

Multi-Agent Safe Policy Learning for Power Management of Networked Microgrids

Qianzhi Zhang¹, Graduate Student Member, IEEE, Kaveh Dehghanpour¹, Member, IEEE,
Zhaoyu Wang¹, Senior Member, IEEE, Feng Qiu¹, Senior Member, IEEE,
and Dongbo Zhao¹, Senior Member, IEEE

AQ1

Abstract—This article presents a supervised multi-agent safe policy learning (SMAS-PL) method for optimal power management of networked microgrids (MGs) in distribution systems. While unconstrained reinforcement learning (RL) algorithms are black-box decision models that could fail to satisfy grid operational constraints, our proposed method considers AC power flow equations and other operational limits. Accordingly, the training process employs the gradient information of operational constraints to ensure that the optimal control policy functions generate safe and feasible decisions. Furthermore, we have developed a distributed consensus-based optimization approach to train the agents' policy functions while maintaining MGs' privacy and data ownership boundaries. After training, the learned optimal policy functions can be safely used by the MGs to dispatch their local resources, without the need to solve a complex optimization problem from scratch. Numerical experiments have been devised to verify the performance of the proposed method.

Index Terms—Safe policy learning, multi-agent framework, networked microgrids, power management, policy gradient.

NOMENCLATURE

Indices

i, j	Indices of buses, $\forall i, j \in \Omega_I$.
ij	Index of branch between bus i and bus j , $\forall ij \in \Omega_{Br}$.
k	Iteration index in distributed optimization, $k \in \{1, \dots, k^{max}\}$.
m	Constraint index, $m \in \{1, \dots, M_c\}$.
n	Agent index, $n \in \{1, \dots, N\}$.
t'	Episode index in training process, $t' \in [t, t+T]$.

Parameters

a_n^f, b_n^f, c_n^f	Coefficients of the DG quadratic cost function for agent n .
-----------------------	--

b_{m, μ_n}	Gradient vector of the constraint return function m w.r.t. the parameters μ_n .	34
b_{m, Σ_n}	Gradient vectors of the constraint return function m w.r.t. the parameters Σ_n .	35
D_n	Dimension of multivariate Gaussian distribution function for agent n .	36
d_m	Upper limit for constraint m .	37
E^{Cap}	Max. capacity of ESS unit.	38
H_n	Fisher information matrix of agent n .	39
g_{μ_n}	Gradient vector of the reward functions w.r.t. the parameters μ_n .	40
g_{Σ_n}	Gradient vector of the reward functions w.r.t. the parameters Σ_n .	41
I_{ij}^M	Max. current limit on branch ij .	42
M_c	Number of constraints.	43
M_c^G	Number of global constraints.	44
M_c^L	Number of local constraints.	45
N	Number of MGs.	46
N_n	Number of neighboring MGs for agent n .	47
$p^{Ch, M}$	Max. ESS charging limits.	48
$p^{Dis, M}$	Max. ESS discharging limits.	49
P^D, Q^D	Active and reactive load power.	50
$p^{DG, M}$	Max. DG active power capacity.	51
$Q^{DG, M}$	Max. DG reactive power capacity.	52
$p^{DG, R}$	Max. DG ramp limit.	53
p^{PV}	PV active power output.	54
$p^{PCC, M}$	Max. active power flow at the PCCs.	55
$Q^{PCC, M}$	Max. reactive power flow at the PCCs.	56
$Q^{PV, M}$	Max. PV reactive power output limit.	57
SOC^M	Max. SOC limits.	58
SOC^m	Min. SOC limits.	59
T	Length of the moving decision window.	60
V_i^M, V_i^m	Max. and min. voltage limit on bus i .	61
$w_n(n')$	Weight parameters assigned of agent n to neighboring agent n' .	62
γ^{Re}, γ^{Im}	Real and imaginary parts of the nodal admittance matrix Y .	63
η^{Ch}, η^{Dis}	Charging and discharging efficiency of ESS.	64
λ^F	Diesel generator fuel price.	65
λ^R	Retail price signals at the PCCs.	66
$\theta_{\mu_n}, \theta_{\Sigma_n}$	Vector of DNN weights and bias of agent n .	67
μ_n, Σ_n	Mean vector and covariance matrices for control action of agent n .	68
δ, ρ_1	Step sizes for updating θ and λ .	69

Manuscript received October 21, 2019; revised February 4, 2020, June 6, 2020, and September 12, 2020; accepted October 26, 2020. This work was supported in part by the U.S. Department of Energy Wind Energy Technologies Office under Grant DE-EE0008956, and in part by the National Science Foundation Grant ECCS 1929975. Paper no. TSG-01584-2019. (Corresponding author: Zhaoyu Wang.)

