

A Learning-Based Power Management Method for Networked Microgrids Under Incomplete Information

Qianzhi Zhang^{1b}, Student Member, IEEE, Kaveh Dehghanpour^{1b}, Member, IEEE,
Zhaoyu Wang^{1b}, Member, IEEE, and Qihua Huang^{1b}, Member, IEEE

Abstract—This paper presents an approximate Reinforcement Learning (RL) methodology for bi-level power management of networked Microgrids (MG) in electric distribution systems. In practice, the cooperative agent can have limited or no knowledge of the MG asset behavior and detailed models behind the Point of Common Coupling (PCC). This makes the distribution systems unobservable and impedes conventional optimization solutions for the constrained MG power management problem. To tackle this challenge, we have proposed a bi-level RL framework in a price-based environment. At the higher level, a cooperative agent performs function approximation to predict the behavior of entities under incomplete information of MG parametric models; while at the lower level, each MG provides power-flow-constrained optimal response to price signals. The function approximation scheme is then used within an adaptive RL framework to optimize the price signal as the system load and solar generation change over time. Numerical experiments have verified that, compared to previous works in the literature, the proposed privacy-preserving learning model has better adaptability and enhanced computational speed.

Index Terms—Distribution systems, networked microgrids, power management, reinforcement learning, adaptive training.

NOMENCLATURE

Indices

- i, j Indices of bus numbers $\forall i, j \in \Omega_I$.
 k Index of line number $\forall k \in \Omega_K$.
 n Index of MG.
 t Index of episode/time instant.

Parameters

- $a_f / b_f / c_f$ Coefficients of the DG quadratic cost function.
 E^{Cap} Max. capacity of ESS unit.
 e_{PV}, e_D Prediction error standard deviations.

AQ1 Manuscript received March 8, 2019; revised June 13, 2019 and July 25, 2019; accepted August 2, 2019. This work was supported by the U.S. Department of Energy Office of Electricity under Grant DE-OE0000839. Paper no. TSG-00346-2019. (Corresponding author: Zhaoyu Wang.)

Q. Zhang, K. Dehghanpour, and Z. Wang are with the Department of Electrical and Computer Engineering, Iowa State University, Ames, IA 50011 USA (e-mail: wzy@iastate.edu).

AQ2 Q. Huang is with Pacific Northwest National Laboratory, Richland, WA 99354 USA (e-mail: qihua.huang@pnnl.gov).

Color versions of one or more of the figures in this article are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TSG.2019.2933502

G/B	Real/imag. parts of the bus admittance matrix.	34
\hat{I}_{PV}	Vectors of solar irradiance estimation.	35
I^{PV}	Real normalized solar irradiance.	36
$P^{Ch/Dis,M}$	Max. ESS charging/discharging limits.	37
P/Q^D	Active/reactive load.	38
$P/Q^{DG,M}$	Max. DG active/reactive power capacity.	39
$p^{DG,R}$	Max. DG ramp limit.	40
P^{PV}	PV active power output.	41
$P/Q^{PCC,M}$	Max. active/reactive power flow at the PCCs.	42
\hat{P}_D	Vectors of aggregate active load estimation.	43
P^D	Real active load.	44
$Q^{PV,M}$	Max. PV reactive power output limit.	45
S	States in Markov decision process.	46
L^M	Max. line flow limit.	47
$SOC^{M/m}$	Max./min. SOC limits.	48
T	Length of the moving decision window.	49
Δt	Time step.	50
α/β	Shape parameters of beta distribution.	51
$\eta^{Ch/Dis}$	Charging/discharging efficiency of ESS unit.	52
λ^F	Diesel generator fuel price.	53
$\lambda^{R,M/m}$	Max./min. retail price limits.	54
λ^W	Wholesale energy price.	55
θ	Vector of regression parameter.	56
θ^*	Vector of converged regression parameter.	57
θ_{Th}/V_{Th}	Threshold value.	58
γ	Discount factor that defines the preference.	59
δ	Step size that defines the rate of learning.	60
μ	Regularization factor.	61
ϕ	Forgetting factor.	62
ϵ	ϵ -greedy exploration factor.	63

Variables

a	Actions in Markov decision process.	65
F	Fuel consumption of DG.	66
SOC	SOC of the battery system.	67
$P^{Ch/Dis}$	Charging/discharging power of ESS unit.	68
P/Q^{DG}	DG active/reactive power outputs	69
P/Q^{ij}	Line active/reactive power flows	70
P/Q^{PCC}	Active/reactive power flow at the PCC.	71
P^W	Exchanged power with the wholesale market.	72
Q^{ESS}	Reactive power outputs of ESS unit.	73
Q^{PV}	PV inverter reactive power outputs.	74
$V/\Delta\theta$	Voltage magnitude and phase angle difference.	75

76	x_p/x_q	MGs power management decision vectors.
77	λ^R	Retail price signals at the PCCs.
78	$u^{Ch/Dis}$	ESS charge/discharge binary variables.

79 Functions

80	$Q_t(S, a)$	State-action value function.
81	$Q_t^*(S, a)$	Optimal state-action value function.
82	$\hat{Q}_t(S, a \theta)$	Parameterized approximate state-action value function.
83		
84	$Q_{S,a}(t \theta)$	Parameterized regression sub-component with state-action interaction.
85		
86	$Q_S(t \theta)$	Parameterized regression sub-component with state values.
87		
88	$Q_a(t \theta)$	Parameterized regression sub-component with action values.
89		
90	$R(t)$	Reward function in Markov decision process.

91 I. INTRODUCTION

92 **A** SMART distribution system consisting of networked
 93 microgrids (MGs), with local Distributed Generators
 94 (DG), Renewable Energy Resources (RES), and Energy
 95 Storage Systems (ESS), can facilitate reliable service provi-
 96 sion to customers in power systems [1]. Smart independent
 97 MGs are considered as a viable solution for electrification
 98 of rural areas, which are excluded from traditional electri-
 99 fication programs due to their remote location and financial
 100 constraints [2]. To ensure the long-term sustainability and
 101 encourage economic development in rural communities, the
 102 feasibility of cooperative business models for rural system
 103 electrification has been analyzed previously [2]–[4]. It has
 104 been shown that a non-profit cooperative can act as an interme-
 105 diary agent between the rural MGs and the wholesale market.
 106 The power is exchanged between the MGs and the cooperative
 107 at a retail rate, and the revenue from electricity sales in the
 108 wholesale market is returned to MGs. The retail energy pricing
 109 program can be used to influence the MGs’ behavior based on
 110 the availability of resources. Real cases of cooperative business
 111 models with rural MGs as participating members can be found
 112 in [3], [4]. The autonomous cooperative business settings in
 113 these cases have been designed to benefit rural communities.

114 Coordinating the real-time behavior of multiple privately-
 115 owned rural MGs in a cooperative business model is a
 116 necessary, yet challenging task [5], [6]. Due to data pri-
 117 vacy and ownership concerns for MGs, the main difficulty
 118 in the way of obtaining a desirable coordination scheme is
 119 the limited access to real-time asset behaviors and models
 120 behind the Point of Common Coupling (PCC) with MGs,
 121 which hinders conventional model-based constrained power
 122 management solvers. This problem becomes more severe as
 123 the penetration of MGs in rural distribution systems grows.
 124 A wide range of methods have been applied in the litera-
 125 ture with the aim of economic operation of the networked
 126 MGs, including methods such as heuristic techniques [7], [8],
 127 centralized decision models [9], [10], constrained hierarchical
 128 control architectures [11]–[13], and distributed optimization
 129 methods [14], [15].

130 However, the functionality of previous models [7]–[15]
 131 highly depends on the full system operator’s knowledge of
 132 MG operation behind the PCC and customers’ private data
 133 at node-level, including nodal demand load consumption,
 134 nodal generation capacities, nodal PV generations, sensitive
 135 cost information, asset constraints, as well as MG network
 136 topology and configuration data. Access to these information
 137 could compromise the data confidentiality and privacy of MGs
 138 and customers that participate in a cooperative business set-
 139 ting. Also, previous methods can be mostly categorized as
 140 “model-based”, since the decision agents depend on detailed
 141 physical models of the distribution systems. One shortcom-
 142 ing of model-based solutions is their inability to adapt to
 143 constantly-changing system conditions when the amount of
 144 measurement data is limited.

145 A promising alternative to model-based optimization
 146 approaches is reinforcement learning (RL), which is a model-
 147 free data-driven technique that can be used to optimize the
 148 behavior of an agent through repeated interactions with its
 149 environment, without full system identification and no *a pri-*
 150 *ori* knowledge of the system. A number of papers have given
 151 examples of how RL techniques can be applied in power
 152 systems. In [16], [17], energy consumption scheduling prob-
 153 lems were solved for single MGs and individual residential
 154 buildings using RL algorithms. However, the above studies
 155 only focus on providing optimal solutions to power man-
 156 agement problems for single entities instead of addressing
 157 coupled decision models for multiple interconnected entities
 158 in a cooperative setting.

