

# A Review of Deep Reinforcement Learning for Smart Building Energy Management

Liang Yu, *Member, IEEE*, Shuqi Qin, Meng Zhang, Chao Shen, *Senior Member, IEEE*,  
Tao Jiang, *Fellow, IEEE*, and Xiaohong Guan, *Fellow, IEEE*

## Abstract

Global buildings account for about 30% of the total energy consumption and carbon emission, raising severe energy and environmental concerns. Therefore, it is significant and urgent to develop novel smart building energy management (SBEM) technologies for the advance of energy-efficient and green buildings. However, it is a nontrivial task due to the following challenges. Firstly, it is generally difficult to develop an explicit building thermal dynamics model that is both accurate and efficient enough for building control. Secondly, there are many uncertain system parameters (e.g., renewable generation output, outdoor temperature, and the number of occupants). Thirdly, there are many spatially and temporally coupled operational constraints. Fourthly, building energy optimization problems can not be solved in real-time by traditional methods when they have extremely large solution spaces. Fifthly, traditional building energy management methods have respective applicable premises, which means that they have low versatility when confronted with varying building environments. With the rapid development of Internet of Things technology and computation capability, artificial intelligence technology find its significant competence in control and optimization. As a general artificial intelligence technology, deep reinforcement learning (DRL) is promising to address the above challenges. Notably, the recent years have seen the surge of DRL for SBEM. However, there lacks a systematic overview of different DRL methods for SBEM. To fill the gap, this paper provides a comprehensive review of DRL for SBEM from the perspective of system scale. In particular, we identify the existing unresolved issues and point out possible future research directions.

## Index Terms

Deep reinforcement learning, artificial intelligence, Internet of things, smart buildings, energy management, uncertainty, building microgrids.

## I. INTRODUCTION

Buildings account for a large portion of total energy consumption and total carbon emission in the world [1]–[5]. For example, global buildings consumed 30% of total energy and generated 28% of total carbon emission in 2019 [6]. Moreover, the energy demand of buildings is expected to increase by 50% in the next 30 years [7] [8]. Under the above background, smart buildings have received more and more attention in recent years, which can provide sustainable, economical, and comfortable operational environments for occupants using many advanced technologies, e.g., Internet of Things (IoT), cloud computing, machine learning, and big data analytics [9]–[11]. For supporting the above features, it is significant and urgent to develop novel smart building energy management (SBEM) technologies [12], which can implement the optimal tradeoff among energy consumption, carbon emission, energy cost, and user comfort [13]–[18] by intelligently scheduling building energy systems.

Although SBEM has many advantages, the following challenges have to be addressed. Firstly, due to the existence of many complex and random factors, it is often intractable to develop an explicit building thermal dynamics model that is accurate and efficient enough for building energy optimization [19]. Secondly, there are many uncertain system parameters [20], e.g., renewable generation output, electricity price, indoor temperature, outdoor temperature, CO<sub>2</sub> concentration, and the number of occupants. Thirdly, there are many temporally and spatially coupled operational constraints related to energy subsystems [21] [22], e.g., heating, ventilation, and air conditioning (HVAC) systems, and energy storage systems (ESSs), which means that the current system decision will affect the future decisions and the decisions among different subsystems should be coordinated. Fourthly, it is difficult to solve large-scale building energy optimization problems in real-time when traditional optimization methods are adopted [23]. To be specific, any time when an optimization is needed, these methods have to compute completely or partially all the possible solutions and choose the best one. When the solution space is very large, the computation process is time-consuming [21]. Finally, it is hard to develop a generalized building energy management method that can be applied in all building environments [18]. In existing SBEM methods, most of them have strong applicable premises [24], e.g., stochastic

L. Yu is with the College of Automation & College of Artificial Intelligence, Nanjing University of Posts and Telecommunications, Nanjing 210003, China, and also with Xi'an Jiaotong University, Xi'an 710049, China. (email: liang.yu@njupt.edu.cn)

S. Qin is with the College of Internet of Things, Nanjing University of Posts and Telecommunications, Nanjing 210003, China.

M. Zhang and X. Guan are with Systems Engineering Institute, Ministry of Education Key Lab for Intelligent Networks and Network Security, Xi'an Jiaotong University, Xi'an 710049, China.

C. Shen is with the School of Cyber Science and Engineering, Xi'an Jiaotong University, Xi'an 710049, China. (email: cshen@sei.xjtu.edu.cn)

T. Jiang is with Wuhan National Laboratory for Optoelectronics, School of Electronic Information and Communications, Huazhong University of Science and Technology, Wuhan 430074, China. (email: Tao.Jiang@ieee.org)

programming and model predictive control (MPC) need the prior or forecasting information of uncertain parameters [25] [26], and Lyapunov optimization techniques require some strict usage conditions [16] [27].

TABLE I  
THE COMPARISON BETWEEN OUR WORK AND RELATED SURVEYS

Literature	Main focus	System type(s)	Involved Methods/Algorithms	DRL methods for SBEM classified	Future directions in DRL-based SBEM provided
Han <i>et al.</i> [57]	Occupant comfort control	HVAC, lighting systems	RL	No	No
Leitão <i>et al.</i> [58]	Building energy optimization	Smart home	LP, NLP, CP, DP, GA, PSO, MPC, RL	No	No
Mason <i>et al.</i> [59]	Building energy optimization	HVAC, EWH, home management systems, smart home	RL	No	No
Wang <i>et al.</i> [60]	Building energy optimization	HVAC, batteries, home appliances, EWH, windows, lighting	RL	No	No
Rajasekhar <i>et al.</i> [61]	Building energy optimization	HVAC	RNN, WNN, RT, SVM, PSO, MPC, FL, RL, DQN	No	No
Zhang <i>et al.</i> [55]	Cyber security, demand response, load forecasting, and microgrid	Smart grid	RL, DQN, DDPG, NAF, A3C	No	No
Yang <i>et al.</i> [56]	Energy and electric system security, operation optimization	Microgrid, ESS, HVAC, home appliances, PV	RL, DQN, DDPG, A3C, DDQN, TRPO	No	No
Our work	Building energy optimization	A single building energy subsystem, multiple energy subsystems in buildings, building microgrids	DQN, DDQN, BDQ, DDPG, PDDPG, MADDPG, FH-DDPG, A2C, A3C, TRPO, PPO, MAPPO, MuZero, MAAC, EB-C-A2C, EB-C-DQN	Yes	Yes

As a general artificial intelligence technology, deep reinforcement learning (DRL) [28] [29] is promising to address the above challenges and has been applied in many fields, e.g., games [88] [31], autonomous driving [32]–[35], autonomous IoT [36], smart buildings [19] [20] [37] [38], smart city [39], wireless networks [40], Internet of energy [41], unmanned aerial vehicles [42], smart microgrids [43], edge computing [44], and manufacturing systems [45]. In 2017, the first work that adopts DRL algorithm for SBEM has been done [19]. To be specific, a deep Q-network (DQN) algorithm has been adopted for the control of building HVAC systems and simulation results have showed the effectiveness of the designed control algorithm in reducing energy cost and maintaining thermal comfort of occupants. Since then, many DRL-based methods for SBEM have been proposed [20] [21] [46] [47]. In general, DRL-based methods have following advantages in dealing with the above-mentioned five challenges:

- **For challenge 1:** Based on the information interacted with actual building environments, DRL agents can learn the optimal control policies by trial and error. Therefore, DRL can support system operation without knowing explicit building thermal dynamics models [19].
- **For challenge 2:** After the training process is finished, the trained DRL agent will be used for performance testing. Given the current state of an actual environment, the DRL agent will generate an action via a mapping function. Since no forecasting or statistics information of building environments is used in the above process, DRL can tackle system uncertainties [20] [48].
- **For challenge 3:** By designing proper reward functions, building energy subsystems can coordinate with each other under the framework of multi-agent DRL. As a result, spatially-coupled operational constraints are guaranteed [21]. Moreover, by choosing reasonable actions or designing efficient reward functions, temporally-coupled operational constraints related to energy subsystems (e.g., HVAC systems and ESSs) can be satisfied [20].
- **For challenge 4:** During the testing phase, the computational complexity of the DRL algorithm is very low since just the forward propagation in deep neural networks (DNNs) is involved. Even if a high-dimensional raw state is given, the optimal control actions can be determined instantly (e.g., few milliseconds) [23] [49] [50].
- **For challenge 5:** Since simulated or real data is used for training agents, the applications of DRL methods do not require rigorous mathematical models and premise conditions. Moreover, the trained DRL agent can still work or even be improved persistently by online learning when confronted with varying building environments [23] [51]. Thus, DRL methods have wide applicable premises in solving SBEM problems.

There are many surveys related to DRL in the literature, e.g., the applications of DRL in power and electric systems,

communications and networking, autonomous driving, autonomous IoT, cyber security, and multi-agent systems can be found in [24], [33]–[36], [52]–[56]. However, they do not consider DRL for SBEM. Although there are several surveys on building energy systems, the involved methods are RL [57]–[60] or other artificial intelligence methods (e.g., MPC and fuzzy logic (FL)) [61]. Based on the above observation, we are motivated to conduct a comprehensive review on DRL for SBEM. For convenience, we provide the comparison between our work and related surveys in Table I. It can be observed that our work mainly focuses on DRL for SBEM from the perspective of different building system scales (i.e., a single building energy subsystem, multiple energy subsystems in buildings, and building microgrids). Moreover, we provide a systematic overview of different DRL methods for SBEM. Above all, we identify the existing unresolved issues and point out possible future directions. We hope that this paper can show some insights in this direction and raise the attention of SBEM research community to explore and exploit DRL, as another alternative or even a better solution for SBEM.

The rest of this paper is organized as follows. In Section II, we give an overview of DRL. In Section III, we introduce the background of SBEM, the procedure of solving SBEM problems using DRL, and the classification of DRL methods for SBEM. In next three sections, we discuss DRL applications in a single building energy subsystem, multiple energy subsystems of buildings, and building microgrids. In Section VII, we identify some unsolved issues and point out the future research directions. Finally, conclusions and lessons learned are provided in Section VIII. For easy understanding, the list of abbreviations in alphabetical order is provided in Table II.

## II. AN OVERVIEW OF DEEP REINFORCEMENT LEARNING

According to the ways of feedback, machine learning can be divided into three types as shown in Fig. 1, i.e., supervised learning, unsupervised learning, and reinforcement learning (RL) [62]. As for supervised learning, an immediate feedback can be obtained by comparing the prediction value with the real value that given by the label data, which will be used to improve predictor. In contrast, no feedback can be received in unsupervised learning since the input data is not labeled. While interacting with an environment, a delayed feedback is involved in RL since the action taken at the current state will affect future states and actions, and the value of taking an action at the current state can not be known immediately but be learned gradually. Typically, supervised learning and unsupervised learning are used to solve single-stage problems (e.g., regression, classification, clustering, and dimension reduction), but RL is specialized in solving multi-stage decision problems [62]. Under the background of SBEM, supervised learning can be used to develop building thermal dynamics models and reward models. Based on these models, RL can reduce the number of interactions with the environment, resulting in a high sampling efficiency.