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

Feng Qiu and Dongbo Zhao are with the Energy Systems Division, Argonne National Laboratory, Lemont, IL 60439 USA (e-mail: fqiu@anl.gov; dongbo.zhao@anl.gov).

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

Digital Object Identifier 10.1109/TSG.2020.3034827

AQ2

78	ρ_2	Penalty factor for constraints violation.
79	Δt	Time step.
80	$\Delta\theta_n$	Threshold for parameter updating.
81	γ	Discount factor.
82	τ	Tightening multiplier.
83	<i>Variables</i>	
84	\mathbf{a}_n	Vector of control actions of agent n .
85	$C_m(\pi)$	Return value of constraint m based on the control policy π .
86		
87	$F_{i,n}$	Fuel consumption of DG at bus i of agent n .
88	I_i^{Re}, I_i^{Im}	Real and imaginary parts of the injected current at bus i .
89		
90	I_{ij}^{Re}, I_{ij}^{Im}	Real and imaginary parts of the branch current at branch ij .
91		
92	\mathbf{O}_t	Vectors of observation variable.
93	p^{Ch}, p^{Dis}	Charging and discharging power of ESS unit.
94	p^{DG}, Q^{DG}	DG active and reactive power outputs
95	p^{PCC}	Active power flow at the PCC.
96	Q^{PCC}	Reactive power flow at the PCC.
97	Q^{ESS}	Reactive power outputs of ESS unit.
98	Q^{PV}	PV inverter reactive power output.
99	SOC	SOC of the battery system.
100	\mathbf{S}_n	Vectors of system state of agent n .
101	V_i^{Re}, V_i^{Im}	Real and imaginary parts of the bus voltage magnitude at bus i .
102		
103	λ_n	Vector of Lagrangian multipliers.
104	<i>Functions</i>	
105	J_{R_n}	Expected reward function of agent n .
106	J_{C_m}	Expected return function of constraint m .
107	π_n	Multivariate distribution function over control actions of agent n .
108		
109	Δ	Kullback Leibler (KL)-divergence function.

I. INTRODUCTION

MICROGRIDS (MGs) are active clusters of distributed energy resources (DERs), loads, energy storage system (ESS), and other onsite electric components. A smart distribution system may consist of multiple MGs and the coordinated control of the networked MGs can offer various benefits, including higher penetration of local DERs, improved controllability, and enhancement of power system resilience and reliability [1], [2]. Solving the power management problem of networked MGs is a complex task. While previous works in this area have provided valuable insight, we have identified two shortcomings in the literature:

(1) *Limitations of model-based optimization methods*: In the existing literature, there are quite a few model-based methods for solving the optimal power management problem of networked MGs, such as centralized decision models [3]–[5] and distributed control frameworks [6]–[8]. However, with increasing number of MGs in distribution networks, these methods have to solve large-scale optimization problems with numerous nonlinear constraints that incur high computational costs and hinder real-time decision making. Furthermore, model-based methods are unable to adapt to the continuously evolving

system conditions, as they need to re-solve the problem at each time step.

(2) *Potential infeasibility of model-free machine learning methods*: To address the limitations of model-based methods, model-free reinforcement learning (RL) techniques have been used to solve the optimal power management problem through repeated interactions between a control agent and its environment. This approach eliminates the need to solve a large-scale optimization problem at each time point and enables the control agent to provide adaptive response to time-varying system states. Existing examples of RL application in power systems include economic dispatch and energy consumption scheduling of individual MGs [9]–[11] and multi-area smart control of generation in interconnected power grids [12], [13]. Further, in our previous paper [14], we have proposed a bi-level power management method for networked MGs, where a centralized RL agent determines retail prices in a cooperative business model for each MG under the incomplete information of physical model. Current RL-based solutions employ control agents to train *black-box* functions to approximate the optimal actions through trial and error. However, the trained black-box functions can fail to satisfy critical operational constraints, such as network nodal voltage and capacity limits, since these constraints have not been encoded in the training process. This can lead to unsafe operational states and control action infeasibility.

However, incorporating constraints into the training process of conventional black-box methods is challenging since these methods have generally relied on adding penalty terms to training objective functions for enforcing constraints, which cannot guarantee the safety of control policies as the number of constraints grows. Inspired by recent advances in constrained policy learning (PL) [15]–[17] and to address the shortcomings in the existing literature, we have cast the power management of networked MGs as a *supervised multi-agent safe PL problem* (SMAS-PL). The various resources inside each MG and the collaborative behavior of MGs are both controlled to optimize the total cost of operation, while satisfying all the local and global constraints. Moreover, we have proposed a multi-agent policy gradient solution strategy, which enables individual MGs learn control policy functions to maximize the social welfare and ensure safety in a distributed way. The proposed method introduces a trade-off between model-free and model-based methods and combines the benefits offered by both sides. The purpose is to leverage the advantages of both model-free and model-based methods, for scalable real-time decision making while also maintaining a user-defined level of safety by considering constraints in the training process. Hence, on one hand, MGs' power management policy functions are modeled using black-box deep neural networks (DNNs); while on the other hand, to ensure decision feasibility, a constrained gradient-based training method is proposed that exploits the derivatives of the constraints and objective functions of the power management problem w.r.t. control actions and learning parameters. The training process employs these gradient factors to provide a convex quadratically constrained linear program (QCLP) approximation to the power management problem at each