159 In this paper, to solve the problem of decision making under
 160 incomplete information while providing decision adaptability,
 161 a bi-level cooperative framework is proposed using an RL-
 162 based method for a distribution system consisting of multiple
 163 networked privately-owned MGs: at Level I of the hierarchy, a
 164 non-profit cooperative agent maximizes the total MGs’ revenue
 165 from power exchange with the wholesale market. This is done
 166 by setting the retail prices, with access only to active/reactive
 167 power measurements at the MG PCCs and *aggregate* load and
 168 solar irradiance information behind the PCCs. The cooperative
 169 agent acts as an intermediary between the MGs and the whole-
 170 sale market, and returns the revenue to the MGs. At Level II
 171 of the hierarchy, each MG Control Center (MGCC) agent
 172 receives the price signal from the cooperative agent and solves
 173 the power-flow-constrained MG power management problem.
 174 The objective at this level consists of the MG operational cost
 175 and the allocated revenue from the cooperative agent. In sum-
 176 mary, the main contributions of this paper can be listed as
 177 follows:

- 178 • The proposed power management system can handle
 179 the current limitations raised from data privacy and
 180 ownership in the cooperative setting. Considering the
 181 model-free nature of our RL-based method, the data pri-
 182 vacy of MGs and the data confidentiality of customers are
 183 maintained. The power management problem is solved
 184 with access to only minimal and aggregated data.
- 185 • The proposed RL solver is faster than conventional
 186 optimization solvers since the learned state-action value
 187 function acts similar to a *memory* that recalls from the

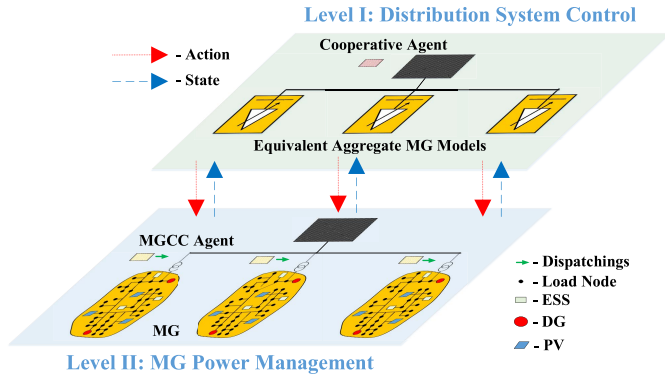


Fig. 1. The architecture of the bi-level networked MGs power management.

cooperative agent's past experiences to estimate new optimal solutions. This is done by updating the state values at each decision window and without re-solving the decision problem.

- The RL framework is trained using a regularized recursive least square methodology with a *forgetting factor*, which enables the decision model to be adaptive to changes in system parameters which are excluded from the cooperative agent's state set.

The remainder of the paper is organized as follows: Section II presents the overall decision hierarchy. Section III elaborates the proposed RL-based framework. Section IV describes the MG power management problem. Simulation results and conclusions are given in Sections V and VI, respectively.

II. OVERALL DECISION HIERARCHY

Fig. 1 gives a general overview of the proposed bi-level power management scheme for a distribution system with multiple MGs, as follows:

Level I - RL-based Distribution System Control: The cooperative agent employs an adaptive model-free RL method, developed using a regularized recursive least square function approximation methodology, to find the optimal retail price signals for the MGs based on the latest system states. This cooperative agent is *non-profit* in the sense that it does not maximize its own profit, but maximizes the social welfare for the whole system, which includes the summation of profits of all the MGs as participating members in the cooperative business model. The price signals are then transmitted to MGCC agents. The RL training process is performed by the cooperative agent through repeated interactions with the MGCC agents. At this level, each MG is modeled as an aggregate controllable load which is price-sensitive. The task of the RL algorithm is to discover the complex relationship between retail price and exchanged power with MGs at PCCs, without direct detailed knowledge of system operation behind the PCCs and only with access to estimations of the solar irradiance and aggregate fixed loads for each MG. Based on the definitions of data privacy and confidentiality in smart grid [18], this approach limits the need for access to local cost and operational constraint data of individual MGs in the

first place. Hence, the proposed method maintains both the privacy of personal information and privacy of behavior for MGs. Moreover, unlike conventional centralized optimization methods, the proposed RL technique does not need customer confidential information at the node-level, such as customer load consumption, as it only uses aggregate data at the MG PCCs for optimal decision making. Furthermore, renewable and load power uncertainty are represented within the learning model state set. To facilitate adaptive conformation to changes in system parameters that are not included in cooperative agent's state set, such as fuel price, a forgetting mechanism has been integrated into the training process to assign higher importance levels to the latest observed data, compared to previous observations.

Level II - MG Power Management: At the second level, the MGCC agents receive the price signal for a look-ahead moving decision window. Based on the received price signals, each MGCC agent solves a constrained Mixed Integer Nonlinear Programming (MINP) to dispatch their local generation/storage assets to maximize their revenue (or equivalently minimize their cost) in the price-based environment, subject to full AC power flow constraints. Each MG's total revenue includes the cost of operation and the allocated revenue received from the cooperative agent. Based on the solution to this problem, each individual MGCC agent determines the exchanged active and reactive power with the distribution system at PCC.

Note that the RL-based reward maximization problem at Level I is subject to the power-flow-constrained response of MGs at Level II. Since the MGs are sensitive to electricity price, the reward value cannot be maximized by setting the price to its highest value. This will lead to the maximum DG generation, which will result in a decline in the cooperative agent's revenue. Hence, optimal price is reached based on a tradeoff between MGs' over-generation (when price is too high) and over-consumption (when price is too low). Also, note that the response of MGs itself is explicitly constrained by network power flow constraints.

III. LEVEL I: ADAPTIVE RL-BASED DISTRIBUTION SYSTEM CONTROL

At the first level of the hierarchy, a non-profit cooperative agent is in charge of setting the retail price of electricity at different times to maximize the revenue from power exchange with wholesale market, which will be allocated between MGs. This problem is formulated and solved over a moving decision window of length T . The difficulty in solving this problem is that the cooperative agent has incomplete knowledge of MGs' asset control and management data. To solve this problem, an RL approach is adopted, in which the decision making cooperative agent observes the response of its environment, consisting of networked MGs, to its actions at different states. Based on the received reward/cost signals from its environment and without explicit modeling, the cooperative agent searches for actions that optimize its expected accumulated received rewards at different system states.

284 A. Proposed RL-Based Method Structure

285 A RL framework consists of a Markov decision process
286 including a set of states ($\mathcal{S} \in \mathcal{S}$), a set of actions ($\mathbf{a} \in$
287 \mathcal{A}), a reward function ($\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$), and a state-
288 action value function corresponding to each state-action pair
289 ($Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$). These components are defined for the
290 problem at hand, as follows.

291 1) *State Set Definition*: In this paper, the system state,
292 which is denoted by $\mathbf{S}(t) = (\mathbf{S}_1(t), \dots, \mathbf{S}_N(t))^\top$ at time t , is
293 a concatenation of MGs' local state vectors ($\mathbf{S}_n(t)$ for the n^{th}
294 MG) defined as:

$$295 \quad \mathbf{S}_n(t) = \left\{ \hat{I}_{PV}(t, n), \hat{P}_D(t, n) \right\} \quad (1)$$

296 where, $\hat{I}_{PV}(t, n)$, $\hat{P}_D(t, n)$ are the vectors of solar irradiance
297 estimation, and aggregate active load power estimation for the
298 n^{th} MG at time t , respectively. Hence, to define the global state,
299 the cooperative agent needs to estimate or predict the uncer-
300 tain aggregate solar irradiance and load at the PCC for each
301 MG. To represent the uncertainty of the prediction process,
302 prediction error values are considered to the actual underlying
303 solar irradiance and load values, as shown below:

$$304 \quad \hat{I}_{PV}(t, n) \sim \text{Beta}(\alpha, \beta) \quad (2a)$$

$$305 \quad \alpha = \frac{\beta(\sum_i I_{i,t,n}^{PV})}{(1 - \sum_i I_{i,t,n}^{PV})} \quad (2b)$$

$$306 \quad \beta = \left(1 - \sum_i I_{i,t,n}^{PV}\right) \left(\frac{\sum_i I_{i,t,n}^{PV}(1 + \sum_i I_{i,t,n}^{PV})}{e_{PV}^2} - 1\right) \quad (2c)$$

$$307 \quad \hat{P}_D(t, n) \sim \mathcal{N}\left(\sum_i P_{i,t,n}^D, e_D^2(t)\right) \quad (2d)$$

308 where, $\sum_i I_{i,t,n}^{PV}$ and $\sum_i P_{i,t,n}^D$ are the real aggregate normalized
309 solar irradiance and load over the decision window, and e_{PV}
310 and e_D are the beta and Gaussian estimation error standard
311 deviations. The values of parameters of beta and Gaussian
312 distributions are adopted from the [19]–[21].

313 2) *Action Set Definition*: Given the definition of model
314 states, the global action vector is similarly defined by the retail
315 price signals at the PCCs with MGs, denoted as $\lambda_{t,n}^R$ for the
316 n^{th} MG, $\mathbf{a}(t) = (\lambda_{t,1}^R, \dots, \lambda_{t,N}^R)^\top$.