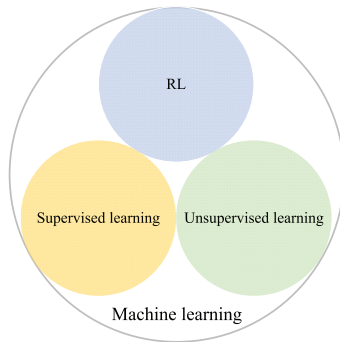


Fig. 1. Classification of machine learning

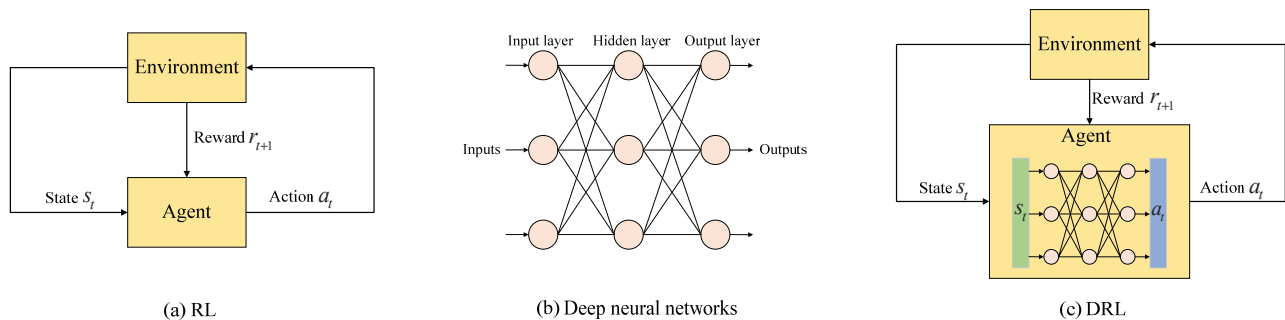


Fig. 2. Illustration of RL, deep learning, and DRL

DRL can be regarded as the combination of deep learning and RL as shown in Fig. 2. To be specific, DNNs are used to approximate the optimal value functions or optimal policies in RL. Therefore, DRL has a powerful representation ability and

TABLE II  
LIST OF ABBREVIATIONS IN ALPHABETICAL ORDER

Abbreviation	Description
A2C	Advantage Actor-Critic
A3C	Asynchronous Advantage Actor-Critic
AHU	Air Handling Unit
BAS	Building Automation System
BDQ	Branching Dueling Q-Network
BEM	Building Energy Model
CD	Clothes Dryer
CNN/DNN	Convolutional/Deep Neural Network
CP	Convex Programming
D-DNFQI	Double Deep Neural Fitted Q Iteration
DDPG	Deep Deterministic Policy Gradient
DDQN	Double Deep Q-Network
DG	Diesel Generator
DP	Dynamic Programming
DQN	Deep Q-Network
DRL	Deep Reinforcement Learning
DW	Dishwasher
EB-C-A2C	Entropy-Based Collective Advantage Actor-Critic
EB-C-DQN	Entropy-Based Collective Deep Q-Network
EHP	Electric Heat Pump
ESS	Energy Storage System
EV	Electric Vehicle
EWH	Electric Water Heater
FH-DDPG	Finite-Horizon Deep Deterministic Policy Gradient
FH-RDPG	Finite-Horizon Recurrent Deterministic Policy Gradient
FL	Fuzzy Logic
GA	Genetic Algorithm
GB	Gas Boiler
HVAC	Heating, Ventilation and Air Conditioning
IoT	Internet of Things
LP	Linear Programming
LSTM	Long Short-Term Memory
MAAC	Multi-Actor Attention-Critic
MADDPG	Multi-Agent Deep Deterministic Policy Gradient
MAPPO	Multi-Agent Proximal Policy Optimization
MCTS	Monte-Carlo Tree Search
MDP	Markov Decision Process
MPC	Model Predictive Control
NAF	Normalized Advantage Functions
NLP	Non-Linear Programming
PDDPG	Prioritized Deep Deterministic Policy Gradient
PILCO	Probabilistic Inference for Learning Control
PPO	Proximal Policy Optimization
PSO	Particle Swarm Optimization
PV	Photovoltaic
RL	Reinforcement Learning
RNN	Recurrent Neural Network
RT	Regression Tree
SBEM	Smart Building Energy Management
SVM	Support Vector Machine
TES	Thermal Energy Storage
TRPO	Trust Region Policy Optimization
VAV	Variable Air Volume
WM	Washing Machine
WNN	Wavelet Neural Network
WT	Wind Turbine



strong decision-making ability under uncertainty [63] [64]. Since DRL algorithms are mainly based on Markov decision process (MDP) framework or its variants (e.g., partially observable MDP [36] and Markov game [21]), we first give the background of MDP. Then, we introduce some terms (e.g., *policy*, *action-value function*, *experience replay*, and *target network*) in RL and DRL, which will be mentioned frequently in next several sections.

### A. MDP

Typically, an MDP is defined by a five-tuple  $(\mathcal{S}, \mathcal{A}, P, \mathcal{R}, \gamma)$ , where  $\mathcal{S}$  and  $\mathcal{A}$  denote the sets of state and action, respectively.  $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  denotes the state-transition probability function  $P(s'|s, a)$  ( $s', s \in \mathcal{S}, a \in \mathcal{A}$ ), which models the uncertainty in the evolution of system states based on the action taken by the agent.  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward function and  $\gamma \in [0, 1]$  is a discount factor. Note that MDP provides a mathematical framework for multi-stage optimal decision problems under uncertainty. In other words, the decision maker (i.e., the agent) observes a state  $s_t$  and takes an action  $a_t$  at each time slot  $t$ . Next, the state of the system (i.e., the environment) evolves into another one. Then, the agent finds itself in a new state  $s_{t+1}$  and receives a reward  $r_{t+1}$ . In addition, the aim of the agent at time slot  $t$  is to maximize the expected return it receives over the future [62], which is given by  $\sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$ .

### B. RL

RL has been widely used in solving MDPs [57], [62], [65]–[70]. In an RL process, the agent learns its optimal policy  $\pi$  by interacting with the environment, where a *policy*  $\pi$  is a mapping from states to the probabilities of selecting every possible action [62]. In particular, the agent observes a state and takes an action at slot  $t$ . Then, it receives a reward and a new state, which are used to update the policy. The above process repeats until the policy converges.

To better illustrate the key idea of RL, Q-learning is introduced in this subsection, which is one of the most classic RL algorithms that learn a deterministic policy indirectly. In other words, Q-function (i.e., *action-value function*) is learned for selecting decisions instead of policy function itself. Let the value of taking action  $a$  in state  $s$  under a policy  $\pi$  be  $Q_\pi(s, a)$ , which is defined by

$$Q_\pi(s, a) \doteq \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} (s_t = s, a_t = a) \right], \quad (1)$$

where  $\mathbb{E}_\pi[\cdot]$  denotes the expected value of a random variable given that the agent follows policy  $\pi$ . Then, the optimal action-value function  $Q^*(s, a)$  is  $\max_\pi Q_\pi(s, a)$  and can be calculated by the Bellman optimality equation in a recursive manner [20] as follows,

$$\begin{aligned} Q^*(s, a) &= \mathbb{E}[r_{t+1} + \gamma \max_{a'} Q^*(s_{t+1}, a') | s_t = s, a_t = a] \\ &= \sum_{s', r} P(s', r | s, a) [r + \gamma \max_{a'} Q^*(s', a')], \end{aligned}$$

where  $s' \in \mathcal{S}$ ,  $r \in \mathcal{R}$ ,  $a' \in \mathcal{A}$ , and  $P(s', r | s, a)$  denotes a conditional probability function. To obtain the value of  $Q^*(s, a)$ , the information of  $P(s', r | s, a)$  must be known, which may be unavailable in practice. To address this challenge, Q-learning algorithm is proposed to approximate  $Q^*(s, a)$  as follows,

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \Delta_t, \quad (2)$$

where  $\Delta_t = \alpha [r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t)]$  and  $\alpha$  is the step size. It is obvious that  $Q(s_t, a_t) = r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a')$  when  $\Delta_t = 0$ . At this time,  $Q(s_t, a_t)$  will not be updated and the learned action-value function  $Q$  directly approximates the optimal action-value function  $Q^*(s, a)$ . Note that Q-learning algorithm is effective when state space is low-dimensional. To support high-dimensional state space, a nonlinear function approximator such as a neural network can be used to represent the action-value function in RL. At this time, RL is known to be unstable or even divergent.

### C. DRL

As the first DRL algorithm, DQN can overcome the above-mentioned drawback of Q-learning by adopting several techniques of stabilizing learning process, e.g., *experience replay* and *target network* [28]. To be specific, experience replay mechanism stores the experience transitions  $(s_t, a_t, s_{t+1}, r_{t+1})$  in a replay memory and draw samples of them uniformly at random for training, which brings greater data efficiency when compared with standard online Q-learning algorithm. Moreover, randomizing the samples contributes to the reductions of their correlations and the variance of updating DNN weights. In addition, target network is adopted to improve the stability of training process by copying a separate network with longer update period for the computation of target value (i.e.,  $r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a')$ ).

In addition to DQN, many DRL algorithms have been proposed in existing works. Generally speaking, these DRL algorithms can be classified into two types, i.e., model-free DRL algorithms and model-based DRL algorithms. In particular, model-free DRL algorithms do not need to know environment models (i.e., state-transition probability function  $P(s'|s, a)$  or  $P(s', r|s, a)$ ) since they learn policies based on the information directly interacted with unknown environments. Different from model-free DRL algorithms, model-based DRL algorithms need to construct environment models. According to the way of learning a

policy, model-free DRL methods can be further divided into value-based methods and policy-based methods. To be specific, the former learns an approximation of optimal value function (i.e., learn a deterministic policy indirectly), while the latter learns an approximation of optimal policy directly. Typically, value-based methods sometimes update value function in an “off-policy” (i.e., the policy to be evaluated and improved is unrelated to the policy used for sampling an action at the next state) manner, which means that the previous collected experience transitions in the same environment can be used for training and a high data efficiency can be achieved. In contrast, “on-policy” means that all of the updates are made using the data from the trajectory distribution induced by the current policy [71]. Therefore, “on-policy” methods are more stable but less data-efficient compared with “off-policy” methods. In addition, according to the number of agents, DRL algorithms can be divided into two types, i.e., single-agent and multi-agent DRL algorithms. In Table III, typical DRL algorithms and their categories are summarized.

TABLE III  
CLASSIFICATION OF TYPICAL DRL ALGORITHMS

DRL Algorithms	Category	Source
DQN	Value-based, off-policy, single-agent	[28]
Double DQN	Value-based, off-policy, single-agent	[72]
Dueling DQN	Value-based, off-policy, single-agent	[73]
Prioritized DQN	Value-based, off-policy, single-agent	[74]
Distributional DQN	Value-based, off-policy, single-agent	[75]
Noisy DQN	Value-based, off-policy, single-agent	[76]
Rainbow	Value-based, off-policy, single-agent	[77]
DDPG	Policy-based, off-policy, single-agent	[78]
MADDPG	Policy-based, off-policy, multi-agent	[79]
MAAC	Policy-based, off-policy, multi-agent	[80]
Soft Actor-Critic	Policy-based, off-policy, single-agent	[81]
A3C	Policy-based, on-policy, single-agent	[82]
TRPO	Policy-based, on-policy, single-agent	[83]
PPO	Policy-based, on-policy, single-agent	[84]
MAPPO	Policy-based, on-policy, multi-agent	[85]
Deep PILCO	Model-based, on-policy, single-agent	[86]
World Model	Model-based, on-policy, single-agent	[87]
MuZero	Model-based, off-policy, single-agent	[88]

### III. DRL-BASED SMART BUILDING ENERGY MANAGEMENT

In this section, we briefly introduce the main research problems in the field of SBEM. Then, we classify the representative DRL algorithms for SBEM by pointing out their respective advantages, disadvantages, and application scenarios, which contributes to the selection of appropriate DRL algorithms for SBEM.

#### A. SBEM

In smart buildings, there are several types of energy equipments, e.g., photovoltaic panels (PVs), wind turbines (WTs), diesel generators (DGs), electric energy storage systems, thermal energy storage systems, HVAC systems, lighting systems, blind systems, window systems, electric water heaters (EWHs), electric vehicles (EVs), washing machines (WMs), gas boilers (GBs), and clothes dryers (CDs). Since the operations of such equipments have considerable economic, environmental, and social impacts on buildings, it is very necessary to schedule them coordinately.

Considering that HVAC systems have high power consumption and can be adjusted flexibly without sacrificing user comfort, they are taken as an example to illustrate a typical research problem in SBEM field. When considering economic and social impacts, a comfort-aware energy cost minimization problem related to an HVAC system in a  $N$ -zone commercial building can be formulated by **P1** as follows [21],

$$(\mathbf{P1}) \quad \min_{m_{i,t}, \sigma_t} \sum_{t=1}^L \mathbb{E}\{C_t(m_{i,t}, \sigma_t)\} \quad (3a)$$

$$s.t. \quad T_i^{\min} \leq T_{i,t} \leq T_i^{\max}, \quad \forall i, t, \quad K_{i,t} > 0, \quad (3b)$$

$$T_{i,t+1} = \mathcal{F}(T_{i,t}, T_{z,t} | \forall z \in \mathcal{N}_i, T_t^{\text{out}}, m_{i,t}, S_{i,t}), \quad (3c)$$

$$O_{i,t} \leq O_i^{\max}, \quad \forall i, t, \quad K_{i,t} > 0, \quad (3d)$$

$$O_{i,t+1} = \mathcal{G}(O_{j,t} | \forall j \in \mathcal{N}, K_{i,t}, m_{i,t}, \sigma_t), \quad (3e)$$

$$m_{i,t} \in \mathcal{M}, \quad \forall i, t, \quad (3f)$$

$$\sigma_t \in \Omega, \quad \forall t, \quad (3g)$$

where  $\mathbb{E}$  denotes the expectation operator, which acts on random system parameters, e.g., outdoor temperature  $T_t^{\text{out}}$ , number of occupants  $K_{i,t}$ . Decision variables of **P1** are air supply rate in each zone  $m_{i,t}$  and damper position in air handling unit

(AHU)  $\sigma_t$ ,  $L$  is the considered total number of time slots.  $C_t(m_{i,t}, \sigma_t)$  is the energy cost at slot  $t$ ,  $T_{i,t}$  and  $O_{i,t}$  are indoor air temperature and indoor CO<sub>2</sub> concentration at slot  $t$ , respectively. It is obvious that they should be controlled within comfortable ranges, which can be captured by (3b) and (3d), respectively. The dynamics of  $T_{i,t}$  and  $O_{i,t}$  are represented by (3c) and (3e), respectively. Note that  $\varsigma_{i,t}$ ,  $\mathcal{N}_i$ , and  $\mathcal{N}$  are thermal disturbance, the set of neighbors related to zone  $i$  ( $1 \leq i \leq N$ ), and the set of zones, respectively. The discrete solution spaces of  $m_{i,t}$  and  $\sigma_t$  are shown in (3f) and (3g), respectively.