190 episode. This enables the proposed method to be both adapt-
 191 able to changes in the inputs of the black-box components, and
 192 feasible with respect to operational constraints, including AC
 193 power flow. Finally, a distributed consensus-based primal-dual
 194 optimization method [18] is adopted to decompose the training
 195 task among MG agents. In summary, compared to existing
 196 decision making solutions, the main advantages of this article
 197 are as follows:

- 198 • Compared to the black-box learning-based methods, the
 199 proposed SMAS-PL leverages the gradient information
 200 of all the operational constraints to devise a tractable
 201 QCLP-based training process to promote the safety and
 202 feasibility of control policies. A backtracking mechanism
 203 is added into the PL framework to perform a final verification
 204 of feasibility before issuing control commands to
 205 the assets.
- 206 • Compared to conventional centralized training methods,
 207 the distributed training process in the SMAS-PL offers
 208 two advantages: it preserves the privacy of MG agents,
 209 including their control policies parameters and structures,
 210 operation cost functions, and local asset constraints;
 211 it also enhances computational efficiency and maintains
 212 scalability as the number of learning parameters grows
 213 into a humongous size.
- 214 • The proposed SMAS-PL method does not need to solve a
 215 complex optimization problem in real-time. The agents'
 216 policy functions, that are trained offline, can be leveraged
 217 online to select optimal control actions in response to
 218 latest system state data.

219 The reminder of this article is organized as follows.
 220 Section II presents the overall framework of the proposed solution.
 221 Section III introduces the SMAS-PL problem and integrates
 222 problem gradients into the solver. Section IV describes
 223 the multi-agent consensus-based training algorithm for SMAS-
 224 PL. Simulation results and conclusions are given in Section V
 225 and Section VI, respectively.

226 II. OVERVIEW OF THE PROPOSED FRAMEWORK

227 The general framework of the proposed SMAS-PL method
 228 is shown in Fig. 1. Note that vectors are denoted in bold letters
 229 throughout this article. The micro-sources within each MG
 230 are controlled by an agent that adopts a private control policy.
 231 Here, the *control policy* for the n 'th agent, π_n , is a parametric
 232 probability distribution function, with parameters θ_n , over the
 233 agent's control actions ($\mathbf{a}_{n,t}$), including active/reactive power
 234 dispatching signals for local diesel generators (DGs), ESS and
 235 solar photo-voltaic (PV) panels. Note that the control policy
 236 π_n is a function of the MG's state variables ($\mathbf{S}_{n,t}$), defined
 237 by the aggregate MG load and solar irradiance. To ensure the
 238 safety of the control policies, MG agents receive the observed
 239 variables from the grid, including network nodal voltages \mathbf{V}_t
 240 and injection currents \mathbf{I}_t , to determine gradient factors of the
 241 problem constraints and objectives w.r.t. to learning parameters,
 242 $\nabla_{\theta} J$. These gradient factors are then integrated into a
 243 multi-agent constrained training algorithm, which employs
 244 local inter-MG communication to satisfy all global and local
 245 operational constraints through exchanging and processing

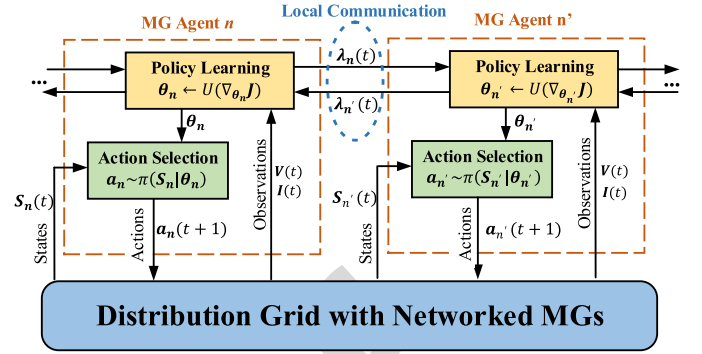


Fig. 1. Structure of the proposed SMAS-PL method for power management of networked MGs.