317 3) *Reward Function Definition*: The reward function at
318 time t represents the discounted accumulated revenue of the
319 cooperative agent over the moving decision window with
320 length T :

$$321 \quad R(t) = \sum_{t'=0}^{T-1} \gamma^{t'} \left(\lambda_{t+t'}^W P_{t+t'}^W - \sum_{n=1}^N \lambda_{t+t',n}^R P_{t+t',n}^{PCC} \right) \quad (3)$$

322 where, γ is a discount factor ($0 \leq \gamma \leq 1$) that
323 defines the cooperative agent's preference for the immediate
324 reward, defined as the revenue at time t , $\pi(t) = \lambda_t^W P_t^W -$
325 $\sum_{n=1}^N \lambda_{t,n}^R P_{t,n}^{PCC}$. Also, λ_t^W denotes the wholesale energy price,
326 P_t^W is the exchanged power with the wholesale market, where
327 $P_t^W \leq 0$ represents power import from the wholesale market.
328 $P_{t,n}^{PCC}$ is the active power transfer between grid and the n^{th}
329 MG through the PCC, where $P_{t,n}^{PCC} \geq 0$ implies export from
330 MGs to grid. The extreme case of $\gamma = 0$ represents a myopic

cooperative agent, which favors only the immediate economic 331
rewards and assigns zero weights to future expected rewards. 332
However, as the discount factor increases the cooperative agent 333
starts to include future expected rewards into its optimal deci- 334
sion problem. Hence, when the discount factor reaches $\gamma = 1$ 335
the cooperative agent assigns equal weights to all the expected 336
reward values for all the time instants in the decision window. 337

4) *State-Action Value Function Parameterization*: To 338
optimize the cooperative agent's action, an auxiliary state- 339
action value function is formed, denoted as $Q(\mathcal{S}, \mathbf{a})$, which 340
can be thought of as a replacement for the explicit system 341
model. The state-action value function determines the long- 342
term accumulated expected reward given the current state and 343
action vectors: 344

$$Q_t(\mathcal{S}, \mathbf{a}) = E \left\{ \sum_{t'=0}^{T-1} \gamma^{t'} \pi(t+t') | \mathcal{S}(t) = \mathcal{S}, \mathbf{a}(t) = \mathbf{a} \right\} \quad (4)$$

where, $Q_t(\mathcal{S}, \mathbf{a})$ is the expected accumulated reward if the ini- 346
tial starting state is $\mathcal{S}(t)$, while the selected initial action is 347
 $\mathbf{a}(t)$, and the latest optimal policy is followed for every other 348
time-step in the future. The expectation operator $E\{\}$ is calcu- 349
lated with respect to the future expected action-states, which 350
in this case are in turn functions of the solar-load uncertain 351
powers. 352

The goal of RL is to learn an optimal state-action value 353
function, $Q_t^*(\mathcal{S}, \mathbf{a})$, that satisfies the Bellman optimality equa- 354
tion [22], as follows: 355

$$Q_t^*(\mathcal{S}, \mathbf{a}) = E \left\{ \pi(t+1) + \gamma \cdot \max_{\mathbf{a}'} Q_t^*(\mathcal{S}(t+1), \mathbf{a}') \right\} \quad (5)$$

Since solving (5) directly is not possible, RL provides a 357
framework to obtain the optimal state-action value function 358
which satisfies (5) using an iterative episodic learning envi- 359
ronment. To implement this framework for the cooperative 360
agent interacting with multiple MGs, the state-action value 361
function is parameterized employing a multivariate polynomial 362
regression approximation technique [22]–[24], defined by \hat{Q}_t , 363
which consists of three multivariate polynomial elements with 364
maximum degree 2: 365

$$Q_t(\mathcal{S}, \mathbf{a}) \approx \hat{Q}_t(\mathcal{S}, \mathbf{a} | \boldsymbol{\theta}) = Q_{\mathcal{S}, \mathbf{a}}(t | \boldsymbol{\theta}) + Q_{\mathcal{S}}(t | \boldsymbol{\theta}) + Q_{\mathbf{a}}(t | \boldsymbol{\theta}) \quad (6)$$

Given the regression parameter vector $\boldsymbol{\theta}$, $Q_{\mathcal{S}, \mathbf{a}}$, $Q_{\mathcal{S}}$, and 367
 $Q_{\mathbf{a}}$ are the parameterized sub-components that quantify the 368
impacts of state-action interaction $Q_{\mathcal{S}, \mathbf{a}}(t | \boldsymbol{\theta})$, state values 369
 $Q_{\mathcal{S}}(t | \boldsymbol{\theta})$, and action values $Q_{\mathbf{a}}(t | \boldsymbol{\theta})$, respectively. These regres- 370
sion sub-components in multivariate polynomial regression 371
model are defined as follows: 372

$$Q_{\mathcal{S}, \mathbf{a}}(t | \boldsymbol{\theta}) = \sum_{n=1}^N \theta_{t,n}^1 \lambda_{t,n}^R \hat{I}_{PV}(t, n) + \sum_{n=1}^N \theta_{t,n}^2 \lambda_{t,n}^R \hat{P}_D(t, n) \quad (7)$$

$$Q_{\mathcal{S}}(t | \boldsymbol{\theta}) = \sum_{n=1}^N \theta_{t,n}^3 \hat{I}_{PV}(t, n) + \sum_{n=1}^N \theta_{t,n}^4 \hat{P}_D(t, n) \quad (8)$$

$$Q_{\mathbf{a}}(t | \boldsymbol{\theta}) = \sum_{n=1}^N \theta_{t,n}^5 \lambda_{t,n}^R + \theta^6 \quad (9)$$

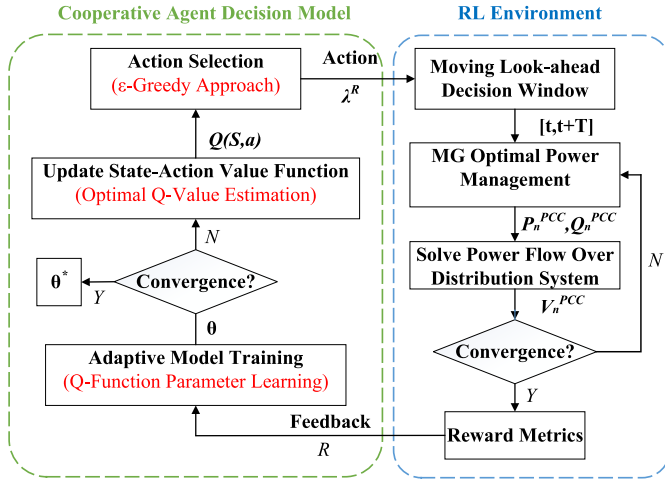


Fig. 2. Proposed RL-based framework.

where, $\theta = \{\theta_{t,n}^k, \theta^k\}$ constitute the parameters of the approximate state-action value function that have to be learned by the cooperative agent through repeated interaction with the MGs. Together these three components form a bilinear regression model to parametrize the state-action value function (i.e., the regression model is linear with respect to each of its arguments.) The reason for selecting a bilinear regression model is the structure of the reward function (3), which also follows a bilinear relationship between the price signal and the aggregate power measured at MG PCCs and the substation. Furthermore, the state-action value parameterization shown in (7)-(9) offers two critical advantages compared to other types of function approximators: 1) Using a bilinear regression model will simplify optimal action selection procedure considerably, as will be shown in Section III-B. For instance, if an artificial neural network is used, optimal action selection becomes intractable. However, using the proposed bilinear regression model, optimal action selection reduces to linear programming, which can be solved easily. 2) A basic challenge in choosing the form of a function approximator is the trade-off between over-parametrization and estimation accuracy. For example, as we increase the degree of the multivariate polynomial approximator the value estimation accuracy for new state-action pairs would also improve; however, at some point the function approximator becomes over-parameterized and will start overfitting to the available data, at which point the performance declines. We observed that by limiting the degree of the multivariate polynomial degree to 2, the best estimation accuracy can be achieved while maintaining a safe margin to avoid overfitting under various practical case studies.

B. Adaptive RL-Based Method Training

To achieve this task we have adopted an adaptive episodic learning mechanism, which is shown in Fig. 2. Each episode in the learning process corresponds to an online decision instant. Hence, as the decision window rolls along time new episodes are perceived by the cooperative agent. The learning process has the following steps.

Step 1 (Initialization): The time index is initialized as $t = t_0$, representing the first episode. The parameters of the state-action value function are initialized, $\theta \leftarrow \theta(t_0)$. The initial state of the system, corresponding to solar irradiance and aggregate load of all the MGs for the decision window $[t_0, t_0 + T]$ is predicted, $\mathcal{S}(t_0), \dots, \mathcal{S}(t_0 + T)$. Note that these predicted states, while representing system uncertainty, are updated continuously as the decision window rolls along time.

Step 2 (ε-greedy Action Selection): Based on the latest state-action value function defined by parameter θ , the optimal actions are estimated for the decision window $[t, t + T]$ to maximize the cooperative agent's accumulated reward, as follows:

$$\begin{aligned} \mathbf{a}_{opt}(t) &= \arg \max_{\mathbf{a}'} Q_t(\mathcal{S}(t'), \mathbf{a}') \\ \text{s.t. } \mathbf{a}' &= (\lambda_{t',1}^R, \dots, \lambda_{t',N}^R)^\top \\ \lambda_{t',m}^{R,m} &\leq \lambda_{t',i}^R \leq \lambda_{t',i}^{R,M}, \forall i = \{1, \dots, N\} \\ \forall t' &= \{t, \dots, t + T\} \end{aligned} \quad (10)$$

where, $\rho_\lambda = [\lambda_{t',m}^{R,m}, \lambda_{t',i}^{R,M}]$ defines the minimum/maximum range of action for retail price. Note that given the parameterization for $Q_t(\mathcal{S}, \mathbf{a})$ in (7)-(9), (10) is basically a set of linear programs, which can be solved efficiently using off the shelf solvers. A critical aspect of (10) is that the obtained optimal action, $\mathbf{a}_{opt}(t)$, is calculated with respect to the latest state-action value function, which could be far from being accurate in the early stages of training. Hence, to reduce the risk of sub-optimality and to strike a balance between exploration and exploitation of decision space, an ϵ -greedy action selection method [22] is adopted, with $0 \leq \epsilon \ll 1$, to select the cooperative agent's action at time t :

$$\mathbf{a}(t) = \begin{cases} \mathbf{a}_{opt}(t) & \text{if } r \geq \epsilon \\ \lambda_{t',i}^R \sim U\{\rho_\lambda\} \forall i & \text{if } r < \epsilon \end{cases} \quad (11)$$

where, r is a random number selected uniformly, $r \sim U\{[0, 1]\}$, with $U\{\mathbf{A}\}$ representing uniform probability distribution over the set \mathbf{A} . The randomization (11) promotes continuous exploration of action space to improve the outcome of the learning process. Upon obtaining the action vector $\mathbf{a}(t)$, retail price signals are sent to each MGCC agent.