To solve SBEM problem **P1** efficiently, several challenges mentioned in Section I have to be addressed. In addition, non-convexity and non-separability of the objective function increase the difficulty of solving **P1**. When taking all challenges into consideration, existing building energy optimization approaches are not applicable. Note that the above example is just related to the management of a single building energy subsystem. With the increase of system scale, more and more challenges are involved, which will be discussed in next four sections.

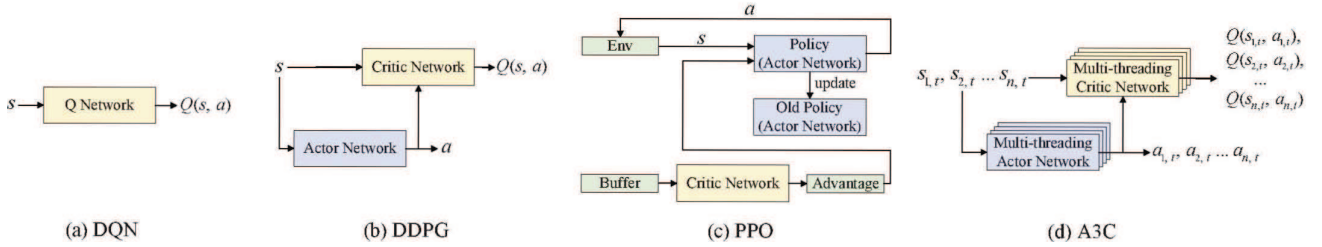


Fig. 3. Illustration of network architectures of four representative model-free DRL algorithms

### B. Procedure of Solving SBEM Problems using DRL

To solve SBEM problems using DRL methods, several steps can be taken as follows:

**Step 1:** Reformulating the original problem (e.g., **P1**) as an MDP problem or a Markov game (i.e., a multi-agent extension of MDP). Take **P1** for example,  $N + 1$  agents are adopted for the purpose of scalability, since there are  $N + 1$  decision variables and the solution space grows rapidly with the increase of zone number. Thus, **P1** should be formulated as a Markov game and its components (e.g., state, action, and reward function) should be designed.

**Step 2:** Designing an appropriate DRL-based algorithm for the reformulated problem. For instance, in order to solve the Markov game related to **P1**, an HVAC control algorithm has been proposed in [21] based on multi-agent DRL with attention mechanism, which is scalable to the number of agents.

**Step 3:** Analyzing the computational complexity of the designed DRL-based algorithm in the training and testing phases. Note that a trained DRL algorithm has very low computational complexity in the testing phase since just forward propagation in DNNs is involved, e.g., DRL agent can make real-time decisions within several milliseconds given a high-dimensional environment state. As a result, computational complexity analysis of the designed DRL-based algorithm in the training phase is the priority. Typical factors that affect computational complexity of a DRL-based energy management algorithm in the training phase are summarized as follows [20] [89], e.g., the number of hidden layers, the number of neurons in each hidden layer, the number of training episodes needed for algorithm convergence, the frequency of updating weights, and batch size.

**Step 4:** Evaluating the performances of the designed DRL-based algorithm from different perspectives, e.g., convergence, effectiveness, scalability, and robustness.

### C. Representative DRL Algorithms for SBEM

In this subsection, we introduce some representative DRL algorithms for SBEM in Table IV. According to the descriptions in Section II-C, these algorithms can be divided into two types, i.e., model-free algorithms and model-based algorithms.

Model-free algorithms do not require explicit building environment models, but they need to collect sufficient experience transitions for training, which may result in a long exploration time and a high exploration cost. In Table IV, four representative model-free DRL algorithms for SBEM are given, i.e., DQN, DDPG, PPO, and A2C/A3C. To be specific, DQN is a value-based off-policy algorithm and only supports discrete action space, while DDPG is a policy-based off-policy algorithm and only supports continuous action space. Since experience replay is adopted by DQN and DDPG, they have higher data efficiency compared with “on-policy” algorithms. However, they tend to overestimate Q-value and generate sub-optimal policies. As two representative policy-based off-policy algorithms, PPO and A2C/A3C can support both discrete and continuous action spaces. Moreover, PPO can support stable learning by controlling the similarity between the current policy and the old policy. Furthermore, it is robust to hyperparameters and network architectures. Although A2C/A3C can support reliable and parallel learning on a single multi-core CPU, it is sensitive to the employed hyperparameters. To show the differences among these algorithms clearly, their network architectures are illustrated in Fig. 3.

Model-based algorithms need to construct models to simulate building environments and use them to generate future episodes for training. Therefore, model-based algorithms outperform model-free algorithms in terms of sample complexity. However,

TABLE IV  
REPRESENTATIVE DRL ALGORITHMS FOR SBEM

Model-free	Description	Agents learn optimal policies based on the interaction information with unknown environment			
	Advantages	No need to know building environment model, high data-efficiency			
	Disadvantages	Require a large number of samples, long exploration time, and high exploration cost			
	Representative algorithms, features and categories	DQN	DDPG	PPO	A2C/A3C
		Support only for discrete action space	Support only for continuous action space	Support stable training and discrete/continuous action space	Support fast training and discrete/continuous action space
Application scenarios	Value-based, off-policy	Policy-based, off-policy	Policy-based, on-policy	Policy-based, on-policy	
Model-based	Description	Agents construct a stimulated environment model and use it to generate future episodes for training			
	Advantages	High sample efficiency			
	Disadvantages	Develop an accurate and useful model is often challenging			
	Representative algorithms, features and categories	MuZero	LSTM-DDPG	Differentiable MPC-PPO	BEM-A3C
		Learn a network model with accurate planning performance	Use LSTM and historical data to learn transition function and reward function	Pre-train a differentiable MPC policy based on imitation learning, which is both sample-efficient and interpretable	Use EnergyPlus and measured data to simulate building environment
Application scenarios	Model-based, off-policy	Model-based, off-policy	Model-based, on-policy	Model-based, on-policy	
	Microgrid energy management [113]	HVAC control [90]	HVAC control [91]	HVAC control [92]	

for model-based algorithms, it is often challenging to obtain an accurate building environment model. In existing works, many model-based DRL methods for SBEM have been proposed by constructing a system dynamics model with historical operational data [90] [91] or calibrating a building environment simulation model with the measured data [92]. Four representative model-based DRL algorithms for SBEM are provided in Table IV, i.e., MuZero, LSTM-DDPG, Differentiable MPC-PPO, and BEM-A3C. To be specific, MuZero is a model-based off-policy algorithm [93], which intends to learn a network model with accurate planning performance. Since it uses tree-based search methods, MuZero may not be good at dealing with continuous action space. LSTM-DDPG is also a model-based off-policy algorithm, which uses LSTM and historical operational data to learn the environment model. Then, the obtained model is used to generate a large number of data for training DRL agents with DDPG algorithm. Differentiable MPC-PPO is a model-based on-policy algorithm, which uses differentiable MPC to learn existing controller via imitation learning and improves the learned controller via PPO algorithm. BEM-A3C is also a model-based on-policy algorithm, which uses EnergyPlus software to develop a building energy model and uses actual operational data to calibrate the model. Finally, the calibrated model can be used for training DRL agents via A3C algorithm.

In the next three sections, we will introduce DRL applications in SBEM considering different system scales as shown in Fig. 4, i.e., a single building energy subsystem, multiple energy subsystems (MES) in buildings, and building energy systems in microgrid environment.

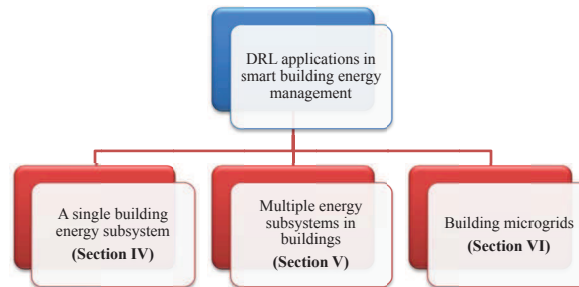


Fig. 4. Taxonomy of DRL Applications in SBEM

#### IV. APPLICATIONS OF DRL IN A SINGLE BUILDING ENERGY SUBSYSTEM

In existing works, DRL techniques have been adopted to optimize the operation cost or energy consumption of a single building energy subsystem. Among all single building energy subsystems, HVAC and EWH [94]–[96] have very flexible power consumption [10]. Therefore, we mainly focus on them in this section. To be specific, model-free DRL methods for HVAC control, model-based DRL methods for HVAC control, and DRL methods for EWHs are introduced in section IV-A,

section IV-B, and section IV-C, respectively. At the end of this section, we give a summary of existing works and provide some insights.

#### A. Model-free DRL Methods for HVAC Control

It is well known that the main purpose of an HVAC system is to maintain thermal comfort for occupants. To achieve this aim, many DRL-based methods have been proposed. For example, Morinibu *et al.* proposed an A2C-based method to decrease the non-uniformity of radiation temperature in a room by flexibly controlling the wind direction of an HVAC system [97]. Simulation results showed that the proposed method has better performance than random control and normal control. Since the operations of HVAC systems place an economic burden on building operators, it is very necessary to minimize energy cost while maintaining thermal comfort for occupants. In [19], Wei *et al.* proposed a DQN-based HVAC control method to save energy cost in office buildings while maintaining the room temperature requirements. When 5 zones are considered, energy cost can be reduced by 35.1%. Similarly, Nagy *et al.* and Gupta *et al.* proposed a model-free DRL-based HVAC control method in a residential building to save energy cost and reduce the loss of occupant comfort based on D-DNFQI and DQN, respectively [51] [98]. In addition to energy cost, energy consumption is also an important metric. To minimize energy consumption while maintaining thermal comfort, some works have been done in [99] [100]. Since the above-mentioned works use value-based DRL methods, they can not deal with continuous actions. To support continuous actions, Gao *et al.* presented a DDPG-based HVAC control method to optimize energy consumption and thermal comfort in a laboratory by jointly adjusting temperature set-point and humidity [18]. Simulation results showed that the proposed method has higher thermal comfort and energy-efficiency than other baselines, e.g., Q-learning and DQN. Due to the importance of indoor air quality, Valladares *et al.* proposed a DDQN-based control algorithm to optimize HVAC energy consumption while maintaining thermal comfort and indoor air quality comfort for occupants [101]. Although the above-mentioned methods are effective, they can not be used for coordinating multiple components in HVAC systems.

To overcome this drawback, Nagarathinam *et al.* [49] proposed a multi-agent DRL based algorithm to minimize HVAC energy consumption without sacrificing user comfort by adjusting both the building and chiller set-points. To be specific, each DDQN-based agent coordinates with others to learn an optimal HVAC control policy. Note that the coordination is achieved by allocating the same reward for each agent. Since a large building may have a few hundreds of AHUs and a few tens of chillers, it is time-consuming to train all agents centrally. To achieve an accelerated learning process, transfer learning (i.e., transferring the knowledge from one task to a related but different task [102] [103]) is adopted. As shown in Fig. 5, the optimal policies obtained by training multiple agents on a sub-set of HVAC systems (including one AHU and one chiller) can be used to pre-train multiple agents related to other HVAC subsystems due to the problem similarity.

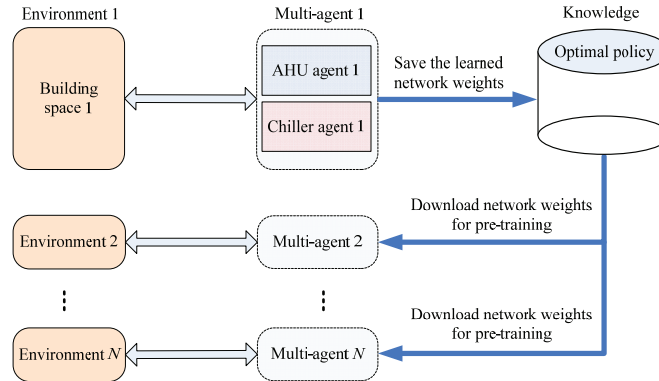


Fig. 5. The proposed transfer learning framework for multi-agent training

#### B. Model-based DRL Methods for HVAC Control

Although the above-mentioned works are effective, there are two drawbacks in the process of training a DRL agent. Firstly, it is impractical to let the DRL agent to explore the state space fully in a real building environment since unacceptably high cost may be incurred [59] [90] [91]. Secondly, it may take a long time for the DRL agent to learn an optimal policy if trained in a real-world environment [90] [91].