dual Lagrangian variables, $\lambda(t)$. The Lagrangian multipliers
 246 embody the interactions among the MGs and capture the
 247 impacts of MGs' decisions on each other. Theoretical analysis
 248 and numerical simulations are conducted to show that the
 249 proposed SMAS-PL method can minimize the MG agents'
 250 operational cost and satisfy operational constraints. Note that
 251 the proposed SMAS-PL is not a purely model-free approach,
 252 since the AC power flow equations are used to calculate gradient
 253 factors and ensure the decision feasibility when training the
 254 DNNs.
 255

256 In this article, the MGs are chosen to be collaborative,
 257 because the satisfaction of the global constraints (i.e., limits
 258 on nodal voltages and line flows) for the whole network
 259 needs coordination among all MGs. Since global constraints
 260 are impacted by the response of all the MGs, we have devised
 261 a collaborative policy learning to ensure that grid-wide operation
 262 remains safe. Specifically, the consensus-based training method
 263 leverages the Lagrange multipliers of the global constraints to
 264 coordinate the policy optimization of the MGs. Thus, each
 265 Lagrange multiplier serves as a penalty factor or a shadow price,
 266 which enforces safety in the data-driven procedure.
 267

268 III. SAFE POLICY LEARNING FOR POWER MANAGEMENT 269 OF NETWORKED MGs

270 To facilitate the discussion, Section III-A introduces a general
 271 power management formulation that is commonly used in
 272 literature [4], [6], [14]. Section III-B defines each component
 273 of the proposed SMAS-PL. In Sections III-C and III-D, we
 274 propose a tractable SMAS-PL method, employing the gradient
 275 factors of reward function and constraint return functions
 276 w.r.t. actions and learning parameters, to solve the power
 277 management of networked MGs.

278 A. Power Management Problem Statement

279 Each MG is assumed to have local DGs, ESS, solar PV
 280 panels and a number of loads. This optimization problem is
 281 solved over a moving look-ahead decision window $t' \in [t, t + T]$,
 282 using the latest estimations of solar and load power at
 283 different instants. Here, n is the MG index ($n \in \{1, \dots, N\}$),
 284 i and j define the node numbers ($\forall i, j \in \Omega_i$), ij defines the
 285 branch numbers ($\forall ij \in \Omega_{Br}$).
 286

286 1) *Problem Objective*: The objective function (1), with
 287 control action vector $[P^{DG}, P^{Ch}, P^{Dis}, Q^{DG}, Q^{PV}, Q^{ESS}] \in$
 288 $(\mathbf{x}_p, \mathbf{x}_q)$, minimizes MGs' total cost of operation, which is
 289 composed of the income/cost from power transfer with the
 290 grid and cost of running local DG. Here, λ_n^F is the DG fuel
 291 price, λ_n is the electricity price, and $P_{n,t'}^{PCC}$ is active power
 292 transfer between grid and the n 'th MG at the point of com-
 293 mon coupling (PCC). The fuel consumption of DG, $F_{i,n,t'}$, can
 294 be expressed as a quadratic polynomial function of its power,
 295 $P_{i,n,t'}^{DG}$, with parameters a_n^f , b_n^f and c_n^f .

$$296 \quad \min_{\mathbf{x}_p, \mathbf{x}_q} \sum_{n=1}^N \sum_{t'=t}^{t+T} \left(-\lambda_n^R P_{n,t'}^{PCC} + \lambda_{i,n}^F F_{i,n,t'} \right) \quad (1)$$

$$297 \quad F_{i,n,t'} = a_n^f \left(P_{i,n,t'}^{DG} \right)^2 + b_n^f P_{i,n,t'}^{DG} + c_n^f. \quad (2)$$

298 2) *Global Constraints*: These constraints are defined over
 299 variables that are impacted by control actions of all the MGs,
 300 including the voltage amplitude limits for the entire nodes,
 301 $[V_i^m, V_i^M]$, and the maximum permissible branch current flow
 302 magnitudes I_{ij}^M throughout the distribution grid and the MGs:

$$303 \quad V_i^m \leq V_{i,t'} \leq V_i^M \quad (3)$$

$$304 \quad -I_{ij}^M \leq I_{ij,t'} \leq I_{ij}^M \quad (4)$$

305 The global constraints (3)-(4) are implicitly determined by
 306 the AC power flow equations, which will be used to calculate
 307 the gradient factors of objective (1) and constraints (3)-(16)
 308 w.r.t. learning parameters as elaborated in Section III-D. Note
 309 that unlike previous centralized optimization solutions that are
 310 generally model-based, our strategy is a combination of both
 311 model-based and model-free approaches. Thus, while power
 312 flow equations appear explicitly in centralized optimization
 313 models, our solution only leverage power flow equations in an
 314 implicit way in the training process to ensure that the learning
 315 modules are generating feasible outcomes.