Step 3 (Networked MG Power Management): Based on the received price signals, $\lambda_{t',n}^R, \forall n, t' = \{t, \dots, t + T\}$, each MGCC agent solves its optimal power management problem (Section IV). Based on the solutions at this stage, the aggregate power injection/withdrawal to/from the grid are obtained at the PCCs with the MGs, denoted as $P_{t',n}^{PCC}$ and $Q_{t',n}^{PCC}, \forall n, t' = \{t, \dots, t + T\}$.

Step 4 (Accumulated Reward Calculation): Based on the outcomes of the MG power managements, the net power exchange with the wholesale market, P_t^W , is determined and used to calculate the discounted accumulated revenue for the decision window $[t, t + T]$, using (3).

Step 5 (Adaptive Model Training): Using the observed reward signal, the regression models defined in (7)-(9) are updated, based on a gradient descent approach to modify the parameters in the direction of improving the generalization

465 capacity of the state-action value function [22]:

$$466 \quad \theta(t+1) \leftarrow \theta(t) + \delta \left\{ R(t) - \hat{Q}_t(\mathcal{S}, \mathbf{a}|\theta) \right\} \nabla_{\theta} \hat{Q}_t(\mathcal{S}, \mathbf{a}|\theta) \quad (12)$$

467 where, δ is the step size that defines the rate of learning. Note
 468 that ideally we require $\hat{Q}_t(\mathcal{S}, \mathbf{a}|\theta) = R(t)$, which implies that
 469 the approximate state-action value function is able to accu-
 470 rately predict the accumulated reward. Accordingly, (12) is
 471 devised to reduce this prediction error over time. To imple-
 472 ment (12), two points have to be taken under consideration:
 473 1) since data acquisition and the training process both depend
 474 on cooperative agent action selection, approximate RL algo-
 475 rithms are known to be prone to overfitting and over-estimation
 476 of the values of state-action pairs [25]. Hence, a regular-
 477 ization mechanism has to be adopted to reduce the risk of
 478 overfitting, 2) the distribution system parameters are subject
 479 to change over time. These time-varying parameters, such as
 480 price of fuel, are not directly captured in the Markov decision
 481 process's state definition. This makes the learned model sus-
 482 ceptible to failure in case considerable changes occur in the
 483 values of these parameters. Hence, the training process needs
 484 to be *adaptive* to enable cooperative agent to quickly conform
 485 to new system conditions. To implement (12) while consider-
 486 ing the above-mentioned points, a regularized recursive least
 487 squares algorithm with exponential forgetting is designed [26].
 488 The regression parameters are updated recursively, as follows:

$$489 \quad \theta(t+1) \leftarrow \theta(t) + \Delta(t)\mathbf{x}(t) \left\{ R(t) - \hat{Q}_t(\mathcal{S}, \mathbf{a}|\theta) \right\} \quad (13)$$

$$490 \quad \Delta(t+1) \leftarrow \hat{\Delta}(t+1) \left(I + \mu \hat{\Delta}(t+1) \right)^{-1} \quad (14)$$

$$491 \quad \hat{\Delta}(t+1) \leftarrow \frac{1}{1-\phi} \left(\Delta(t) - \frac{\Delta(t)\mathbf{x}(t)\mathbf{x}(t)^{\top}\Delta(t)}{1+\mathbf{x}(t)^{\top}\Delta(t)\mathbf{x}(t)} \right) \quad (15)$$

492 where, $\mathbf{x}(t) = (\mathcal{S}(t), \mathbf{a}(t))^{\top}$ represents the latest cooperative
 493 agent's observation, Δ is an auxiliary matrix mimicking the
 494 regression pseudo-inverse matrix, μ is the regularization fac-
 495 tor which is used for re-scaling the model covariance, and
 496 $0 \leq \phi < 1$ is the forgetting factor. The regularization fac-
 497 tor acts as a weight for penalizing the Euclidean norm of
 498 parameter vector (i.e., $\|\theta\|_2$) in a ridge regression setting to
 499 prevent overfitting. The forgetting factor enables the coop-
 500 erative agent to "forget" its earlier experiences in favor of
 501 the newer observations by assigning lower weights to the
 502 previously learned parameters. Hence, the forgetting factor
 503 introduces an exponential attenuation of data history over
 504 time.

505 *Step 6 (State Transition):* The decision window is moved
 506 forward to the new episode, $t \leftarrow t+1$. The new system state
 507 for the decision window, $[t, t+T]$ is predicted and denoted as
 508 $\{\mathcal{S}(t), \dots, \mathcal{S}(t+T)\}$.

509 IV. LEVEL II: MGCC AGENT POWER MANAGEMENT

510 At Level II, each MG receives the price signals from the
 511 cooperative agent to solve the constrained optimal power
 512 management problem within a moving decision window indi-
 513 vidualy, as shown in the paper Appendix, (16)-(40). Each MG
 514 is comprised of local DGs, ESS, solar Photo-Voltaic (PV) pan-
 515 els and a number of loads. Hence, to account for the impacts

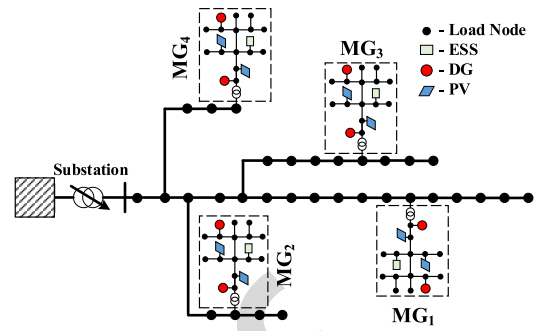


Fig. 3. Test system under study.

TABLE I
RL-BASED METHOD PARAMETERS

Parameters	γ	δ	μ	ϕ	ϵ
Values	0.99	0.01	1×10^{-5}	0.01	0.1

of MGs on each other, the MG-level optimal power flow solver
 is based on an interactive non-linear programming algorithm.
 The steps of the interactive power flow solution are as follows:

Step I (Receive Input Signals From Level I): The MGs
 receive the retail price signals at the PCCs, $\lambda_{t,n}^R$, from the
 cooperative agent.

*Step II (Solve Individual MG Optimal Power Management
 Problem):* Given $\lambda_{t,n}^R$ and the estimated voltage at PCC, the
 power management problem (16)-(40) is solved independently
 by each MGCC, and the exchanged active and reactive powers
 at the PCCs are obtained for each MG.

*Step III (Solve Power Flow Problem Over Distribution
 System):* Treating MGs as fixed PQ loads in the external
 distribution system, power flow is solved over the network
 connecting the MGs. The total substation exchanged power,
 P_t^W , and voltage values at PCCs, $V_{t,n}^{PCC}$, are updated based on
 the power flow solution.

Step IV (Check Convergence): Go back to Step III to update
 PQ values corresponding to each MG, until the changes in
 voltage values at MG PCCs are smaller than a threshold
 value V_{Th} .

To summarize, the pseudo-code of the proposed bi-level RL-
 based framework has been shown in Algorithm 1.

539 V. NUMERICAL RESULTS

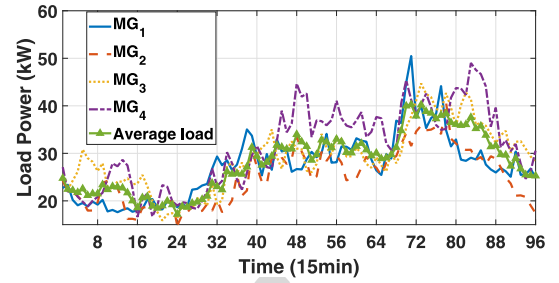
540 The proposed method is tested on a modified medium volt-
 541 age 33-bus distribution network [27], which has been widely
 542 used for studies pertaining to distribution system [28]. The
 543 case study consists of four MGs as shown in Fig. 3. Each MG
 544 is modeled as a modified IEEE 13-bus network at a low voltage
 545 level [29]. Hence, the system has a total number of 85 nodes.
 546 To represent a realistic model, we simulated an unbalanced
 547 system, where the loads and generators are almost uniformly
 548 distributed across phases. Note that the proposed model-free
 549 power management technique applies to both balanced and
 550 unbalanced systems. Table I presents all setting parameters
 551 for the proposed RL-based method in this paper.

