A stackelberg game approach. *IEEE Trans. on Industrial Informatics*, 13(3):1088–1098, 2017. doi: 10.1109/TII.2017.2654302.
- [3] M Khoshjahan, R Baembitov, and M Kezunovic. Impacts of weather-related outages on der participation in the wholesale market energy and ancillary services. In *2021 CIGRE Grid of the Future Symposium, Providence, Rhode Island*,

USA, 2021.

- [4] Leong Kit Gan, PengFei Zhang, Jaehwa Lee, Michael A Osborne, and David A Howey. Data-driven energy management system with gaussian process forecasting and mpc for interconnected microgrids. *IEEE Trans. on Sustainable Energy*, 12(1):695–704, 2020.
- [5] Nikolaos G. Paterakis, Ozan Erdinç, Anastasios G. Bakirtzis, and João P. S. Catalão. Optimal household appliances scheduling under day-ahead pricing and load-shaping demand response strategies. *IEEE Transactions on Industrial Informatics*, 11(6):1509–1519, 2015. doi: 10.1109/TII.2015.2438534.
- [6] Hesam Farzaneh, Mohammad Shokri, Hamed Kebriaei, and Farrokh Aminifar. Robust energy management of residential nanogrids via decentralized mean field control. *IEEE Trans. on Sustainable Energy*, 11(3):1995–2002, 2019.
- [7] Wann-Jiun Ma, Jianhui Wang, Vijay Gupta, and Chen Chen. Distributed energy management for networked microgrids using online admm with regret. *IEEE Transactions on Smart Grid*, 9(2):847–856, 2018. doi: 10.1109/TSG.2016.2569604.
- [8] Marina Dorokhova, Yann Martinson, Christophe Ballif, and Nicolas Wyrsch. Deep reinforcement learning control of electric vehicle charging in the presence of photovoltaic generation. *Applied Energy*, 301:117504, 2021.
- [9] Jianhong Wang, Wangkun Xu, Yunjie Gu, Wenbin Song, and Tim Green. Multi-agent reinforcement learning for active voltage control on power distribution networks. *Advances in Neural Information Processing Systems*, 34, 2021.
- [10] Bala Suraj Pedasingu, Easwar Subramanian, Yogesh Bichpuriya, Venkatesh Sarangan, and Nidhisha Mahilong. Bidding strategy for two-sided electricity markets: A reinforcement learning based framework. In *Proceedings of the 7th ACM International Conference on Systems for Energy-Efficient Buildings, Cities, and Transportation*, pages 110–119, 2020.
- [11] Esmat Samadi, Ali Badri, and Reza Ebrahimpour. Decentralized multi-agent based energy management of microgrid using reinforcement learning. *International Journal of Electrical Power & Energy Systems*, 122:106211, 2020.
- [12] José R Vázquez-Canteli and Zoltán Nagy. Reinforcement learning for demand response: A review of algorithms and modeling techniques. *Applied energy*, 235:1072–1089, 2019.
- [13] Hanchen Xu, Hongbo Sun, Daniel Nikovski, Shoichi Kitamura, Kazuyuki Mori, and Hiroyuki Hashimoto. Deep reinforcement learning for joint bidding and pricing of load serving entity. *IEEE Trans. on Smart Grid*, 10(6), 2019.
- [14] Yanchang Liang, Chunlin Guo, Zhaohao Ding, and Huichun Hua. Agent-based modeling in electricity market using deep deterministic policy gradient algorithm. *IEEE Trans. on Power Systems*, 35(6):4180–4192, 2020.
- [15] Elham Foruzan, Leen-Kiat Soh, and Sohrab Asgarpoor. Reinforcement learning approach for optimal distributed energy management in a microgrid. *IEEE Trans. on Power Systems*, 33(5):5749–5758, 2018.
- [16] Aicha Dridi, Hossam Afifi, Hassine Moungla, and Jordi Badosa. A novel deep reinforcement approach for iiot microgrid energy management systems. *IEEE Transactions on Green Communications and Networking*, 2021.
- [17] Biao Wang, Yan Li, Weiyu Ming, and Shaorong Wang. Deep reinforcement learning method for demand response management of interruptible load. *IEEE Trans. on Smart Grid*, 2020.
- [18] Hwei-Ming Chung, Sabita Maharjan, Yan Zhang, and Frank Eliassen. Distributed deep reinforcement learning for intelligent load scheduling in residential smart grid. *IEEE Trans. on Industrial Informatics*, 2020.
- [19] Qinglai Wei, Zehua Liao, and Guang Shi. Generalized actor-critic learning optimal control in smart home energy management. *IEEE Trans. on Industrial Informatics*, 17(10):6614–6623, 2021. doi: 10.1109/TII.2020.3042631.
- [20] F. Ruelens, B. J. Claessens, S. Vandael, B. De Schutter, R. Babuška, and R. Belmans. Residential demand response of thermostatically controlled loads using batch reinforcement learning. *IEEE Trans. on Smart Grid*, 8(5):2149–2159, 2017. doi: 10.1109/TSG.2016.2517211.
- [21] Xu Xu, Youwei Jia, Yan Xu, Zhao Xu, Songjian Chai, and Chun Sing Lai. A multi-agent reinforcement learning based data-driven method for home energy management. *IEEE Trans. on Smart Grid*, 2020.
- [22] Jianwen Sun, Yan Zheng, Jianye Hao, Zhaopeng Meng, and Yang Liu. Continuous multiagent control using collective behavior entropy for large-scale home energy management. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 922–929, 2020.
- [23] Yaodong Yang, Jianye Hao, Yan Zheng, and Chao Yu. Large-scale home energy management using entropy-based collective multiagent deep reinforcement learning framework. In *IJCAI*, pages 630–636, 2019.
- [24] Yujian Ye, Dawei Qiu, Jonathan Ward, Marcin Abram, and C Bessiere. Model-free real-time autonomous energy management for a residential multi-carrier energy system: A deep reinforcement learning approach. In *IJCAI*, pages 339–346, 2020.
- [25] Pedram Samadi, Hamed Mohsenian-Rad, Vincent WS Wong, and Robert Schober. Real-time pricing for demand response based on stochastic approximation. *IEEE Trans. on Smart Grid*, 5(2):789–798, 2014.
- [26] Liyan Jia and Lang Tong. Dynamic pricing and distributed energy management for demand response. *IEEE Trans. on Smart Grid*, 7(2):1128–1136, 2016.
- [27] Mengmeng Yu, Seung Ho Hong, Yuemin Ding, and Xun Ye. An incentive-based demand response (dr) model considering composited dr resources. *IEEE Trans. on Industrial Electronics*, 66(2):1488–1498, 2018.