To reduce the number of interactions with a real building environment, many model-based DRL control methods have been developed [90] [92]. For example, Zhang *et al.* [92] proposed and implemented a building energy model (BEM)-based DRL control framework for a novel radiant heating system in an existing office building. The proposed framework consists of four steps as shown in Fig. 6, i.e., building energy modeling, model calibration, DRL training, and real deployment. To be specific, EnergyPlus software is used to develop a building energy model for the office building. Next, based on the observed data,



the building energy model can be calibrated. Then, the calibrated model is used as the simulator of environment for training the DRL agent off-line based on A3C algorithm. Finally, the learned control policy will be deployed in building automation system (BAS) for generating control signals in real-time. Experimental results showed that the obtained control strategy can reduce heating demand by 16.7% compared with the rule-based control strategy.

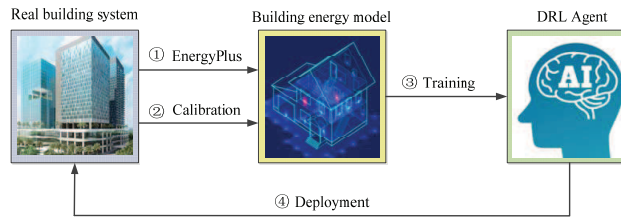


Fig. 6. BEM-based DRL control framework

In [92], the real-world HVAC operational data in three months are used for calibrating building energy models, which will affect model accuracy. To overcome this drawback, Zou *et al.* [90] proposed a DRL-based HVAC control framework to minimize energy consumption while maintaining thermal comfort levels for occupants based on operational data within two years. The proposed framework is composed of two parts as shown in Fig. 7, i.e., creating DRL training environment and training DRL agent based on the created environment. To be specific, LSTM models are built based on BAS historical data, which can approximate HVAC operations. Note that the inputs of LSTM models are current state and action, while their outputs are next state and reward. After LSTM networks are trained, they can be used to create training environment. Next, DRL agent interacts with the training environment until it converges to an optimal HVAC control policy. Finally, the optimal control policy can be deployed for controlling AHUs in real-time. Moreover, DRL agent contains an actor network and a critic network, which are trained using DDPG algorithm. Algorithmic testing results showed that DRL agents can save energy by 27% to 30% while maintaining the predicted percentage of discomfort at 10%.

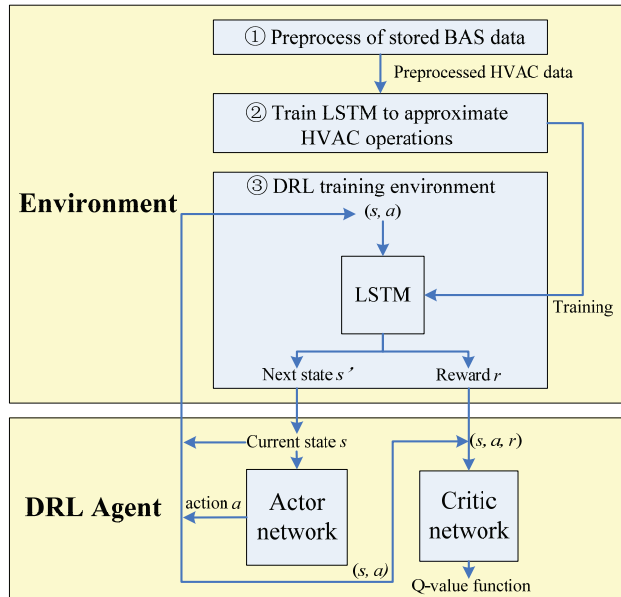


Fig. 7. LSTM-based DRL control framework

Similar to [90], Chen *et al.* proposed a PPO-based approach for HVAC control by utilizing historical data so that practical deployment can be achieved [91]. To be specific, the framework of the proposed approach is shown in Fig. 8. Firstly, historical data from existing HVAC controllers are used to pre-train a differentiable MPC policy based on imitation learning. Note that the pre-trained policy can encode domain knowledge into planning and system dynamics, making it both sample-efficient and interpretable. Secondly, the pre-trained control policy is improved continually in the process of interacting with the real building environment using online learning algorithm. Since PPO is robust to hyperparameters and network architectures, it is adopted to improve the pre-trained policy. Practical experimental results showed that the proposed approach can save 16.7% of cooling demand compared with the existing controller and track temperature set-point better.

TABLE V  
SUMMARY OF EXISTING WORKS ON DRL FOR A SINGLE BUILDING ENERGY SUBSYSTEM

Research work	Object(s)	Subsystem	Primary objective	Secondary objective(s)	DRL algorithm, function estimator	Performance improvement	Practical implementation
Morinibu <i>et al.</i> [97]	Smart home	HVAC	Non-uniformity of radiant temperature	Thermal comfort	A2C, CNN&LSTM	—	No
Wei <i>et al.</i> [19]	Office	HVAC	Energy cost	Thermal comfort	DQN, DNN	19.1%~71.2%	No
Nagy <i>et al.</i> [51]	Residential buildings	HVAC	Energy cost	Thermal comfort	D-DNFQI, DNN	5.5%~10%	No
Gupta <i>et al.</i> [98]	Residential buildings	HVAC	Energy cost	Thermal comfort	DQN, DNN	5%~12%	No
Yoon <i>et al.</i> [99]	Office	HVAC	Energy consumption	Thermal comfort	DQN, DNN	12.4%~32.2%	No
Sakuma <i>et al.</i> [100]	Residential buildings	HVAC	Energy consumption	Thermal comfort	DQN, DNN	34.5%	No
Gao <i>et al.</i> [18]	Laboratory	HVAC	Energy cost	Thermal comfort	DDPG, DNN	4.31%~9.15%	No
Valladares <i>et al.</i> [101]	Laboratory and classroom	HVAC	Energy cost	Thermal comfort, air quality	DDQN, DNN	4%~5%	No
Nagarathinam <i>et al.</i> [49]	A campus building	HVAC	Energy consumption	Thermal comfort	Multi-agent DDQN, DNN	17%	No
Zhang <i>et al.</i> [92]	Office	HVAC	Energy cost	Thermal comfort	BEM-A3C, DNN	7.06%~16.7%	Yes
Zou <i>et al.</i> [90]	Office	HVAC	Energy consumption	Thermal comfort	LSTM-DDPG, LSTM	27%~31.27%	No
Chen <i>et al.</i> [91]	Office	HVAC	Energy cost	Thermal comfort	Differentiable MPC-PPO, Linear model	16.7%	Yes
Kazmi <i>et al.</i> [94]	Residential buildings	EWH	Energy consumption	Thermal comfort, exploration bonus	Deep PILCO, DNN	20%	Yes
Ruelens <i>et al.</i> [95]	Residential buildings	EWH	Energy cost	Thermal comfort	Fitted Q-iteration, CNN/LSTM	5.5%~10.2%	No
Peirelinck <i>et al.</i> [96]	Residential buildings	EWH	Energy cost	Thermal comfort	DDQN, DNN	8.8%~32.2%	No

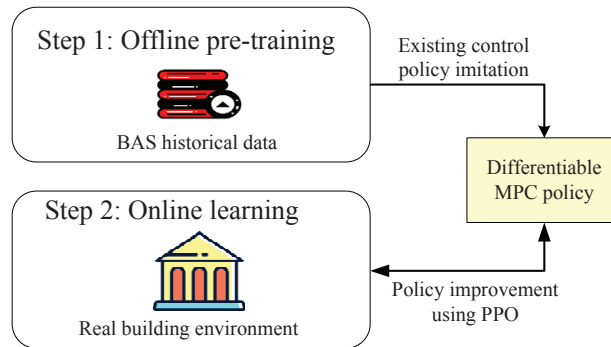


Fig. 8. Differentiable MPC policy based HVAC control framework

### C. DRL Methods for EWH Control

In an EWH, there are many separate layers of water and each layer has a unique temperature in practice, measuring the temperature within the EWH using a single sensor will lead to sparse observations. In [95], Ruelens *et al.* proposed an effective method to tackle sparsely observed control problem related to EWHs based on fitted Q-iteration and LSTM. The key idea of the proposed method is to store sequences of past observations and actions in the state vector so that relevant features for finding near-optimal control policies can be extracted based on RL. Simulation results showed that LSTM has better performance than CNN and DNN when they are used as function estimators in RL.

Since training a DRL agent without knowing system dynamics of EWHs requires a large number of interactions with the actual environment, model-based DRL methods are preferred due to their high sample efficiency. For example, Kazmi *et al.* proposed a model-based DRL method to optimize the hot water production based on Deep PILCO, which can reduce energy consumption by about 20% and has been applied to a set of 32 houses in the Netherlands. The key idea of the proposed

method is summarized as follows. Firstly, executing actions under the current policy and collecting experience transition data for training system dynamics model. Next, generating trajectories based on the current policy and the obtained dynamics model. Then, trajectories are used to update policy. Finally, the updated policy will be used in next loop.

When prior knowledge about system dynamics is available, learning a model from observations could be avoided. In [96], Peirelinck *et al.* investigated an energy cost minimization problem related to EWHs given a known system dynamics model. Since the model expression is very complex, it is challenging to find the optimal policy for EWH operation. Therefore, DRL is used to obtain the optimal policy since it merely cares about the input and output of the model. To reduce the training time in target domain (i.e., practical environment), domain randomization is used as shown in Fig. 9. To be specific, model parameters are randomized in the source domain (i.e., training environment). Based on the model with uncertain parameters, a generalized policy in the source domain is trained. Then, the trained policy is transferred to target domain for initializing DNNs. Simulation results showed that pre-training is helpful for reducing energy cost by 8.8%.

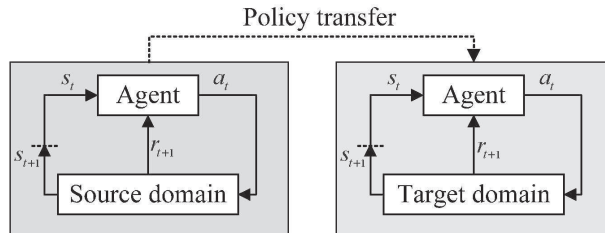


Fig. 9. Domain randomization based pre-training method

**Summary:** In this section, we review existing works on DRL for a single building energy subsystem. For easy reading, the specific details including objectives, DRL algorithms, and implementation methods are summarized in Table V. It can be observed that most existing works focus on HVAC control since HVAC systems have the largest energy consumption among all single building energy subsystems [10]. Moreover, most of optimization objectives are related to energy cost/consumption and thermal comfort. By controlling an HVAC system intelligently based on DRL methods, its energy cost can be reduced by 4%-71.2% and energy consumption can be decreased by 12.4%-34.5% without sacrificing thermal comfort. In addition, nearly all model-free DRL methods are evaluated by simulations and several model-based DRL methods have been deployed in practice.

## V. APPLICATIONS OF DRL IN MULTIPLE ENERGY SUBSYSTEMS OF BUILDINGS

In this section, we will introduce the applications of DRL in multiple energy subsystems of residential buildings and commercial buildings, respectively. To be specific, section V-A focuses on the coordination of home energy management system, HVAC systems, ESSs, EVs, WMs, PVs, and EWHs in residential buildings, while section V-B focuses on the coordination of HVAC systems, lighting systems, blind systems, window systems, and personal electric devices in commercial buildings. Moreover, we give a summary of existing works and provide some insights in the last paragraph of this section.

### A. Multiple Energy Subsystems in Residential Buildings

As the smallest unit in a residential building, smart home has many kinds of appliances, e.g., HVAC systems, EVs, ESSs, and PVs. To implement the coordination of different appliances, many DRL-based methods have been proposed to save energy cost. For example, Yu *et al.* proposed a DDPG-based home energy management algorithm to minimize energy cost for the joint scheduling of HVAC systems and ESSs [20]. Simulation results showed that the proposed algorithm can reduce energy cost by 8.1%-15.21% through the utilization of temporal diversity of dynamic prices [104] [105]. Similar works can be found in [106] and [107]. To be specific, DDQN and TRPO based methods have been proposed to minimize energy cost of a smart home with the consideration of occupant satisfaction degree or thermal comfort, respectively.

As for residential buildings, other objectives may be pursued when optimizing its energy use, e.g., peak demand [23], transformer capacity violation [22], and revenue of excess renewable energy [108]. For example, Mocanu *et al.* proposed two algorithms to minimize energy cost and peak load of residential buildings with the consideration of HVAC systems, EVs, and dishwashers (DWs) [23]. Simulation results illustrated that the proposed algorithms based on DQN and deep policy gradient can efficiently cope with the inherent uncertainty and variability in renewable energy generation and power demand. Although the proposed algorithms in [23] are effective, they neglect the physical constraints related to residential buildings, e.g., transformer capacity. To deal with this limitation, Zhang *et al.* investigated a multi-household energy management problem for residential units connected to the same transformer. Since violating the transformer capacity is harmful to its lifetime, an efficient approach was designed based on cooperative multi-agent DRL [22]. Simulation results indicated that the energy cost of residential households can be reduced by 59.77% without violating the transformer capacity.