Algorithm 1 Bi-Level RL-Based Power Management Method

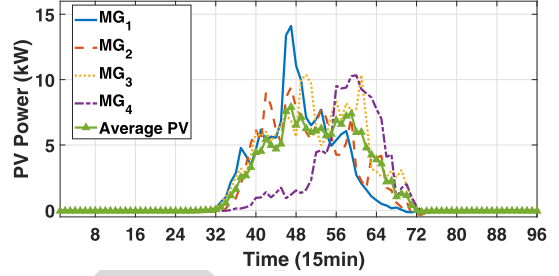
```

1: Select  $T, \gamma, \delta, \mu, \phi, \epsilon, \theta(t_0)$ 
2: procedure LEVEL I: RL ACTION SELECTION( $\theta$ )
3:    $t \leftarrow 1$ 
4:    $\mathbf{S} \leftarrow [\mathbf{S}(t), \dots, \mathbf{S}(t+T)]$ 
5:    $Q_t(\mathbf{S}, \mathbf{a}) \leftarrow \hat{Q}_t(\mathbf{S}, \mathbf{a}|\theta)$ 
6:    $\mathbf{a}_{opt}(t) \leftarrow$  Solve linear program (10)
7:    $\lambda_{t,i}^R \sim U\{\rho_\lambda\}$ 
8:    $r \sim U\{[0, 1]\}$ 
9:   if  $r \geq \epsilon$  then
10:     $\mathbf{a}(t) \leftarrow \mathbf{a}_{opt}(t)$ 
11:   else
12:     $\mathbf{a}(t) \leftarrow \lambda_{t,i}^R$ 
13:   end if
14: end procedure
15: procedure LEVEL II: MGCC AGENT POWER MANAGEMENT( $\mathbf{a}$ )
16:    $k \leftarrow 1$ 
17:    $\lambda^R \leftarrow \mathbf{a}(t), V_n(k) \leftarrow V_{t,n}^{PCC}$ 
18:    $P_{t,n}^{PCC}, Q_{t,n}^{PCC} \leftarrow$  Solve (16)-(40)  $\forall n$  with  $V_n(k)$ 
19:    $V_n(k) \leftarrow$  Solve power flow with  $\{P_{t,n}^{PCC}, Q_{t,n}^{PCC}\}$ 
20:   if  $\Delta|V_n| \geq V_{Th}$  then
21:      $k \leftarrow k + 1$ 
22:     Go back to Step 18
23:   else
24:     Go to Step 27
25:   end if
26: end procedure
27: procedure LEVEL I: RL UPDATE STATE-ACTION VALUE FUNCTION( $P^{PCC}, P^W, \mathbf{S}, \mathbf{a}, \theta$ )
28:    $R(t) \leftarrow \sum_{t'=0}^{T-1} \gamma^{t'} (\lambda_{t+t',n}^R P_{t+t',n}^W - \sum_{n=1}^N \lambda_{t+t',n}^R P_{t+t',n}^{PCC})$ 
29:    $\hat{Q}_t(\mathbf{S}, \mathbf{a}|\theta) \leftarrow Q_{S\mathbf{a}}(t|\theta) + Q_S(t|\theta) + Q_{\mathbf{a}}(t|\theta)$ 
30:    $\theta(t+1) \leftarrow \theta(t) + \delta\{R(t) - \hat{Q}_t(\mathbf{S}, \mathbf{a}|\theta)\} \nabla_{\theta} \hat{Q}_t(\mathbf{S}, \mathbf{a}|\theta)$ 
31:   if  $\|\theta(t+1) - \theta(t)\| \geq \theta_{Th}$  then
32:      $t \leftarrow t + 1$ 
33:     Go back to Step 4
34:   else
35:      $\theta^* \leftarrow \theta(t+1)$ 
36:     Output  $\theta^*$ 
37:   end if
38: end procedure

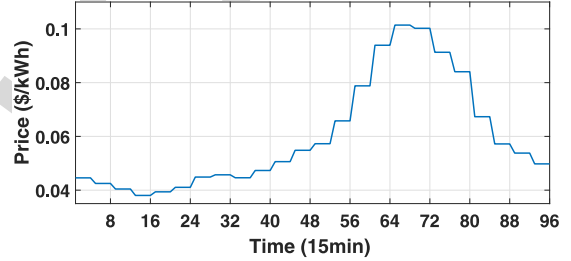
```



(a) Aggregate active load profile of the MGs



(b) Aggregate PV power of the MGs



(c) Wholesale market price

Fig. 4. Input data for the case study.

main grid under optimal price actions, which are the responses of each MG to the actions, are shown in Fig. 6. These figures show the correlation between MGs' behavior and the retail price signal. This demonstrates the mutual impacts of the two levels of the decision model. As the wholesale price increases, the cooperative agent increases the retail prices to encourage the MGs to produce more power to reduce the costs of power purchase from the wholesale market. It can be observed that, most of the time, the cooperative agent exports power to the heavily loaded MGs to maintain power balance in the system. The reason for this is that MGs cannot provide their local demand consumption by their own local generation and have to purchase power from the cooperative service provider. The overall operational costs of MGs have been compared with and without a cooperative agent as an intermediary between the wholesale market and MGs. As can be seen from Fig. 7, the total operational costs of each MG are reduced due to the returned revenue from the cooperative service provider. Therefore, as an intermediary between the MGs and the wholesale market, the cooperative agent can help MGs to reduce their overall operational cost. Hence, it is in the interest of the MGs to participate in the wholesale market through the non-profit cooperative agent.

552 *A. System Operation Outcomes*

553 The *aggregate* active load profiles of all the MGs and
554 the average load are presented in Fig. 4(a). The *aggregate*
555 solar active generations in each MGs have been shown in
556 Fig. 4(b). Both load demands and PV generations data with
557 15 minutes time resolution are obtained from smart meters to
558 provide realistic numerical experiments. The wholesale market
559 prices used in the numerical case study have been shown
560 in Fig. 4(c), which are adopted from the historical whole-
561 sale electricity market data from U.S. Energy Information
562 Administration [30].

563 The retail price signals for the MGs, which are the optimal
564 actions from Level I of the proposed RL-based model, are
565 presented in Fig. 5. Power exchange between MGs and the

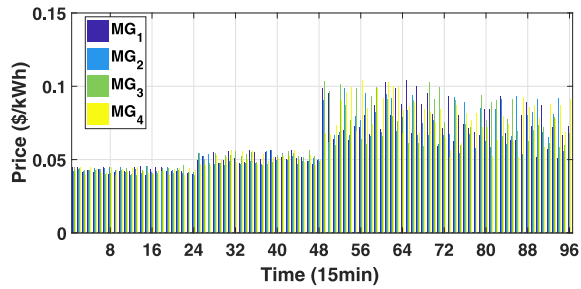


Fig. 5. Optimal retail price signals (Level I actions).

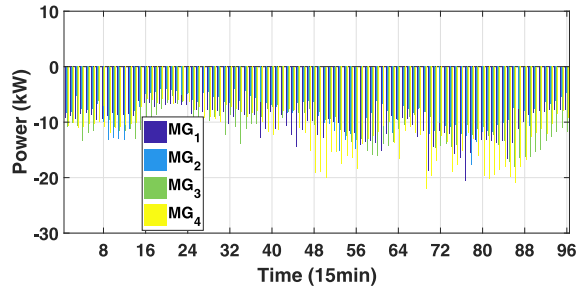


Fig. 6. Optimal power transfer through PCC of MGs (Level II responses to optimal actions).

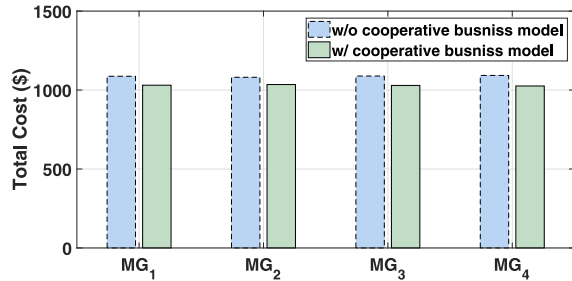


Fig. 7. Comparison of total operational cost of MGs.

TABLE II
COMPARISON WITH A CENTRALIZED OPTIMIZATION METHOD

	RL-based method	Centralized Opt.
Social welfare (\$)	4232.264	4212.372
Computational time (s)	9.64	116.35
MG privacy maintenance	Yes	No

of three different decision windows have been used for a new 611
 decision window without re-training. In Fig. 8, optimal power 612
 transfers are compared for four scenarios representing four 613
 distinct decision windows: in each scenario the RL training 614
 is performed for one of the decision windows from random 615
 initial conditions, while the updated aggregate MG solar gen- 616
 eration and load demand from that decision window are simply 617
 inserted into the learned state-action value functions obtained 618
 from the other three decision windows. Then, the optimal 619
 actions are calculated for each decision window. As can be 620
 seen, for all scenarios the optimal solutions are close to each 621
 other and almost identical. This shows that the state-action 622
 value function learned from other decision windows can be 623
 used reliably in new situations using updated state information. 624
 Hence, the RL model does not necessarily need to be trained 625
 from scratch, and the latest learned function approximator can 626
 be simply used to update the cooperative agent's decisions. In 627
 practice, however, the re-training process has to be performed 628
 with a user-defined frequency depending on the rate of change 629
 of system parameters. 630

Therefore, the RL-based method has two fundamental 631
 advantages over centralized optimization method: 1) RL is 632
 model-free; hence, unlike centralized optimization approaches, 633
 it does not require detailed private knowledge of MG systems 634
 to reach the optimal solution. 2) RL is much faster com- 635
 pared to centralized solvers since the learned state-action value 636
 function, which acts similar to a *memory*, is able to leverage 637
 the cooperative agents past experiences to obtain new optimal 638
 solutions by generalizing to new unseen states. 639

590 B. Benefits of RL-Based Method

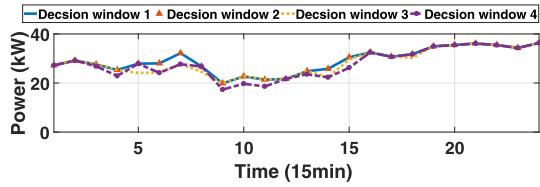
591 A numerical comparison between a centralized off the shelf
 592 solver [31] versus the proposed method for the multiple MGs
 593 power management problem is shown in Table II. In this table,
 594 the total social welfare is defined as the summation of the
 595 cooperative agent's accumulated reward and the operational
 596 cost of all the MGs. Ideally both of the solvers should output
 597 the global optimal solution to the problem. As can be seen, the
 598 difference between the solutions obtained by the centralized
 599 solver with complete system information, and the proposed
 600 RL method under incomplete information is less than 0.5%
 601 of the total achieved welfare. Note that while the initial RL
 602 training stage can be time-consuming, the decision time is
 603 much smaller than that of a centralized optimization method,
 604 upon convergence. This is due to the fact that the proposed RL-
 605 based method is able to receive continual updates over time,
 606 which enables the decision framework to reach a solution in
 607 real-time without the need to solve a large-scale optimization
 608 problem at each time instant.