316 3) *Local Constraints*: These constraints are defined over the
 317 local control actions of each MG. Constraints (5)-(6) ensure
 318 that the DG active/reactive power outputs, $P_{i,n}^{DG}/Q_{i,n}^{DG}$, are
 319 within the DG power capacity $P_{i,n}^{DG,M}/Q_{i,n}^{DG,M}$, and (7) enforces
 320 the maximum DG ramp limit, $P_{i,n}^{DG,R}$. PV reactive power out-
 321 put, $Q_{i,n}^{PV}$, is constrained by its maximum limit $Q_{i,n}^{PV,M}$ per (8).
 322 The active power transfer $P_{n,t'}^{PCC}$ and the reactive power trans-
 323 fer $Q_{n,t'}^{PCC}$ at the PCCs are bounded with the constraints (9)
 324 and (10), respectively.

$$325 \quad 0 \leq P_{i,n,t'}^{DG} \leq P_{i,n}^{DG,M} \quad (5)$$

$$326 \quad 0 \leq Q_{i,n,t'}^{DG} \leq Q_{i,n}^{DG,M} \quad (6)$$

$$327 \quad \left| P_{i,n,t'}^{DG} - P_{i,n,t'-1}^{DG} \right| \leq P_{i,n}^{DG,R} \quad (7)$$

$$328 \quad \left| Q_{i,n,t'}^{PV} \right| \leq Q_{i,n}^{PV,M} \quad (8)$$

$$329 \quad \left| P_{n,t'}^{PCC} \right| \leq P_{n,t'}^{PCC,M} \quad (9)$$

$$330 \quad \left| Q_{n,t'}^{PCC} \right| \leq Q_{n,t'}^{PCC,M} \quad (10)$$

331 The operational ESS constraints are described by (11)-(16),
 332 where (11) determines the state of charge (SOC) of ESSs,

$SOC_{i,n}$. $E_{i,n}^{Cap}$ denotes the maximum capacity of ESSs. To
 ensure safe ESS operation, the SOC and charging/discharging
 power of ESS, $P_{i,n}^{Ch}$, $P_{i,n}^{Dis}$, are constrained as shown in
 (12)-(16). 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, respectively. Constraint (15) indicates that
 ESSs cannot charge and discharge at the same time instant.
 And η_{Ch}/η_{Dis} represents the charging/discharging efficiency.
 The reactive power of ESS, $Q_{i,n}^{ESS}$, is kept within maximum
 limit, $Q_{i,n}^{ESS,M}$, through constraint (16).

$$343 \quad SOC_{i,n,t'} = SOC_{i,n,t'-1} + \Delta t \frac{\left(P_{i,n,t'}^{Ch} \eta_{Ch} - P_{i,n,t'}^{Dis} / \eta_{Dis} \right)}{E_{i,n}^{Cap}} \quad (11)$$

$$344 \quad SOC_{i,n}^m \leq SOC_{i,n,t'} \leq SOC_{i,n}^M \quad (12)$$

$$345 \quad 0 \leq P_{i,n,t'}^{Ch} \leq P_{i,n}^{Ch,M} \quad (13)$$

$$346 \quad 0 \leq P_{i,n,t'}^{Dis} \leq P_{i,n}^{Dis,M} \quad (14)$$

$$347 \quad P_{i,n,t'}^{Ch} P_{i,n,t'}^{Dis} = 0 \quad (15)$$

$$348 \quad \left| Q_{i,n,t'}^{ESS} \right| \leq Q_{i,n}^{ESS,M} \quad (16)$$

349 Note that the distribution system and networked MGs are
 350 operated in normal condition, which means the switch opera-
 351 tion and the network topology are assumed to be unchanged
 352 during the operation period.
 353

354 B. Safe Policy Learning Setup

355 In this section, the optimal power management of networked
 356 MGs is transformed into a SMAS-PL problem. The purpose
 357 of the SMAS-PL is to provide a framework for control agents
 358 to collaboratively find control policies to maximize their total
 359 accumulated reward while satisfying all problem constraints.
 360 To do this, we have provided formulations to ensure that the
 361 outcome of the SMAS-PL also corresponds to the solution of
 362 optimal power management of networked MGs (1)-(16). To
 363 show this, first we provide a description of the components of
 364 the SMAS-PL method:

365 1) *Control Agents*: The problem consists of N autonomous
 366 control agents, where each agent is in charge of dispatching
 367 the resources within an individual MG. The MGs are *collabo-*
 368 *rative*, in the sense that they depend on local communication
 369 with each other to optimize their behaviors.