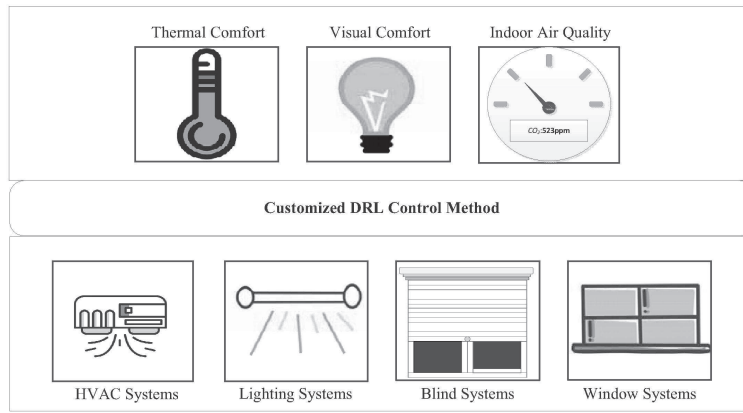


Fig. 10. The architecture of the proposed control framework

**B. Multiple Energy Subsystems in Commercial Buildings**

In existing works, some DRL-based approaches have been proposed to reduce energy consumption in commercial buildings [90] [92] [101]. Although some advances have been made, these works only consider a single subsystem in buildings (e.g., an HVAC system) without noticing that other subsystems can also affect energy consumption and user comfort in terms of thermal, air quality, and illumination conditions. In fact, some research results showed that jointly controlling HVAC systems and other building energy subsystems (e.g., blind systems, lighting systems, and window systems) has great potential of saving energy [109] [110]. For example, HVAC energy consumption can be reduced by 17%-47% if window-based natural ventilation is adopted [110]. Based on the above observation, Ding *et al.* proposed a DRL-based framework as shown in Fig. 10 for efficiently controlling four building energy subsystems (including HVAC systems, lighting systems, blind systems, and window systems [111]) so that the total energy consumed by all subsystems can be minimized while still maintaining user comfort. To solve the high-dimensional action problem, a branching dueling Q-Network (BDQ) algorithm was used. Moreover, a calibrated EnergyPlus simulation model was adopted to generate enough data for the training of the DRL agent. Simulation results showed that the proposed framework can save energy by 14.26% compared with the rule-based method while maintaining human comfort within a desired range.

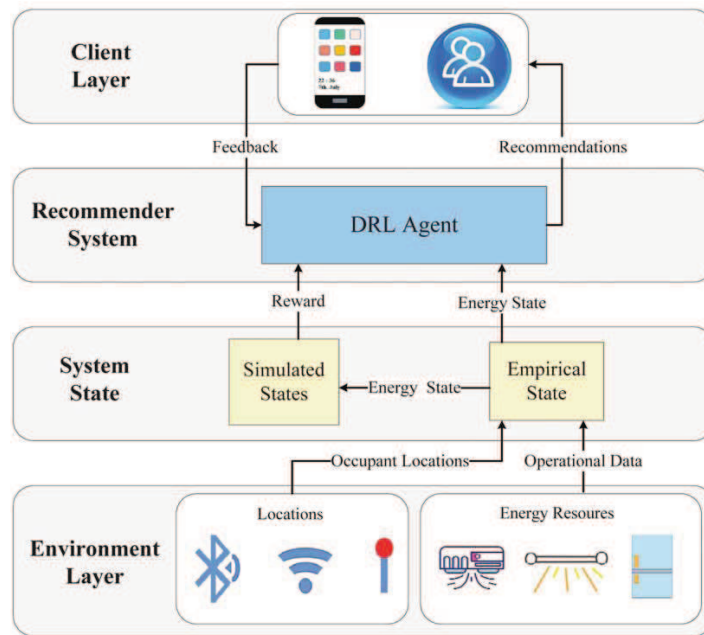


Fig. 11. The system architecture of the designed recommender

However, the above-mentioned works mainly focus on building energy system itself and treat occupants as immovable objects, which may decrease the potential of reducing energy consumption and be illustrated by the following example. Suppose that a space is sparsely occupied by some occupants. If they are recommended to vacate this space and move to another occupied

space, their comforts may be not sacrificed while the energy consumption in the vacated space can be reduced. Therefore, it is very necessary to investigate the potential of saving energy by shaping occupant behavior. To this end, Wei *et al.* [112] designed a DRL-based recommender system in commercial buildings, which can learn actions with high energy saving potential and distribute recommendations to occupants. Based on the feedback from occupants, better recommendations can be learned. The system architecture of the designed recommender is shown in Fig. 11, which consists of four layers, i.e., *environment layer*, *system state layer*, *recommender system*, and *client layer*. To be specific, *environment layer* measures building environment (e.g., occupant locations and energy consumption information) and sends such information to *system state layer*. *System state layer* contains two components, i.e., an empirical state, which maintains the current building state, and simulated states, which are used to represent the next state after the potential energy saving actions are taken. *The recommender system layer* learns the potential of different recommendation actions (including *move recommendation*, *schedule change*, *reduce personal resources*, and *reduce service in spaces*). *The client layer* receives recommendations and allows clients to provide feedback (e.g., accept or reject the recommendation). A four-week user study showed that the designed recommender system can reduce building energy consumption by 19% to 26% compared with a passive-only strategy.

TABLE VI  
SUMMARY OF EXISTING WORKS ON DRL FOR MULTI-ENERGY SUBSYSTEMS IN BUILDINGS

Research work	Object(s)	Energy subsystems	Primary objective	Secondary objective(s)	DRL algorithm, function estimator	Performance improvement	Practical implementation
Yu <i>et al.</i> [20]	Smart home	PV, ESS, HVAC	Energy cost	Thermal comfort	DDPG, DNN	8.10%~15.21%	No
Liu <i>et al.</i> [106]	Smart home	PV, ESS, HVAC, EV, Heater, DW	Energy cost	Consumers' satisfaction degree	DDQN, DNN	41.8%~59%	No
Li <i>et al.</i> [107]	Smart home	HVAC, EV, EWH, DW, WM	Energy cost	Thermal comfort and range anxiety	TRPO, DNN	31.6%	No
Mocanu <i>et al.</i> [23]	Residential buildings	PV, HVAC, EV, DW	Energy cost	Peak demand, load operational time or condition	DQN, DNN	14.1%~27.4%	No
Zhang <i>et al.</i> [22]	Residential buildings	PV, ESS, EV	Energy cost	Transformer capacity violation	PPO, DNN	59.77%	No
Ye <i>et al.</i> [108]	Residential buildings	PV, ESS, TES, EHP, GB	Energy cost	Excess energy sale revenue	PDDPG, DNN	6.28%~10.21%	No
Ding <i>et al.</i> [111]	Commercial buildings	HVAC, lighting, blind and window	Energy consumption	Thermal comfort, IAQ, lighting comfort	BDQ, DNN	14.26%	No
Wei <i>et al.</i> [112]	Commercial buildings	HVAC, lighting, plug load	Energy consumption	Safety, comfort, productivity	DQN, DNN	19%~26%	Yes

**Summary:** In this section, we review existing works on DRL applications in multiple energy subsystems of buildings. For easy understanding, the research objects, considered energy subsystems, research objectives, DRL algorithms, performance improvement, and implementation methods in existing works are summarized in Table VI. It can be observed that there is a great potential in reducing energy cost of buildings by scheduling multiple energy subsystems coordinately, e.g., relative energy cost reduction is up to 59% while maintaining comfort of occupants. Compared with the optimal HVAC control in Table V, more advanced DRL algorithms are adopted to deal with more complex problems, e.g., PDDPG, BDQ, and TRPO. In addition, most of DRL methods are evaluated by simulations.

## VI. APPLICATIONS OF DRL IN BUILDING MICROGRIDS

In this section, we review the existing works on DRL-based energy optimization for building microgrids. To be specific, section VI-A introduces DRL-based energy management algorithms for microgrids with uncontrollable building loads, while section VI-B introduces DRL-based microgrid optimization algorithms considering the flexibility of building loads. Moreover, we summarize the existing works and point out some insights at the end of this section.

### A. Microgrid Optimization without Considering Controllable Building Loads

In existing works, many DRL-based methods have been proposed for residential microgrids [113]–[117], where a microgrid is a low voltage distribution network comprising various distributed generation, storage devices, and responsive loads [118]. For example, Francois-Lavet *et al.* proposed a DQN-based control algorithm for a residential microgrid with the consideration of battery and hydrogen storage device to minimize the leveled energy cost [114]. Similar work can be found in [116]. However, power demand is assumed to be satisfied in [114]. For an isolated residential microgrid, load shedding may happen when the



total power supply is smaller than the total power demand. At this time, non-served power demand should be penalized. Based on this observation, Dominguez-Barbero *et al.* proposed a DQN-based microgrid optimization algorithm to minimize the sum of DG generation cost and the penalty of non-served power demand [115]. Different from the above works, Chen *et al.* investigated a peer-to-peer energy trading problem among multiple microgrids [117]. Moreover, a DQN-based energy trading strategy was proposed to maximize the utility function in a microgrid, which is related to trading profit, retail profit, battery wear cost, demand penalty, and virtual penalty. Simulation results based on one-year real generation and demand data showed the effectiveness of the proposed strategy. Although some advances have been made in above efforts, the proposed DQN-based methods can not deal with DRL problems with continuous actions (e.g., the generation output of DGs [43]).

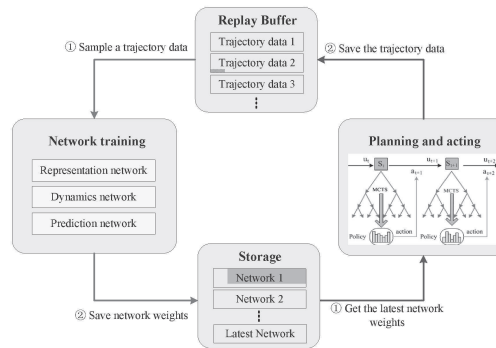


Fig. 12. The training process of the network model

To support continuous actions, DDPG-based methods could be adopted. For example, Lei *et al.* proposed a FH-DDPG based energy management algorithm for an isolated microgrid to minimize the sum of power generation cost and the power unbalance penalty [43]. Since model-free based DRL algorithms in existing works have low sample efficiency, Shuai *et al.* proposed a model-based DRL algorithm (i.e., *MuZero*) for the online scheduling of a residential microgrid under uncertainties [113] based on Monte-Carlo tree search (MCTS) strategy with a learned network model. Note that the off-line learning process of the network model can be depicted by Fig. 12, where four components can be identified, i.e., *network training*, *replay buffer*, *storage*, *planning and acting*. Firstly, the latest network weights are obtained from a storage and used for planning implemented by MCTS. Next, an action is sampled from the search policy, which is proportional to the visit count for each action from the root node. Then, the environment returns a new state and a reward. At the end of the episode, the trajectory data is stored into a replay buffer. When conducting network training, a trajectory data will be randomly sampled from the replay buffer and the updated network weights will be saved in a storage device. It is obvious that network training and trajectory data generation are two independent processes, which can be implemented in parallel. Once the training process of the network model (including three components, i.e., *representation*, *dynamics*, and *prediction*) is completed, the learned network model can be used as the simulator of the real environment. Based on MCTS and the network model, the optimal policy can be learned. Note that the proposed algorithm can operate without relying on any forecasting information and statistic distribution information of the system.

### B. Microgrid Optimization with Controllable Building Loads

In above-mentioned works, building loads are regarded as uncontrollable resources in microgrids. In fact, the energy cost of a microgrid could be reduced by scheduling loads flexibly. For example, Yang *et al.* proposed a DDPG-based scheduling algorithm for a data center microgrid with renewable sources to reduce energy cost by choosing the execution time and the quantity of served workloads flexibly [119]. Simulation results showed that energy cost can be reduced by 6.42%. However, the microgrid optimization problem would be intractable if the number of controllable resources is large due to the increased action space. At this time, multi-agent DRL may be a good choice, which can coordinate all agents effectively. For example, Yang *et al.* proposed an entropy-based collective multi-agent DRL algorithm to schedule EVs and ESSs in large-scale households. Simulation results based on real-world traces showed the effectiveness of the proposed algorithm in reducing the operating cost and the peak load [50]. Similarly, Lee *et al.* proposed an MAPPO-based algorithm to solve the demand response problem in a microgrid of residential district [120]. The proposed algorithm intends to train multiple household agents centrally. Once an optimal policy is learned by each household agent, it can schedule household appliances without knowing specific information about other households.