- [28] Amrit Paudel, Kalpesh Chaudhari, Chao Long, and Hoay Beng Gooi. Peer-to-peer energy trading in a prosumer-based community microgrid: A game-theoretic model. *IEEE Trans. on Industrial Electronics*, 66(8), 2019. doi: 10.1109/TIE.2018.2874578.
- [29] Kaveh Dehghanpour, M Hashem Nehrir, John W Sheppard, and Nathan C Kelly. Agent-based modeling of retail electrical energy markets with demand response. *IEEE Trans. on Smart Grid*, 9(4), 2016.
- [30] Renzhi Lu and Seung Ho Hong. Incentive-based demand response for smart grid with reinforcement learning and deep neural network. *Applied energy*, 236:937–949, 2019.
- [31] B. Kim, Y. Zhang, M. van der Schaar, and J. Lee. Dynamic pricing and energy consumption scheduling with reinforcement learning. *IEEE Trans. on Smart Grid*, 7(5):2187–2198, 2016.
- [32] Renzhi Lu, Seung Ho Hong, and Xiongfeng Zhang. A dynamic pricing demand response algorithm for smart grid: Reinforcement learning approach. *Applied Energy*, 220:220–230, 2018.
- [33] Mostafa Goodarzi and Qifeng Li. Evaluate the capacity of electricity-driven water facilities in small communities as virtual energy storage. *Applied Energy*, 309:118349, 2022.
- [34] Weerakorn Ongsakul and Vo Ngoc Dieu. *Artificial intelligence in power system optimization*. Crc Press, 2019.
- [35] Chengcheng Zhao, Jianping He, Peng Cheng, and Jiming Chen. Consensus-based energy management in smart grid with transmission losses and directed communication. *IEEE Trans. on Smart Grid*, 8(5), 2017. doi: 10.1109/TSG.2015.2513772.
- [36] Mohamed Massaoudi, Haitham Abu-Rub, Shady S Refaat, Ines Chih, and Fakhreddine S Oueslati. Deep learning in smart grid technology: A review of recent advancements and future prospects. *IEEE Access*, 9, 2021.
- [37] California ISO. Current and forecasted demand. URL <http://www.caiso.com/TodaysOutlook/Pages/default.aspx>.