370 2) *State Set*: The state vector for the n 'th MG agent at
 371 time t is defined as $\mathbf{S}_{n,t}$ over the time window $[t, t+T]$, as
 372 $\mathbf{S}_{n,t} = [\hat{\mathbf{I}}_{n,t'}^{PV}, \hat{\mathbf{P}}_{n,t'}^D]_{t'=t}^{t+T}$, where $\hat{\mathbf{I}}_{n,t'}^{PV}$ and $\hat{\mathbf{P}}_{n,t'}^D$ are the vectors
 373 of predicted aggregate internal load power and solar irradiance
 374 of the n 'th MG at time t' , respectively. The prediction
 375 errors follow random distributions with zero mean and the
 376 standard deviations selected from the beta and Gaussian dis-
 377 tributions adopted from [19]–[21]. Note that the parameters of
 378 forecasting error distributions are different for different MG
 379 agents.

380 3) *Action Set*: The control action vector for the n 'th agent at
 381 time t is denoted as $\mathbf{a}_{n,t} \in \mathbb{R}^{D_n}$ and consists of the dispatching
 382 decision variables for the n 'th MG over the time window $[t, t+T]$,
 383 as $\mathbf{a}_{n,t} = [P_{n,t'}^{DG}, P_{n,t'}^{Ch}, P_{n,t'}^{Dis}, Q_{n,t'}^{DG}, Q_{n,t'}^{PV}, Q_{n,t'}^{ESS}]_{t'=t}^{t+T}$.

384 4) *Observation Set*: The observation variable vector for the
 385 agents at time t is denoted as \mathbf{O}_t , and includes grid's nodal
 386 voltages \mathbf{V}_t and current injections \mathbf{I}_t at that time, $\mathbf{O}_t = [\mathbf{V}_t, \mathbf{I}_t]$.
 387 Note that the observations are implicitly determined by the
 388 agents' control actions, and thus, cannot be predicted indepen-
 389 dently of the agents' policies. However, unlike the observation
 390 variables, the state variables are independent of the agents'
 391 control actions and can be predicted for the whole decision
 392 window without the need to consider agents' policies. In
 393 the power management problem, nodal sensors or distribution
 394 grid's state estimation module will provide the latest values of
 395 observations.

396 5) *Control Policy*: In this work, the control policies are
 397 modelled as multivariate Gaussian distributions due to several
 398 reasons: (i) Gaussian distributions allow for explicit learning of
 399 both expectations and uncertainties of control policies, which
 400 are directly represented by the parameters of the distribution.
 401 Most of other distributions are parameterized by unintuitive
 402 parameters that make the decision model harder to interpret
 403 and verify. (ii) The gradients of Gaussian policy functions with
 404 respect to actions and learning parameters are easy to com-
 405 pute (see Appendices A and B). (iii) Gaussian policy functions
 406 have been adopted and suggested by [22] and [23]. Thus, the
 407 control policy for the n 'th agent, denoted as π_n , is defined as
 408 a D_n -dimensional multivariate Gaussian distribution over con-
 409 trol actions \mathbf{a}_n . The policy function determines the probability
 410 of the agent's optimal control action after training, as follows:

$$411 \mathbf{a}_n \sim \pi_n(\mathbf{a}_n | \theta_n) = \frac{1}{\sqrt{|\Sigma_n|} (2\pi)^{D_n}} e^{-\frac{1}{2}(\mathbf{a}_n - \boldsymbol{\mu}_n)^\top \Sigma_n^{-1} (\mathbf{a}_n - \boldsymbol{\mu}_n)} \quad (17)$$

412 where $\boldsymbol{\mu}_n \in \mathbb{R}^{D_n \times 1}$ is the mean vector and $\Sigma_n \in \mathbb{R}^{D_n \times D_n}$ is
 413 the covariance matrix of of multivariate Gaussian distribution
 414 for the n 'th agent. The Gaussian policy function explicitly
 415 determines the expected value and uncertainties of optimal
 416 control actions for each agent. Each agent's learning parameter
 417 vector, θ_n , consists of two parametric subsets $\theta_{\boldsymbol{\mu}_n}$ and θ_{Σ_n} ,
 418 corresponding to the mean vector and the covariance matrix of
 419 the agent's policy function. To do this, two DNNs are used for
 420 each MG agent as parametric learning functions to represent
 421 control policy components. These DNNs receive the agent's
 422 states, \mathbf{S}_n , as input to fully quantify the sufficient statistics of
 423 optimal control policies of MGs, i.e., the mean vector and the
 424 covariance matrix of the agent's actions, as follows:

$$425 \boldsymbol{\mu}_n = \text{DNN}(\mathbf{S}_n | \theta_{\boldsymbol{\mu}_n}) \quad (18)$$

$$426 \Sigma_n = \text{DNN}(\mathbf{S}_n | \theta_{\Sigma_n}) \quad (19)$$

427 The DNNs are maintained, continuously updated, and
 428 deployed in real-time by local control agents of each MG. Note
 429 that the proposed SMAS-PL method introduces a trade-off
 430 between model-free and model-based methods and combines
 431 the benefits offered by both sides. Thus, the reasons for the
 432 use of DNN-based distributions for modeling actions are as
 433 follows: (i) we have leveraged the model information to train
 434 safe policy functions that guarantee feasibility (i.e., the model-
 435 based aspect of the solution); (ii) the trained policy functions
 436 are deployed online for action selection, simply by inserting
 437 the latest data samples into the DNN-based policy functions
 438 (i.e., the model-free aspect of the solution).