**Summary:** In this section, we review existing works on DRL applications in building microgrids and summarize the details of existing works in Table VII. It can be observed that existing works mainly focus on economic impacts of building microgrids and the proposed DRL-based methods can indeed bring economic benefits for microgrid operators. However, most of them neglect the control of building loads and all of them are not implemented in practice.

TABLE VII  
SUMMARY OF EXISTING WORKS ON DRL FOR MICROGRIDS

Research work	Microgrid type	Energy systems	Controllable building load considered	Optimization objective(s)	DRL algorithm, function estimator	Cost reduction	Practical implementation
Francois-Lavet <i>et al.</i> [114]	A residential microgrid	PV, Battery, hydrogen storage device	No	The overall levelized energy cost	DQN, DNN	5%~12%	No
Ji <i>et al.</i> [116]	A residential microgrid	PV, WT, DG, ESS	No	Daily operating cost	DQN, DNN	20.75%	No
Dominguez-Barbero <i>et al.</i> [115]	An isolated residential microgrid	PV, DG, Battery, hydrogen storage device	No	Operating cost	DQN, DNN	58.5%~67.20%	No
Chen <i>et al.</i> [117]	A residential microgrid	PV, ESS	No	Profit minus cost	DQN, DNN	>30%	No
Lei <i>et al.</i> [43]	An isolated microgrid	PV, DG, ESS	No	Power generation cost, power unbalance	FH-DDPG/RDPG, LSTM	80%	No
Shuai <i>et al.</i> [113]	A residential microgrid	PV, WT, ESS	No	Operating cost	MuZero, DNN&LSTM	9.28%~28.93%	No
Yang <i>et al.</i> [119]	A data center microgrid	PV, ESS, servers	Yes, servers	Energy cost	DDPG, DNN	6.24%	No
Yang <i>et al.</i> [50]	A residential microgrid	PV, ESS, EV	Yes, EV	Energy cost and peak load	Multi-agent EB-C-A2C/DQN, DNN	24.69%	No
Lee <i>et al.</i> [120]	A residential microgrid	WM, CD, WH, DW and refrigerator.	Yes, household appliances	Energy cost and peak load	Multi-agent PPO, DNN	—	No

## VII. OPEN ISSUES AND FUTURE RESEARCH DIRECTIONS

Although recent years have witnessed the rapid development of DRL for SBEM, there are still some unsolved issues that need better solutions. In this section, we highlight open issues and point out future research directions.

### A. Data-efficient building energy optimization

As mentioned in Sections IV-VI, most model-free DRL methods for SBEM are still not be implemented in practice. The main reason for this phenomenon is that DRL agents have to interact with the building environment directly so as to collect enough data for training, which is a time-consuming process. Moreover, in the process of interaction, actions are taken by trial and error, resulting in a high exploration cost. For example, random selection of an HVAC temperature set-point may lead to thermal discomfort and high energy consumption. When mitigating these issues, there are several opportunities. To be specific, with the development of IoT technologies, many sensing equipments can be deployed to collect building operational data. Then, the collected data could be used to train DRL agents in an offline way. Moreover, deep meta reinforcement learning [121] could be used to implement fast learning using only a few data and training episodes. In addition, the collected operational data can be used to learn an environment model and consequently model-based DRL methods can be used to reduce exploration cost.

### B. Multi-timescale building energy optimization

Most existing DRL-based methods focus on single-timescale building energy optimization problems. In fact, there are many multi-timescale decision problems in the field of building energy optimization. For example, supply air temperature and the ratio of re-use air in a commercial building HVAC system can be adjusted once every hour since the frequent adjustment can cause damage to HVAC components [122]. In contrast, supply air rate in each zone can be changed every 10-15 minutes [123]. When confronted with multi-timescale decision problems, existing DRL-based methods are not applicable. A possible way is to design energy optimization algorithms based on the framework of hierarchical DRL [124], which can support multi-timescale DRL problems with delayed rewards. In hierarchical DRL, actions can be divided into two types with different timescales. To be specific, actions with long timescale are first taken in the upper level based on system state. Then, actions with short timescale are taken in the lower level based on system state and the chosen actions in the upper level. By coordinating the actions of upper level and lower level, hierarchical DRL-based methods can explore the environments efficiently.

### C. Multi-objective building energy optimization

As shown in Section III, multiple objectives are pursued by SBEM, e.g., energy cost/consumption minimization, carbon emission minimization, and comfort maximization. Moreover, such objectives are often conflicting with each other. A typical

way of dealing with conflicting objectives in existing DRL-based methods is to design a synthetic reward function as a weighted sum of different objectives. Since the weight parameters related to different objectives typically have different units and/or scales, it is very challenging to decide their proper values beforehand. Moreover, the learned policies based on the above-mentioned way can not support flexible operation of building energy systems, e.g., switching flexibly between low-energy-cost mode and high-comfort mode. To avoid deciding weighted parameters for multiple objectives and support flexible operations, a possible way is to design building energy optimization algorithms based on some advanced DRL frameworks (e.g., multi-objective DRL [125], and multi-objective meta-DRL [126]).

#### D. Multi-zone building energy optimization

In existing works on building HVAC systems, the proposed DRL-based control methods mainly focus on a single-zone building. In [19], Wei *et al.* proposed a heuristic algorithm for variable air volume (VAV) HVAC control in a multi-zone office building and DRL agent for each zone was trained separately. Although the proposed algorithm was effective when 5 zones were considered, it was not scalable due to the lack of multi-zone coordination. In [47], Hu *et al.* proposed a MADDPG-based method to decide temperature and humidity setpoints in a four-zone building. Since the input of each critic in MADDPG is the concatenation of state and action information from all agents, the scalability of the MADDPG-based method was not very high. In [21], Yu *et al.* proposed an MAAC-based VAV HVAC control method for a multi-zone commercial building with the consideration of thermal comfort, indoor air quality comfort, and random occupancy, which can operate effectively when 30 zones were considered. Although the above methods are effective when the number of zones is not large, more scalable multi-agent DRL algorithms are expected since the number of zones in a practical commercial building may exceed one hundred or even larger.

#### E. Efficient training of DRL agents in multi-building energy optimization

As introduced in Section III, model-based DRL methods for building energy optimization are sample-efficient. However, a large amount of historical data should be required when learning building thermal dynamics models. For some buildings, especially brand-new buildings, historical data are very limited. At this time, how to speed up the training of building DRL agents is a very challenging task. To improve this situation, a possible way is to combine DRL with transfer learning [127]. In the field of building energy management, the transferred knowledge may be building thermal dynamics models [128] or control strategies [129]. Although some efforts have been made in existing works, they mainly focus on transfer learning problems with simple scenarios, where a small similarity gap exists between source MDP and target MDP related to DRL-based SBEM. When the similarity gap is large (e.g., the dimensions of state spaces and action spaces in two MDPs are different), how to design efficient inter-task mapping function and select proper form of the transferred knowledge is very challenging, especially for multi-agent DRL-based SBEM problems.

#### F. DRL-based energy optimization for building microgrids

Due to the high thermal inertia, buildings can be regarded as thermal energy storage units. By incorporating building thermal dynamics into microgrid scheduling [118] or planning [130], the operation cost or total annualized cost can be reduced. However, explicit building thermal dynamics models are required in the above works. Although DRL-based methods can operate without knowing them, several challenges have to be addressed. Firstly, both discrete and continuous decision variables (e.g., discrete variables are used for describing the operational states of WMs, HVAC loads, and distributed generators, while continuous variables are used for describing the EV charging/discharging power) exist in the optimal operation problem related to building microgrids, which means that discrete-continuous hybrid actions should be supported by the designed DRL-based algorithms. Secondly, multi-agent DRL energy management algorithms with complex reward components should be designed to efficiently promote the coordination among the microgrid controller and all building energy management systems, since each building has their respective objectives (e.g., comfort requirements) and also needs to participate in optimizing the objective of the microgrid.

**Remarks:** The above-mentioned DRL techniques for SBEM can be supported by existing IEC energy management standards (e.g., ISO/IEC 15067-3-3-2019). For example, ISO/IEC 15067-3-3-2019 defines some energy management agents and provides their operational modes, e.g., single-agent mode, mesh mode, hierarchical mode, and mixed hierarchical and mesh mode. Correspondingly, algorithms based on single-agent DRL, multi-agent DRL, hierarchical DRL, and multi-agent hierarchical DRL can be adopted and implemented by energy management agents.

### VIII. CONCLUSIONS AND LESSONS LEARNED

In this paper, we reviewed the DRL applications in SBEM with the consideration of different system scales comprehensively. In particular, we summarized the features of different DRL methods for SBEM. Moreover, we provided some insights, identified some unsolved issues, and pointed out potential directions for future research. A few major lessons that we learned from this review are summarized as follows. Firstly, nearly all model-free DRL-based building energy optimization methods are still not

implemented in practice due to a long exploration time and a high exploration cost. Secondly, model-based DRL approaches for building energy optimization are more practical than model-free DRL approaches since the former can generate enough training data for DRL agents and reduce the number of interactions with the real environment. When the amount of historical data is not enough in the current environment, transfer learning can be used to pre-train a building thermal dynamics model or policy based on the large amount of historical data in a related, but different building environment. Thirdly, compared with some traditional methods, DRL-based energy management methods have the potential of improving some building performance metrics (e.g., energy cost, peak load, and occupant dissatisfaction degree) simultaneously. Finally, although some advances have been made in existing works, there are still many challenges caused by low data efficiency, multiple timescales, multiple optimization objectives, multiple zones, multiple buildings, and building microgrids.