609 To further demonstrate this, we have performed numerical
 610 experiments in which the trained state-action value functions

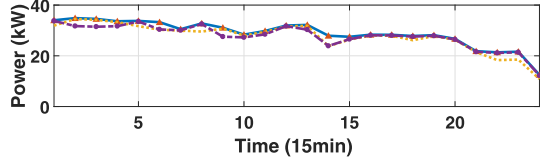
C. Adaptive RL Results 640

To verify the functionality of the RL framework, the esti- 641
 mated reward obtained from the multiple linear regression is 642
 compared with the actual reward at each episode, as shown 643
 in Fig. 9. As can be seen, at the earlier stages of the learn- 644
 ing process, the difference between the estimated reward and 645
 the real reward is relatively high. However, as the number of 646
 episodes increases, this difference drops to within an accept- 647
 able range. The results imply that the cooperative agent is able 648
 to accurately estimate the response of MGs to control actions. 649
 Hence, using the proposed RL approach the cooperative agent 650
 is able to track the behavior of MGs and maximize the reward 651
 through continuous interactions. 652

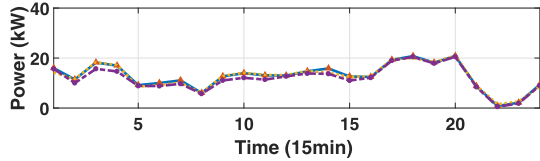
To test the adaptability of the learning framework against 653
 changes in parameters that have not been included in the 654
 definition of state set and are not directly observed by the 655
 cooperative agent, a numerical scenario is devised. At a point 656
 in time (episode $t = 250 h$), the DG fuel price is doubled. The 657
 reward estimation mean absolute percentage error (MAPE) 658



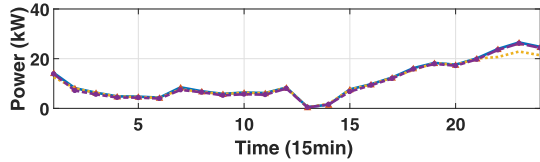
(a) Optimal action for decision window 1, using the trained models of decision windows 2, 3, and 4 for comparison



(b) Optimal action for decision window 2, using the trained models of decision windows 1, 3, and 4 for comparison



(c) Optimal action for decision window 3, using the trained models of decision windows 1, 2, and 4 for comparison



(d) Optimal action for decision window 4, using the trained models of decision windows 1, 2, and 3 for comparison

Fig. 8. Verifying the accuracy of previously-learned models under new state scenarios from different decision windows (memory effect).

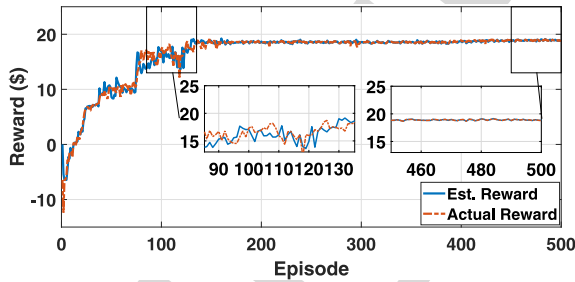
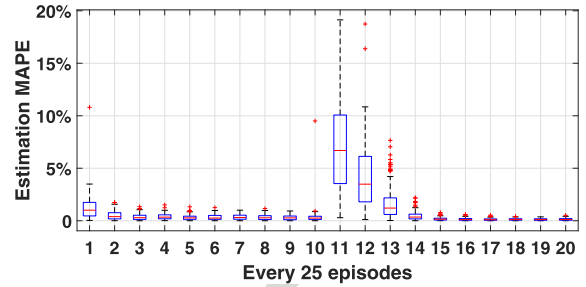
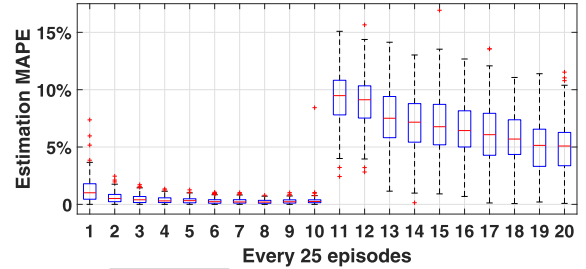


Fig. 9. Performance of the proposed reward function approximation.



(a) Estimation MAPE with forgetting factor (Highly adaptive)



(b) Estimation MAPE without forgetting factor (Slow adaptation)

Fig. 10. Adaptability of the proposed RL-based method.

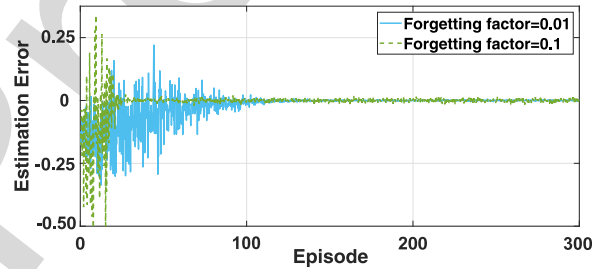


Fig. 11. Impact of forgetting factor on learning convergence.

in Fig. 10(b). As can be seen, compared to the proposed adaptive RL-based method with forgetting factor, the conventional RL-based method without forgetting factor shows slow adaptation to changes in parameters. For this case, our RL-method is able to achieve 25% improvement in the convergence constant over conventional RL.

In Fig. 11, the impact of forgetting factor on the convergence of the RL framework is demonstrated. This figure shows the RL-based reward estimation error for the cooperative agent under two different forgetting factor values. As the forgetting factor increases from 0.01 to 0.1, the convergence speed of the RL framework has been improved. Hence, the forgetting factor controls the rate of adaptiveness to new conditions. However, a tradeoff exists between the rate of convergence and the accuracy of the solution. As can be seen, higher forgetting factors also lead to higher variances in the estimation error signal.

VI. CONCLUSION

Smart distribution systems with networked MGs in a cooperative setting can facilitate reliable power delivery to customers in future rural power grids. However, cooperatives can have incomplete knowledge of MG members' operational parameters due to data privacy and ownership concerns,

with forgetting factor is shown in Fig. 10(a). As can be seen, upon the occurrence of the sudden change in fuel price, the learning MAPE temporarily jumps to a very high value since the cooperative agent is now facing a new unknown environment, as the price of fuel is not included within the cooperative agent's Markov decision process. However, as the learning process with forgetting proceeds, the MAPE drops to within acceptable range once more. The cooperative agent can still track the actual underlying reward signal as the number of episodes increases with the sudden parameter changes. The reward estimation MAPE without forgetting factor is shown

692 which is an obstacle in the way of optimal decision mak-
693 ing. Motivated by the shortcomings of model-based multiple
694 MG power management in distribution systems with lim-
695 ited observability, this paper presents an adaptive RL-based
696 methodology for bi-level power management of cooperatives
697 consisting of multiple networked MGs.

698 We have shown that: 1) using the proposed decision method,
699 a cooperative agent is able to accurately track the behavior
700 of multiple networked MGs under incomplete knowledge of
701 operation variables behind the PCCs. This can be used to indi-
702 rectly control the response of participants in a price-based
703 environment. 2) The proposed RL-based method is able to gen-
704 eralize from its past experiences to estimate optimal solutions
705 in new situations without re-training from random initial con-
706 ditions (i.e., fast response under evolving system conditions).
707 This immensely speeds up the power management compu-
708 tational process. 3) The framework is shown to be adaptive
709 against the changes happening to unobserved parameters that
710 are excluded from cooperative agent's state set. The learning
711 model has been tested and verified using extensive numerical
712 scenarios. To summarize, the proposed decision model shows
713 better adaptability, solution quality, and computational time
714 compared to conventional centralized optimization methods.

715 The current RL-based decision model is limited to the
716 power management of a single cooperative service provider
717 with multiple MGs. However, in more realistic cases, there
718 could also be multiple cooperative service providers in an
719 interconnected rural area, which implies that the impact
720 of cooperative service providers on each other and on the
721 wholesale price could not be ignored. Hence, an optimal coordi-
722 nation scheme needs to be designed to enable collaboration
723 among multiple entities. In future work, we will extend the
724 proposed RL method to address this challenge.