6) *Reward Function*: The reward function for the n 'th MG
 is defined as the discounted negative accumulated operational
 cost of individual MG over the decision window $[t, t + T]$,
 $R_{n,t} = -[\sum_{t'=t}^{t+T} (-\lambda_n^R P_{n,t'}^{PCC} + \lambda_{i,n}^F F_{i,n,t'})]$, obtained from the
 objective functions of the networked MGs power management
 problem, (1), as follows:

$$J_{R_n}(\pi_n) = E_{\pi_n} \left[\sum_{t'=t}^{t+T} \gamma^{t'} R_{n,t'} \right], \forall n \in \{1, \dots, N\} \quad (20)$$

where, $\gamma \in [0, 1)$ is a discount factor that determines each
 MG agent's bias towards rewards received at different time
 instances. An agent with $\gamma = 0$ is a purely-myopic decision
 maker, which favors immediate reward at the expense of later
 expected reward values. On the other hand, $\gamma = 1$ represents
 an unbiased agent, which assigns equal weights to the reward
 received at all time instants. This parameter is user-defined and
 depends on each MG's economic priorities. The expectation
 operation $E_{\pi_n}\{\cdot\}$ is used to calculate reward with respect to the
 future expected action-states, which are in turn impacted by
 the uncertainties of states and observations.

7) *Constraint Return*: The SMAS-PL consists of a total of
 M constraints, including M_c^L local and M_c^G global constraints,
 defined by (3)-(4) and (5)-(16), respectively, and denoted as
 $C_m(\pi) \leq d_m, m \in \{1, \dots, M_c\}$, where $C_m(\pi)$ represents the
 return value of m 'th constraint under the control policy π and
 d_m is the upper-boundary of the m 'th constraint. Note that
 all constraints in the power management problem have been
 transformed into this format (equality constraint (15) can be
 transformed into two inequality constraints). Constraint satis-
 faction is encoded into the SMAS-PL using the discounted
 constraint return values of agents' policies π as:

$$J_{C_m}(\pi) = E_{\pi} \left[\sum_{t'=t}^{t+T} \gamma^{t'} C_{m,t'} \right] \leq d_m, \forall m \in \{1, \dots, M_c\} \quad (21)$$

where, expectation operation has been leveraged in (21) to
 handle the state and observation uncertainties.

C. Safe Policy Learning Formulation

Given the definitions of the components of the SMAS-
 PL (Section III-B), the power management problem of the
 networked MGs (1)-(16) is transformed into an iterative
 SMAS-PL problem, where the control policies of the agents
 are updated at time t , around their latest values, by maximiz-
 ing a reward function (22), while satisfying constraint return
 criteria:

$$\boldsymbol{\pi}^{t+1} = \arg \max_{\pi_1, \dots, \pi_N} \sum_{n=1}^N J_{R_n}(\pi_n) \quad (22)$$

$$s.t. \quad \mathbf{a}_n \sim \pi_n(\mathbf{S}_n) \quad (23)$$

$$J_{C_m}(\pi) \leq d_m, \forall m \quad (24)$$

$$\Delta(\pi_n, \pi_n^t) \leq \delta, \forall n \quad (25)$$

where, $\boldsymbol{\pi} = \{\pi_1, \dots, \pi_n\}$ denotes the set of control policies
 of all agents. In (23), the agent's policy is a function of the
 state vector, \mathbf{S}_n . In (24), the expected constraint return value
 are used to ensure the satisfaction of m 'th constraint based on

control policies. In (25), $\Delta(\cdot, \cdot)$ is the Kullback Leibler (KL)-divergence function [15] that serves as a distance measure between the previous policy, π_n^t , and the updated policy, π_n^{t+1} , and is constrained by a step size, δ . Note that (25) ensures that consecutive policies are within close distance from each other.