## REFERENCES

- [1] S. Hu, C. Hore, P. Raftery, and J. O'Donnell, "Environmental and energy performance assessment of buildings using scenario modelling and fuzzy analytic network process," *Applied Energy*, vol. 255, pp. 113788-113799, 2019.
- [2] S. Hu, E. Corry, M. Horrigan, C. Hore, M. Dos Reis, and J. O'Donnell, "Building performance evaluation using OpenMath and Linked Data," *Energy and Buildings*, vol. 174, pp. 484-494, 2018.
- [3] J. Park, T. Dougherty, H. Fritz, and Z. Nagy, "Multi-agent deep reinforcement learning for zero energy communities," *Building and Environment*, vol. 147, pp. 397-414, 2019.
- [4] X. Dong, Y. Liu, Z. Xu, J. Wu, J. Liu, and X. Guan, "Optimal scheduling of distributed hydrogen-based multi-energy systems for building energy cost and carbon emission reduction," *Proc. of IEEE 16th International Conference on Automation Science and Engineering (CASE)*, 2020.
- [5] E. Haghi, Q. Kong, M. Fowler, K. Raahemifar, and M. Qadrdan, "Assessing the potential of surplus clean power in reducing GHG emissions in the building sector using game theory; a case study of Ontario, Canada," *IET Energy Systems Integration*, vol. 1, no. 3, pp. 184-193, 2019.
- [6] The Global Alliance for Buildings and Construction (GABC), "The Global Status Report 2020," [https://globalabc.org/sites/default/files/inline-files/2020%20Buildings%20GSR\\_FULL%20REPORT.pdf](https://globalabc.org/sites/default/files/inline-files/2020%20Buildings%20GSR_FULL%20REPORT.pdf)
- [7] S. Sharma, Y. Xu, A. Verma, and B. Panigrahi, "Time-coordinated multi-energy management of smart buildings under uncertainties," *IEEE Trans. Industrial Informatics*, vol. 15, no. 8, pp. 4788-4798, 2019.
- [8] N. Zhou, N. Khanna, W. Feng, J. Ke, and M. Levine, "Scenarios of energy efficiency and CO<sub>2</sub> emissions reduction potential in the buildings sector in China to year 2050," *Nature Energy*, vol. 3, no. 978-984, 2018.
- [9] B. Qolomany, A. Al-Fuqaha, A. Gupta, D. Benhaddou, S. Alwajidi, J. Qadir, and A. Fong, "Leveraging machine learning and big data for smart buildings: a comprehensive survey," *IEEE Access*, vol. 7, pp. 90316-90356, 2019.
- [10] D. Minoli, K. Sohraby, and B. Occhiogrosso, "IoT considerations, requirements, and architectures for smart buildings-energy optimization and next-generation building management systems," *IEEE Internet of Things Journal*, vol. 4, no. 1, pp. 269-283, Feb. 2017.
- [11] X. Zhang, M. Pipattanasomporn, T. Chen, and S. Rahman, "An IoT-based thermal model learning framework for smart buildings," *IEEE Internet of Things Journal*, vol. 4, no. 1, pp. 269-283, Feb. 2017.
- [12] W. Feng, Z. Wei, G. Sun, Y. Zhou, H. Zang, and S. Chen, "A conditional value-at-risk-based dispatch approach for the energy management of smart buildings with HVAC systems," *Electric Power Systems Research*, vol. 188, pp. 196535-106534, 2020.
- [13] B. Yang, X. Li, Y. Hou, A. Meier, X. Cheng, J. Choi, F. Wang, H. Wang, A. Wagner, D. Yan, A. Li, T. Olofsson, and H. Li, "Non-invasive (non-contact) measurements of human thermal physiology signals and thermal comfort/discomfort poses-A review," *Energy and Buildings*, vol. 224, pp. 110261-110270, 2020.
- [14] F. Wang, L. Zhou, H. Ren, X. Liu, S. Talari, and M. Shafie-khah, "Multi-objective optimization model of source-load-storage synergetic dispatch for a building energy management system based on TOU price demand response," *IEEE Trans. Industry Applications*, vol. 54, no. 2, pp. 1017-1028, 2018.
- [15] A. Pallante, L. Adacher, M. Botticelli, S. Pizzuti, G. Comodi, and A. Monteriu, "Decision support methodologies and day-ahead optimization for smart building energy management in a dynamic pricing scenario," *Energy and Buildings*, vol. 216, pp. 109963-109973, 2020.
- [16] A. Ahmad and J. Khan, "Real-time load scheduling, energy storage control and comfort management for grid-connected solar integrated smart buildings," *Applied Energy*, vol. 259, pp. 114208-114226, 2020.
- [17] R. Zhang, T. Jiang, G. Li, X. Li, and H. Chen, "Stochastic optimal energy management and pricing for load serving entity with aggregated TCLs of smart buildings: A stackelberg game approach," *IEEE Trans. Industrial Informatics*, DOI: 10.1109/TII.2020.2993112, 2020.
- [18] G. Gao, J. Li, and Y. Wen, "DeepComfort: Energy-efficient thermal comfort control in smart buildings via deep reinforcement learning," *IEEE Internet of Things Journal*, vol. 7, no. 9, pp. 8472-8484, 2020.
- [19] T. Wei, Y. Wang, and Q. Zhu, "Deep reinforcement learning for building HVAC control," *Proc. of DAC*, 2017.
- [20] L. Yu, W. Xie, D. Xie, Y. Zou, D. Zhang, Z. Sun, L. Zhang, Y. Zhang, and T. Jiang, "Deep reinforcement learning for smart home energy management," *IEEE Internet of Things Journal*, vol. 7, no. 4, pp. 2751-2762, 2020.
- [21] L. Yu, Y. Sun, Z. Xu, C. Shen, D. Yue, T. Jiang, and X. Guan, "Multi-agent deep reinforcement learning for HVAC control in commercial buildings," *IEEE Transactions on Smart Grid*, vol. 12, no. 1, pp. 407-419, 2021.
- [22] C. Zhang, S. Kuppannagari, C. Xiong, R. Kannan, and V. Prasanna, "A cooperative multi-agent deep reinforcement learning framework for real-time residential load scheduling," *Proc. of the ACM/IEEE Conference on Internet of Things Design and Implementation*, 2019.
- [23] E. Mocanu, D. Mocanu, P. Nguyen, A. Liotta, M. Webber, M. Gibescu, and J. Sloopweg, "On-line building energy optimization using deep reinforcement learning," *IEEE Trans. Smart Grid*, vol. 10, no. 4, pp. 3698-3708, 2019.
- [24] Z. Zhang, D. Zhang, and R. Qiu, "Deep reinforcement learning for power system: an overview," *CSEE Journal of Power and Energy Systems*, vol. 6, no. 1, pp. 213-225, 2020.
- [25] Y. Ma, J. Matuško, and F. Borrelli, "Stochastic model predictive control for building HVAC systems: complexity and conservatism," *IEEE Trans. Control Systems Technology*, vol. 23, no. 1, pp. 101-116, 2015.
- [26] X. Guan, Z. Xu, and Q. Jia, "Energy-efficient buildings facilitated by microgrid," *IEEE Trans. Smart Grid*, vol. 1, no. 3, pp. 243-252, Dec. 2010.
- [27] L. Yu, D. Xie, T. Jiang, Y. Zou, and K. Wang, "Distributed real-time HVAC control for cost-efficient commercial buildings under smart grid environment," *IEEE Internet of Things Journal*, vol. 5, no. 1, pp. 44-55, 2018.
- [28] V. Mnih, et al. "Human-level control through deep reinforcement learning," *Nature*, vol. 518, pp. 529-541, 2015.
- [29] Y. Li, "Deep reinforcement learning: an overview," <http://arXiv:1701.07274v5>, 2017.
- [30] J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, et al., "Mastering Atari, Go, Chess and Shogi by planning with a learned model," <http://arXiv:1911.08265v2>, 2020.
- [31] K. Shao, Z. Tang, Y. Zhu, N. Li, and D. Zhao, "A survey of deep reinforcement learning in video games," <https://arxiv.org/pdf/1912.10944.pdf>, 2019.
- [32] J. Chen, B. Yuan and M. Tomizuka, "Model-free deep reinforcement learning for urban autonomous driving," <https://arxiv.org/abs/1904.09503v2>, 2019.
- [33] B. Kiran, I. Sobh, V. Talpaert, P. Mannion, et al., "Deep reinforcement learning for autonomous driving: A survey," *IEEE Transactions on Intelligent Transportation Systems*, DOI:10.1109/TITS.2021.3054625, 2021.



- [34] S. Aradi, "Survey of deep reinforcement learning for motion planning of autonomous vehicles," *IEEE Transactions on Intelligent Transportation Systems*, DOI:10.1109/TITS.2020.3024655, 2021.
- [35] A. Haydari, and Y. Yilmaz, "Deep reinforcement learning for intelligent transportation systems: A survey," *IEEE Transactions on Intelligent Transportation Systems*, DOI:10.1109/TITS.2020.3008612, 2021.
- [36] L. Lei, Y. Tan, K. Zheng, S. Liu, K. Zhang, X. Shen, "Deep reinforcement learning for autonomous Internet of things: model, applications and challenges," *IEEE Communications Surveys & Tutorials*, DOI:10.1109/COMST.2020.2988367, 2019.
- [37] S. Lee and D. Choi, "Energy management of smart home with home appliances, energy storage system and electric vehicle: a hierarchical deep reinforcement learning approach," *Sensors*, vol. 20, pp. 2157-2178, 2019.
- [38] X. Zhang, D. Biagioni, M. Cai, P. Graf, and S. Rahman, "An edge-cloud integrated solution for buildings demand response using reinforcement learning," *IEEE Trans. Smart Grid*, DOI: 10.1109/TH.2020.3007167, 2020.
- [39] M. Mohammadi, A. Al-Fuqaha, M. Guizani, and J. Oh, "Semisupervised deep reinforcement learning in support of IOT and smart city services," *IEEE Internet of Things Journal*, vol. 5, no. 2, pp. 624-635, 2018.
- [40] Y. Sun, M. Peng, and S. Mao, "Deep reinforcement learning-based mode selection and resource management for green fog radio access networks," *IEEE Internet of Things Journal*, vol. 6, no. 2, pp. 1960-1971, 2019.
- [41] L. Lin, X. Guan, Y. Peng, N. Wang, S. Maharjan, and T. Ohtsuki, "Deep reinforcement learning for economic dispatch of virtual power plant in Internet of energy," *IEEE Internet of Things Journal*, vol. 7, no. 7, pp. 6288-6301, 2020.
- [42] C. Wang, J. Wang, J. Wang, and X. Zhang, "Deep reinforcement learning-based autonomous UAV navigation with sparse rewards," *IEEE Internet of Things Journal*, vol. 7, no. 7, pp. 6180-6190, 2020.
- [43] L. Lei, Y. Tan, G. Dahlenburg, W. Xiang, and K. Zheng, "Dynamic energy dispatch based on deep reinforcement learning in IoT-driven smart isolated microgrids," *IEEE Internet of Things Journal*, DOI:10.1109/JIOT.2020.3042007, 2020.
- [44] X. Chen, H. Zhang, C. Wu, S. Mao, Y. Ji, M. Bennis, "Optimized computation offloading performance in virtual edge computing systems via deep reinforcement learning," *IEEE Internet of Things Journal*, vol. 6, no. 3, pp. 4005-4018, 2019.
- [45] R. Lu, Y. Li, Y. Li, J. Jiang and Y. Ding, "Multi-agent deep reinforcement learning based demand response for discrete manufacturing systems energy management," *Applied Energy*, vol. 276, pp. 115473-115483, 2020.
- [46] C. Lork, W. Li, Y. Qin, Y. Zhou, C. Yuen, W. Tushar and T.K. Saha, "An uncertainty-aware deep reinforcement learning framework for residential air conditioning energy management," *Applied Energy*, Vol. 276, pp. 115426-115438, 2020.
- [47] W. Hu, "Transforming thermal comfort model and control in the tropics: a machine-learning approach," Nanyang Technological University, 2020.
- [48] Y. Hu, W. Li, K. Xu, T. Zahid, F. Qin, and C. Li, "Energy management strategy for a hybrid electric vehicle based on deep reinforcement learning," *Applied Sciences*, vol. 8, no. 187, pp. 1-15, 2018.
- [49] S. Nagarathinam, V. Menon, A. Vasan, and A. Sivasubramaniam, "MARCO-Multi-agent reinforcement learning based control of building HVAC systems," *Proc. of the Eleventh ACM International Conference on Future Energy Systems*, 2020.
- [50] Y. Yang, J. Hao, Y. Zheng, and C. Yu, "Large-scale home energy management using entropy-based collective multiagent deep reinforcement learning framework," *Proc. of the Twenty-Eighth International Joint Conference on Artificial Intelligence*, 2019.
- [51] A. Nagy, H. Kazmi, F. Cheaib, and J. Driesen, "Deep reinforcement learning for optimal control of space heating," <https://arxiv.org/abs/1805.03777>, 2018.
- [52] N. Luong, D. Hoang, S. Gong, D. Niyato, P. Wang, Y. Liang, and D. Kim, "Applications of deep reinforcement learning in communications and networking: a survey," *IEEE Communications Surveys & Tutorials*, vol. 21, no. 4, pp. 3133-3174, 2019.
- [53] T. Nguyen and V. Reddi, "Deep reinforcement learning for cyber security," <http://arXiv:1906.05799v2>, 2020.
- [54] T. Nguyen, N. Nguyen, and S. Nahavandi, "Deep reinforcement learning for multi-agent systems: A review of challenges, solutions and applications," *IEEE Trans. Cybernetics*, vol. 50, no. 9, 3826-3839, 2020.
- [55] D. Zhang, X. Han, and C. Deng, "Review on the research and practice of deep learning and reinforcement learning in smart grids," *CSEE Journal of Power and Energy Systems*, vol. 4, no. 3, pp. 362-370, 2018.
- [56] T. Yang, L. Zhao, W. Li, and A. Zomaya, "Reinforcement learning in sustainable energy and electric systems: a survey," *Annual Reviews in Control*, vol. 49, pp. 145-163, 2020.
- [57] M. Han, R. May, X. Zhang, X. Wang, S. Pan, D. Yan, Y. Jin, and L. Xu, "A review of reinforcement learning methodologies for controlling occupant comfort in buildings" *Sustainable Cities and Society*, vol. 51, pp. 101748-101762, 2019.
- [58] J. Leitão, P. Gil, B. Ribeiro, and A. Cardoso, "A survey on home energy management," *IEEE Access*, vol. 8, pp. 5699-5722, 2020.
- [59] K. Mason and S. Grijalva, "A review of reinforcement learning for autonomous building energy management," *Computers & Electrical Engineering*, vol. 78, pp. 300-312, 2019.
- [60] Z. Wang and T. Hong, "Reinforcement learning for building controls: the opportunities and challenges," *Applied Energy*, vol. 269, pp. 115036-115056, 2020.
- [61] B. Rajasekhar, W. Tushar, C. Lork, Y. Zhou, C. Yuen, N. Pindoriya, and K. Wood, "A survey of computational intelligence techniques for air-conditioners energy management," *IEEE Trans. Emerging Topics in Computational Intelligence*, vol. 4, no. 4, pp. 555-570, 2020.
- [62] R.S. Sutton and A.G. Barto, "Reinforcement learning: an introduction," The MIT Press, London, England, 2018.
- [63] Z. Wang, T. Hong, and M. Piette, "Data fusion in predicting internal heat gains for office buildings through a deep learning approach," *Applied Energy*, vol. 240, pp. 386-398, 2019.
- [64] N. Raman, A. Devraj, P. Barooah, and S. Meyn, "Reinforcement learning for control of building HVAC systems," *Proc. of American Control Conference (ACC)*, 2020.
- [65] X. Deng, Y. Jiang, L. Yang, L. Yi, J. Chen, Y. Liu, and X. Li, "Learning automata based confident information coverage barriers for smart ocean internet of things," *IEEE Internet of Things Journal*, doi: 10.1109/JIOT.2020.2989696.
- [66] J. Vázquez-Canteli, and Z. Nagy, "Reinforcement learning for demand response: A review of algorithms and modeling techniques," *Applied Energy*, vol. 235, pp. 1072-1089, 2019.
- [67] S. Kim and H. Lim, "Reinforcement learning based energy management algorithm for smart energy buildings," *Energies*, vol. 11, pp. 1-19, 2018.
- [68] M. Ahrarinouri, M. Rastegar, and A. Seif, "Multiagent reinforcement learning for energy management in residential buildings," *IEEE Trans. Industrial Informatics*, vol. 17, no. 1, pp. 659-666, 2021.
- [69] R. Lu, S. Hong, and M. Yu, "Demand response for home energy management using reinforcement learning and artificial neural network," *IEEE Trans. Smart Grid*, DOI: 10.1109/TSG.2019.2909266, 2019.
- [70] F. Ruelens, B. Claessens, S. Vandael, B. Schutter, R. Babuška, and R. Belmans, "Residential demand response of thermostatically controlled loads using batch reinforcement learning," *IEEE Trans. Smart Grid*, vol. 8, no. 5, pp. 2149-2159, Sept. 2017.
- [71] S. Gu, T. Lillicrap, Z. Ghahramani, R.E. Turner, B. Schölkopf, and S. Levine, "Interpolated policy gradient: merging on-policy and off-policy gradient estimation for deep reinforcement learning," *Prof. of NIPS*, 2017.
- [72] H. Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double Q-learning," *Prof. of AAAI*, 2016.
- [73] Z. Wang, T. Schaul, M. Hessel, H. Hasselt, M. Lanctot, N. Freitas, "Dueling network architectures for deep reinforcement learning," *Prof. of ICML*, 2016.
- [74] T. Schaul, J. Quan, I. Antonoglou and D. Silver, "Prioritized experience replay," *Prof. of ICLR*, 2016.
- [75] M. Bellemare, W. Dabney, and R. Munos, "A Distributional Perspective on Reinforcement Learning," *Prof. of ICML*, 2017.
- [76] M. Fortunato, et. al., "Noisy networks for exploration," *Prof. of ICLR*, 2018.