APPENDIX

MG OPTIMAL POWER MANAGEMENT FORMULATION

725 A moving look-ahead decision window $[t, t + T]$ is defined
726 using the latest estimations of solar and load power at dif-
727 ferent instants, where n is the MG index ($n \in \{1, \dots, N\}$),
728 i and j define the bus numbers for each MG ($\forall i, j \in \Omega_I$),
729 and k denotes the line index ($\forall k \in \Omega_K$). It has deci-
730 sion vector $\mathbf{x}_p = (P_{i,t,n}^{DG}, P_{i,t,n}^{PCC}, P_{i,t,n}^{Ch}, P_{i,t,n}^{Dis})^\top$ and $\mathbf{x}_q =$
731 $(Q_{i,t,n}^{DG}, Q_{i,t,n}^{PCC}, Q_{i,t,n}^{PV}, Q_{i,t,n}^{ESS})^\top$.

$$734 \min_{\mathbf{x}_p, \mathbf{x}_q} \sum_t^{T+t} \left(-\lambda_{i,t,n}^R P_{i,t,n}^{PCC} + \lambda_{i,t,n}^F F_{i,t,n} \right) \quad (16)$$

$$735 \text{s.t. } F_{i,t,n} = a_f \left(P_{i,t,n}^{DG} \right)^2 + b_f P_{i,t,n}^{DG} + c_f \quad (17)$$

$$736 \left| P_{i,t,n}^{PCC} \right| \leq P_{i,t,n}^{PCC,M} \quad (18)$$

$$737 \left| Q_{i,t,n}^{PCC} \right| \leq Q_{i,t,n}^{PCC,M} \quad (19)$$

$$738 0 \leq P_{i,t,n}^{DG} \leq P_{i,t,n}^{DG,M} \quad (20)$$

$$739 0 \leq Q_{i,t,n}^{DG} \leq Q_{i,t,n}^{DG,M} \quad (21)$$

$$740 \left| P_{i,t,n}^{DG} - P_{i,t-1,n}^{DG} \right| \leq P_{i,t,n}^{DG,R} \quad (22)$$

$$P_{i,t,n}^{ij} = V_{i,t,n}^i \left(V_{i,t,n}^i G_n^{ij} - V_{i,t,n}^j \left(G_n^{ij} \cos(\Delta\theta_{i,t,n}^{ij}) \right. \right. \quad 741 \\ \left. \left. + B_n^{ij} \sin(\Delta\theta_{i,t,n}^{ij}) \right) \right) \quad (23) \quad 742$$

$$Q_{i,t,n}^{ij} = -V_{i,t,n}^i \left(V_{i,t,n}^i B_n^{ij} + V_{i,t,n}^j \left(G_n^{ij} \cos(\Delta\theta_{i,t,n}^{ij}) \right. \right. \quad 743 \\ \left. \left. - B_n^{ij} \sin(\Delta\theta_{i,t,n}^{ij}) \right) \right) \quad (24) \quad 744$$

$$\left(P_{i,t,n}^{ij} \right)^2 + \left(Q_{i,t,n}^{ij} \right)^2 \leq \left(L_{i,t,n}^{ij,M} \right)^2 \quad (25) \quad 745$$

$$\sum_{ij \in k} P_{i,t,n}^{ij} = \sum_{j \in k} P_{i,t,n}^{ji} - p_{i,t,n} \quad (26) \quad 746$$

$$\sum_{i,j \in k} Q_{i,t,n}^{ij} = \sum_{j \in k} Q_{i,t,n}^{ji} - q_{i,t,n} \quad (27) \quad 747$$

$$p_{i,t,n} = P_{i,t,n}^{D,e} - P_{i,t,n}^{DG} - P_{i,t,n}^{PV,e} + P_{i,t,n}^{Ch} - P_{i,t,n}^{Dis} \quad (28) \quad 748$$

$$P_{i,t,n}^D = P_{i,t,n}^{D,e} - \varepsilon_{i,t,n}^D \quad (29) \quad 749$$

$$P_{i,t,n}^{PV} = P_{i,t,n}^{PV,e} - \varepsilon_{i,t,n}^{PV} \quad (30) \quad 750$$

$$q_{i,t,n} = Q_{i,t,n}^D - Q_{i,t,n}^{DG} - Q_{i,t,n}^{PV} + Q_{i,t,n}^{ESS} \quad (31) \quad 751$$

$$V_{i,t,n}^{PCC} = V_{i,t,n}^{PCC,E} \quad (32) \quad 752$$

$$V_{i,t,n}^m \leq V_{i,t,n} \leq V_{i,t,n}^M \quad (33) \quad 753$$

$$\left| Q_{i,t,n}^{PV} \right| \leq Q_{i,t,n}^{PV,M} \quad (34) \quad 754$$

$$SOC_{i,t,n} = SOC_{i,t-1,n} \quad 755$$

$$+ \Delta t \left(P_{i,t,n}^{Ch} \eta_{Ch} - P_{i,t,n}^{Dis} / \eta_{Dis} \right) / E_{i,t,n}^{Cap} \quad (35) \quad 756$$

$$SOC_{i,t,n}^m \leq SOC_{i,t,n} \leq SOC_{i,t,n}^M \quad (36) \quad 757$$

$$0 \leq P_{i,t,n}^{Ch} \leq u_{i,t,n}^{Ch} P_{i,t,n}^{Ch,M} \quad (37) \quad 758$$

$$0 \leq P_{i,t,n}^{Dis} \leq u_{i,t,n}^{Dis} P_{i,t,n}^{Dis,M} \quad (38) \quad 759$$

$$0 \leq u_{i,t,n}^{Ch} + u_{i,t,n}^{Dis} \leq 1 \quad (39) \quad 760$$

$$u_{i,t,n}^{Ch}, u_{i,t,n}^{Dis} \in \{0, 1\} \quad (40) \quad 761$$

762 The objective function (16) minimizes each MG's total
763 cost of operation, which is composed of two terms: the neg-
764 ative of revenue from power transfer with the cooperative
765 agent and the cost of running local DGs. Here, $\lambda_{i,t,n}^F$ is the
766 diesel generator fuel price in $\$/L$ adopted from [32]. The
767 fuel consumption $F_{i,t,n}$ of diesel generator can be expressed
768 as a quadratic polynomial function (17), with coefficients
769 $a_f = 0.0001773 L/kW^2$, $b_f = 0.1709 L/kW$, and $c_f = 14.67L$
770 adopted from [33]. Constraints (18)-(19) describe the power
771 exchange limit between the MG and the upstream distribution
772 grid with the maximum active/reactive power exchange lim-
773 its, $P_{i,t,n}^{PCC,M}$, $Q_{i,t,n}^{PCC,M}$. Constraints (20)-(21) ensure that the DG
774 active/reactive power outputs, $P_{i,t,n}^{DG}$, $Q_{i,t,n}^{DG}$, are within the DG
775 power capacity $P_{i,t,n}^{DG,M}$, $Q_{i,t,n}^{DG,M}$, and (22) enforces the maxi-
776 mum DG ramp limit, $P_{i,t,n}^{DG,R}$. Internal AC power flow model
777 of the MG is considered here with the network topology con-
778 straints, with (23) and (24) determining the active and reactive
779 power flows of each branch, where G^{ij} and B^{ij} are the cor-
780 responding real and imaginary parts of the bus admittance
781 matrix, and $V_{i,t,n}^i$ and $\Delta\theta_{i,t,n}^{ij}$ are the nodal voltage magnitude and
782 phase angle difference, respectively. Constraint (25) denotes
783 the power flow limits for each branch. Equations (26)-(31)
784 are the nodal active/reactive power balances at MG buses.

The difference between the predicted and actual PV/load values are modeled using Gaussian error variables as shown in equations (29) and (30), where $P_{i,t,n}^{D,e}$ denotes the estimated active load, and $P_{i,t,n}^{PV,e}$ is the estimated active power output of PV. Also, $\varepsilon_{i,t,n}^D, \varepsilon_{i,t,n}^{PV} \sim N(0, \sigma)$ denote the Gaussian estimation errors for active load and PV power, respectively. Constraint (32) sets the voltage at the PCC of the MG according to the estimated input voltage, $V_{i,t,n}^{PCC,E}$. Constraint (33) sets the limits for nodal bus voltage amplitude, $[V_{i,n}^m, V_{i,n}^M]$. PV reactive power output, $Q_{i,n}^{PV}$, is constrained by its maximum limit $Q_{i,n}^{PV,M}$ in (34). Operational ESS constraints are described by (35)-(40). Adopted from [34], constraint (35) determines the state of charge (SOC) of ESSs, $SOC_{i,t,n}$. The SOC and charging/discharging power of ESS, $P_{i,t,n}^{Ch}, P_{i,t,n}^{Dis}$, are constrained in (36)-(40). Here, $[SOC_{i,n}^m, SOC_{i,n}^M], P_{i,n}^{Ch,M}$ and $P_{i,n}^{Dis,M}$ define the permissible range of SOC, and maximum charging and discharging power, with $u_{i,t,n}^{Ch}$ and $u_{i,t,n}^{Dis}$ denoting the charge/discharge binary indicator variables, and η_{Ch}/η_{Dis} representing the charging/discharging efficiency. $E_{i,n}^{Cap}$ denotes the maximum capacity of ESSs.

REFERENCES

[1] S. A. Arefifar, Y. A.-R. I. Mohamed, and T. H. M. EL-Fouly, "Optimum microgrid design for enhancing reliability and supply-security," *IEEE Trans. Smart Grid*, vol. 4, no. 3, pp. 1567–1575, Sep. 2013.

[2] K. Ubilla *et al.*, "Smart microgrids as a solution for rural electrification: Ensuring long-term sustainability through cadastre and business models," *IEEE Trans. Sustain. Energy*, vol. 5, no. 4, pp. 1310–1318, Oct. 2014.

[3] P. Chen, "A roadmap to the new rural electric cooperative business model," M.S. thesis, Nicholas School Environ., Duke Univ., Durham, NC, USA, 2017.

[4] *Understanding Value is Critical to Microgrid Projects*. [Online]. Available: <https://www.ncelectriccooperatives.com/who-we-are/spotlight/understanding-value-is-critical-to-microgrid-projects/>

[5] R. Lasseter *et al.*, *Integration of Distributed Energy Resources. The CERTS Microgrid Concept*. Berkeley, CA, USA: Lawrence Berkeley Nat. Lab., 2002.

[6] M. Manbachi and M. Ordenez, "AMI-based energy management of islanded AC/DC microgrids utilizing energy conversion and optimization," *IEEE Trans. Smart Grid*, to be published.

[7] M. Nemati, M. Braun, and S. Tenbohlen, "Optimization of unit commitment and economic dispatch in microgrids based on genetic algorithm and mixed integer linear programming," *Appl. Energy*, vol. 210, pp. 944–963, Jan. 2018.

[8] N. Nikmehr and S. N. Ravadanegh, "Optimal power dispatch of multi-microgrids at future smart distribution grids," *IEEE Trans. Smart Grid*, vol. 6, no. 4, pp. 1949–1957, Jul. 2015.

[9] S. A. Arefifar, M. Ordenez, and Y. A. I. Mohamed, "Energy management in multi-microgrid systems—Development and assessment," *IEEE Trans. Power Syst.*, vol. 32, no. 2, pp. 910–952, Mar. 2017.

[10] A. Ouammi, H. Dagdougui, and R. Sacile, "Optimal control of power flows and energy local storages in a network of microgrids modeled as a system of systems," *IEEE Trans. Control Syst. Tech.*, vol. 62, no. 4, pp. 2551–2559, Apr. 2015.

[11] Y. Zhang, L. Xie, and Q. Ding, "Interactive control of coupled microgrids for guaranteed system-wide small signal stability," *IEEE Trans. Smart Grid*, vol. 7, no. 2, pp. 1088–1096, Mar. 2016.

[12] Z. Wang, B. Chen, J. Wang, and J. Kim, "Decentralized energy management system for networked microgrids in grid-connected and islanded modes," *IEEE Trans. Smart Grid*, vol. 7, no. 2, pp. 1097–1105, Mar. 2016.

[13] Z. Wang, B. Chen, J. Wang, M. Begovic, and C. Chen, "Coordinated energy management of networked microgrids in distribution systems," *IEEE Trans. Smart Grid*, vol. 6, no. 1, pp. 45–53, Jan. 2015.

[14] P. Kou, D. Liang, and L. Gao, "Distributed EMPC of multiple microgrids for coordinated stochastic energy management," *Appl. Energy*, vol. 185, pp. 939–952, Jan. 2017.

[15] D. Gregoratti and J. Matamoros, "Distributed energy trading: The multiple-microgrid case," *IEEE Trans. Ind. Electron.*, vol. 62, no. 4, pp. 2551–2559, Apr. 2015.

[16] W. Liu, P. Zhuang, H. Liang, J. Peng, and Z. Huang, "Distributed economic dispatch in microgrids based on cooperative reinforcement learning," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 29, no. 2, pp. 2192–2203, Jun. 2018.

[17] E. Mocanu *et al.*, "On-line building energy optimization using deep reinforcement learning," *IEEE Trans. Smart Grid*, vol. 10, no. 4, pp. 3698–3708, Jul. 2019.

[18] *Guidelines for Smart Grid Cybersecurity Volume 2: Privacy and the Smart Grid*. [Online]. Available: <https://nvlpubs.nist.gov/nistpubs/ir/2014/NIST.IR.7628r1.pdf>

[19] G. Mokryani, "Active distribution networks planning with integration of demand response," *Solar Energy*, vol. 122, pp. 1362–1370, Dec. 2015.

[20] M. Wabbah, T. H. M. El-Fouly, B. Zahawi, and S. Feng, "Hybrid beta-KDE model for solar irradiance probability density estimation," *IEEE Trans. Sustain. Energy*, to be published.

[21] H. Wu, M. Shahidehpour, Z. Li, and W. Tian, "Chance-constrained day-ahead scheduling in stochastic power system operation," *IEEE Trans. Power Syst.*, vol. 29, no. 4, pp. 1583–1591, Jul. 2014.

[22] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. London, U.K.: MIT Press, 2017.

[23] F. L. Lewis and K. G. Vamvoudakis, "Reinforcement learning for partially observable dynamic processes: Adaptive dynamic programming using measured output data," *IEEE Trans. Syst., Man, Cybern. B, Cybern.*, vol. 41, no. 1, pp. 14–25, Feb. 2011.

[24] L. Busoniu, R. Babuska, B. D. Schutter, and D. Ernst, *Reinforcement Learning and Dynamic Programming Using Function Approximators*. Boca Raton, FL, USA: CRC Press, 2017.

[25] W. D. Smart and L. P. Kaelbling, "Practical reinforcement learning in continuous spaces," in *Proc. Int. Conf. Mach. Learn.*, Jun. 2000, pp. 903–910.

[26] I. Houtzager, J. W. van Wingerden, and M. Verhaegen, "Recursive predictor-based subspace identification with application to the real-time closed-loop tracking of flutter," *IEEE Trans. Control Syst. Technol.*, vol. 20, no. 4, pp. 934–949, Jul. 2012.

[27] M. E. Baran and F. F. Wu, "Network reconfiguration in distribution systems for loss reduction and load balancing," *IEEE Trans. Power Del.*, vol. 4, no. 2, pp. 1401–1407, Apr. 1989.

[28] R. S. Rao, K. Ravindra, K. Satish, and S. V. L. Narasimham, "Power loss minimization in distribution system using network reconfiguration in the presence of distributed generation," *IEEE Trans. Power Syst.*, vol. 28, no. 1, pp. 317–325, Feb. 2013.

[29] W. H. Kersting, "Radial distribution test feeder," in *Proc. IEEE Power Eng. Soc. Winter Meeting*, vol. 2, 2001, pp. 908–912.

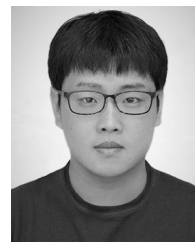
[30] *Wholesale Electricity and Natural Gas Market Data*. [Online]. Available: <https://www.eia.gov/electricity/wholesale/#history>

[31] B. Anthony, D. Kendrick, and A. Meeraus, "GAMS, a user's guide," *ACM SIGNUM Newsl.*, vol. 23, nos. 3–4, pp. 10–11, 1988.

[32] *Gasoline and Diesel Fuel Update*. [Online]. Available: <https://www.eia.gov/petroleum/gasdiesel/>

[33] S. A. Pourmousavi, M. H. Nehrir, and R. K. Sharma, "Multi-timescale power management for islanded microgrids including storage and demand response," *IEEE Trans. Smart Grid*, vol. 6, no. 3, pp. 1185–1195, May 2015.

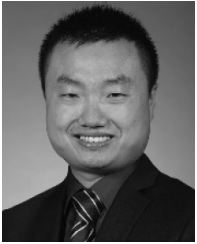
[34] Y. Xuan, J. Zhang, and Q. Zhang, "Combinational linear programming approach for daily optimal operation of customers with DG/ESS under TOU pricing for customer EMS software applications," in *Proc. IEEE Innov. Smart Grid Technol. Asia (ISGT Asia)*, 2018, pp. 300–305.



Qianzhi Zhang (S'17) received the B.S. degree in electrical and computer engineering from the Shandong University of Technology in 2012, and the M.S. degree in electrical and computer engineering from Arizona State University in 2015. He is currently pursuing the Ph.D. degree with the Department of Electrical and Computer Engineering, Iowa State University, Ames, IA, USA. From 2015 to 2016, he was a Research Engineer with Huadian Electric Power Research Institute, Hangzhou, China. His research interests include power distribution systems, microgrids, applications of distributed optimization, and machine learning in power systems.

925
926
927
928
929
930
931
932
933
934

Kaveh Dehghanpour (S'14–M'17) received the B.Sc. and M.S. degrees in electrical and computer engineering from the University of Tehran in 2011 and 2013, respectively, and the Ph.D. degree in electrical engineering from Montana State University in 2017. He is currently a Post-Doctoral Research Associate with Iowa State University. His research interest includes application of machine learning and data-driven techniques in power system monitoring and control.

935
936
937
938
939
940
941
942
943
944
945

Zhaoyu Wang (S'13–M'15) received the B.S. and M.S. degrees in electrical engineering from Shanghai Jiao Tong University in 2009 and 2012, respectively, and the M.S. and Ph.D. degrees in electrical and computer engineering from the Georgia Institute of Technology in 2012 and 2015, respectively. He is the Harpole-Pentair Assistant Professor with Iowa State University. He was a Research Aid with Argonne National Laboratory in 2013 and an Electrical Engineer Intern with Corning Inc., in 2014. His research interests include power distribu-

tion systems, microgrids, renewable integration, power system resilience, and data-driven system modeling. He is the Principal Investigator for a multitude of projects focused on these topics and funded by the National Science Foundation, the Department of Energy, National Laboratories, PSERC, and Iowa Energy Center. He is the Secretary of the IEEE Power and Energy Society Award Subcommittee. He is an Editor of the IEEE TRANSACTIONS ON POWER SYSTEMS, the IEEE TRANSACTIONS ON SMART GRID, and IEEE PES Letters, and an Associate Editor of *IET Smart Grid*.



Qihua Huang (S'14–M'16) received the B.E. and M.Sc. degrees from the South China University of Technology in 2012 and 2009, respectively, and the Ph.D. degree in electrical engineering from Arizona State University in 2016. He is currently a Senior Power System Research Engineer with the Electricity Infrastructure Group, Pacific Northwest National Laboratory, USA. He has been a Principal Investigator, a Project Manager, and task lead in several DOE funded projects. His research interests include power system modeling, simulation and control, and application of advanced computing, and machine learning technologies in power systems. He was a recipient of the IEEE PES Prize Paper Award in 2019 and the R&D 100 Award in 2018. He is an Associate Editor of IEEE ACCESS and the *CSEE Journal of Power and Energy Systems*.

954
955
956
957
958
959
960
961
962
963
964
965
966
967
968

IEEE PROOF