- [77] M. Hessel, et. al., "Rainbow: combining improvements in deep reinforcement learning," *Prof. of AAAI*, 2018.
- [78] T. P. Lillicrap, et. al., "Continuous control with deep reinforcement learning," *Proc. of ICLR*, 2016.
- [79] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," *Proc. of NIPS*, 2017.
- [80] S. Iqbal and F. Sha, "Actor-attention-critic for multi-agent reinforcement learning," *Proc. of ICML*, 2019.
- [81] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor," *Proc. of ICML*, 2018.
- [82] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," *Proc. of ICML*, 2016.
- [83] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," *Proc. of ICML*, 2015.
- [84] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal Policy Optimization Algorithms," <https://arxiv.org/abs/1707.06347>, 2017.
- [85] C. Yu, A. Velu, E. Vinitzky, Y. Wang, A. Bayen, and Y. Wu, "The surprising effectiveness of MAPPO in cooperative, multi-agent games," *Proc. of NIPS*, 2017.
- [86] Y. Gal, R. McAllister and C. Rasmussen, "Improving PILCO with bayesian neural network dynamics models," *ICML Workshop on Data-Efficient Machine Learning*, 2016.
- [87] D. Ha and J. Schmidhuber, "World models," <https://arxiv.org/abs/1803.10122>, 2018.
- [88] J. Schrittwieser, et. al., "Mastering Atari, Go, Chess and Shogi by planning with a learned model," <http://arXiv:1911.08265v2>, 2020.
- [89] H. Chung, S. Maharjan, Y. Zhang, and F. Eliassen, "Distributed deep reinforcement learning for intelligent load scheduling in residential smart grids," *IEEE Trans. Industrial Informatics*, DOI: 10.1109/TII.2020.3007167, 2020.
- [90] Z. Zou, X. Yu, and S. Ergun, "Towards optimal control of air handling units using deep reinforcement learning and recurrent neural network," *Building and Environment*, <https://doi.org/10.1016/j.buildenv.2019.106535>, 2020.
- [91] B. Chen, Z. Cai, and M. Berges, "Gnu-RL: A precocial reinforcement learning solution for building HVAC control using a differentiable MPC policy," *Proc. of BuildSys*, 2019.
- [92] Z. Zhang, A. Chong, Y. Pan, C. Zhang, and K. Lam, "Whole building energy model for HVAC optimal control: A practical framework based on deep reinforcement learning", *Energy and Buildings*, vol. 199, pp. 472-490, 2019.
- [93] S. Latif, H. Cuayáhuitl, F. Pervez, F. Shamshad, H. Shehbaz Ali, and E. Cambria, "A survey on deep reinforcement learning for audio-based applications," <https://arxiv.org/pdf/2101.00240.pdf>, 2021.
- [94] H. Kazmi, F. Mehmood, S. Lodeweyckx, and J. Driesen, "Gigawatt-hour scale savings on a budget of zero: Deep reinforcement learning based optimal control of hot water systems," *Energy*, vol. 144, pp. 159-168, 2018.
- [95] F. Ruelens, B.J. Claessens, P. Vranckx, F. Spiessens, and G. Deconinck, "Direct load control of thermostatically controlled loads based on sparse observations using deep reinforcement learning," *CSEE Journal of Power and Energy Systems*, vol. 5, no. 4, 423-432, 2019.
- [96] T. Peirelinck, C. Hermans, F. Spiessens, and G. Deconinck, "Domain randomization for demand response of an electric water heater," *IEEE Trans. Smart Grid*, vol. 12, no. 2, pp. 1370-1379, 2021.
- [97] T. Morinibu, T. Noda, and S. Tanaka, "Application of deep reinforcement learning in residential preconditioning for radiation temperature," *International Congress on Advanced Applied Informatics*, 2019.
- [98] A. Gupta, Y. Badr, A. Negahban, and R.G. Qiu, "Energy-efficient heating control for smart buildings with deep reinforcement learning," *Journal of Building Engineering*, <https://doi.org/10.1016/j.jobee.2020.101739>, 2020.
- [99] Y. Yoon and H. Moon, "Performance based thermal comfort control (PTCC) using deep reinforcement learning for space cooling," *Energy and Buildings*, vol. 203, pp. 109420-109430, 2019.
- [100] Y. Sakuma and H. Nishi, "Airflow direction control of air conditioners using deep reinforcement learning," *Proc. of 2020 SICE International Symposium on Control Systems*, Tokushima, Japan, pp. 61-68, 2020.
- [101] X. Valladares, M. Galindo, J. Gutiérrez, W. Wu, K. Liao, J. Liao, K. Lu, and K. Wang, "Energy optimization associated with thermal comfort and indoor air control via a deep reinforcement learning algorithm," *Building and Environment*, vol. 155, pp. 105-117, 2019.
- [102] M. Taylor and P. Stone, "Transfer learning for reinforcement learning domains: A survey," *Journal of Machine Learning Research*, vol. 10, no. 7, pp. 1633-1685, 2009.
- [103] X. Zhang, X. Xin, C. Tripp, D.J. Biagioni, P. Graf, and H. Jiang, "Transferable reinforcement learning for smart homes," *International Workshop on Reinforcement Learning for Energy Management in Buildings & Cities*, 2020.
- [104] Z. Wan, H. Li, and H. He, "Residential energy management with deep reinforcement learning," *International Joint Conference on Neural Networks (IJCNN)*, 2018.
- [105] H. Kumar, P. Mammen, and K. Ramamritham, "Explainable AI reinforcement learning agents for residential cost savings," <http://arXiv:1910.08719v2>, 2019.
- [106] Y. Liu, D. Zhang, and H. Gooi, "Optimization strategy based on deep reinforcement learning for home energy management," *CSEE Journal of Power and Energy Systems*, DOI:10.17775/CSEEJPES.2019.02890, 2019.
- [107] H. Li, Z. Wan, and H. He, "Real-time residential demand response," *IEEE Trans. Smart Grid*, vol. 11, no. 5, pp. 4144-4154, 2020.
- [108] Y. Ye, D. Qiu, X. Wu, G. Strbac, and J. Ward, "Model-free real-time autonomous control for a residential multi-energy system using deep reinforcement learning," *IEEE Trans. on Smart Grid*, DOI:10.1109/TSG.2020.2976771, 2020.
- [109] Z. Cheng, Q. Zhao, F. Wang, Y. Jiang, L. Xia, and J. Ding, "Satisfaction based Q-learning for integrated lighting and blind control," *Energy and Buildings*, vol. 127, pp. 43-55, 2016.
- [110] L. Wang and S. Greenberg, "Window operation and impacts on building energy consumption," *Energy and Buildings*, vol. 92, pp. 313C321, 2015.
- [111] X. Ding, W. Du, and A. Cerpa, "OCTOPUS: deep reinforcement learning for holistic smart building control", *Proc. of BuildSys'19*, 2019.
- [112] P. Wei, S. Xia, R. Chen, J. Qian, C. Li, and X. Jiang, "A deep reinforcement learning based recommender system for occupant-driven energy optimization in commercial buildings," *IEEE Internet of Things Journal*, DOI:10.1109/JIOT.2020.2974848, 2020.
- [113] H. Shuai, H. He, and J.Wen, "Online scheduling of a residential microgrid via monte-carlo tree search and a learned model," <http://arXiv:2005.06161v2>, 2020.
- [114] V. Francois-Lavet, D. Taralla, D. Ernst, and R. Fonteneau, "Deep reinforcement learning solutions for energy microgrids management," *European Workshop on Reinforcement Learning*, 2016.
- [115] D. Dominguez-Barbero, J. Garcia-Gonzalez, M.A. Sanz-Bobi, and E.F. Sanchez-ubeda, "Optimising a microgrid system by deep reinforcement learning techniques," *Energies*, vol. 13, pp. 2830-2847, 2020.
- [116] Y. Ji, J. Wang, J.Xu, X. Feng, and H. Zhang, "Real-time energy management of a microgrid using deep reinforcement learning," *Energies*, vol. 12, pp. 2291-2311, 2019.
- [117] T. Chen and S. Bu, "Realistic peer-to-peer energy trading model for microgrids using deep reinforcement learning," *IEEE PES Innovative Smart Grid Technologies Europe (ISGT-Europe)*, 2019.
- [118] G. Liu, M. Starke, B. Xiao, X. Zhang and K. Tomovic, "Community microgrid scheduling considering building thermal dynamics," *IEEE Power & Energy Society General Meeting*, 2017.
- [119] X. Yang, Y. Wang, H. He, C. Sun, and Y. Zhang, "Deep reinforcement learning for economic energy scheduling in data center microgrids," *IEEE Power & Energy Society General Meeting (PESGM)*, 2019.
- [120] J. Lee, W. Wang, and D. Niyato, "Demand-side scheduling based on deep actor-critic learning for smart grids," <http://arXiv:2005.01979v1>, 2020.

- [121] R. Huang, Y. Chen, T. Yin, Q. Huang, J. Tan, W. Yu, X. Li, A. Li, and Y. Du, "Learning and fast adaptation for grid emergency control via deep meta reinforcement learning," <https://arxiv.org/pdf/2101.05317.pdf>, 2021.
- [122] A. Aswani, N. Master, J. Taneja, A. Krioukov, D. Culler, and C. Tomlin, "Energy-efficient building HVAC control using hybrid system LBMPC," *IFAC Proceedings Volumes*, vol. 45, no. 17, pp. 496-501, 2012.
- [123] R.K. Kalaimani, S. Keshav, and C. Rosenberg, "Multiple time-scale model predictive control for thermal comfort in buildings," *Proc. of e-Energy*, 2016.
- [124] T. Kulkarni, K. Narasimhan, A. Saeedi, and J. Tenenbaum, "Hierarchical deep reinforcement learning: integrating temporal abstraction and intrinsic motivation," *Proc. of NIPS*, 2016.
- [125] K. Li, T. Zhang, and R. Wang, "Deep reinforcement learning for multiobjective optimization," *IEEE Trans. Cybernetics*, DOI:10.1109/TCYB.2020.2977661, 2020.
- [126] X. Chen, A. Ghadirzadeh, M. Björkman, and P. Jensfelt, "Meta-learning for multi-objective reinforcement learning," <https://arxiv.org/abs/1811.03376>, 2019.
- [127] Z. Zhu, K. Lin, and J. Zhou, "Transfer learning in deep reinforcement learning: A survey," <https://arxiv.org/abs/2009.07888>, 2020.
- [128] Z. Jiang and Y. Lee, "Deep transfer learning for thermal dynamics modeling in smart buildings," *Proc. of IEEE International Conference on Big Data (Big Data)*, 2019.
- [129] S. Xu, Y. Wang, Y. Wang, Z. O'Neill and Q. Zhu, "One for many: transfer learning for building HVAC control," *Proc. of BuildSys'20*, 2020.
- [130] X. Zhang, D. Bian, D. Shi, Z. Wang and G. Liu, "Community microgrid planning considering building thermal dynamics," *IEEE Sustainable Power and Energy Conference*, 2019.