The intractable non-convex PL formulation, (22)-(25), can be solved in principle using a trust region policy optimization (TRPO) method [15]; however, in this article we apply a further approximation to TRPO to transform the problem into a tractable convex iterative QCLP, which enables learning the PL parameters, $\theta = \{\theta_1, \dots, \theta_N\}$, in a more scalable and efficient manner. Our solution leverages the linear approximations of the objective and constraint returns around the latest parameter values θ^t :

$$\theta^{t+1} = \arg \max_{\theta_1, \dots, \theta_N} \sum_{n=1}^N \mathbf{g}_n^T (\theta_n - \theta_n^t) \quad (26)$$

$$s.t. \quad J_{c_m}(\theta^t) + \mathbf{b}_m^T (\theta - \theta^t) \leq d_m, \quad \forall m \quad (27)$$

$$\frac{1}{2} (\theta_n - \theta_n^t)^T H_n (\theta_n - \theta_n^t) \leq \delta, \quad \forall n \quad (28)$$

where, $\mathbf{g}_n = \nabla_{\theta} J_R$ and $\mathbf{b}_m = \nabla_{\theta} J_{C_m}$ are the *gradient factors* of the reward and constraint return functions w.r.t. the learning parameters. Constraint (25) is transformed into (28) using the Fisher information matrix (FIM) of the policy functions, π_n , denoted by H_n . The FIM is a positive semi-definite matrix, whose (c, d) 'th entry for policy functions with a Gaussian structure is determined as follows [24]:

$$\begin{aligned} H_n(c, d) &= E \left[\frac{\partial \log \pi_n(\mathbf{a}_n | \theta_n)}{\partial \theta_n(c)} \frac{\partial \log \pi_n(\mathbf{a}_n | \theta_n)}{\partial \theta_n(d)} \right] \\ &= 2 \left(\frac{\partial \mu_n^H}{\partial \theta_n(c)} \Sigma_n^{-1} \frac{\partial \mu_n}{\partial \theta_n(d)} \right) \\ &\quad + \text{Tr} \left\{ \Sigma_n^{-1} \frac{\partial \Sigma_n}{\partial \theta_n(c)} \Sigma_n^{-1} \frac{\partial \Sigma_n}{\partial \theta_n(d)} \right\} \quad (29) \end{aligned}$$

Note that (26)–(28) provides a convexified constrained gradient-based method for training the policy functions' parameters of the MG agents; using this QCLP-based strategy the agents do not need to learn an action-value function explicitly. Instead, the power-flow-based gradient factors, \mathbf{g}_n and \mathbf{b}_m , have to be determined for the two sets of learning parameters, $[\theta_{\mu_n}, \theta_{\Sigma_n}]$. This process is outlined in Section III-D.

D. Gradient Factor Determination

To determine gradient factors, the following information are used: (i) the observation variables, \mathbf{O}_t , including nodal voltage \mathbf{V} and current injections \mathbf{I} ; (ii) the latest system states $\mathbf{S}_{n,t}$ for each MG agent; (iii) the latest control actions \mathbf{a}_n of each MG agent; (iv) the latest learning parameters $\theta_n = [\theta_{\mu_n}, \theta_{\Sigma_n}]$; (v) network parameters, including the nodal admittance matrix, \mathbf{Y} . Using information (i)-(v) and chain rule, $\mathbf{g}_n = [\mathbf{g}_{\mu_n}, \mathbf{g}_{\Sigma_n}]$ and $\mathbf{b}_m = [\mathbf{b}_{m, \mu_n}, \mathbf{b}_{m, \Sigma_n}]$ in (26) and (27) can be written as:

$$\mathbf{g}_{\mu_n} = \frac{\partial J_{R_n}}{\partial \mathbf{a}_n} \frac{\partial \mathbf{a}_n}{\partial \pi_n} \frac{\partial \pi_n}{\partial \mu_n} \frac{\partial \mu_n}{\partial \theta_{\mu_n}} \quad (30a)$$

$$\mathbf{b}_{m, \mu_n} = \frac{\partial J_{C_m}}{\partial \mathbf{a}_n} \frac{\partial \mathbf{a}_n}{\partial \pi_n} \frac{\partial \pi_n}{\partial \mu_n} \frac{\partial \mu_n}{\partial \theta_{\mu_n}} \quad (30b)$$

$$\mathbf{g}_{\Sigma_n} = \frac{\partial J_{R_n}}{\partial \mathbf{a}_n} \frac{\partial \mathbf{a}_n}{\partial \pi_n} \frac{\partial \pi_n}{\partial \Sigma_n} \frac{\partial \Sigma_n}{\partial \theta_{\Sigma_n}} \quad (31a)$$

$$\mathbf{b}_{m, \Sigma_n} = \frac{\partial J_{C_m}}{\partial \mathbf{a}_n} \frac{\partial \mathbf{a}_n}{\partial \pi_n} \frac{\partial \pi_n}{\partial \Sigma_n} \frac{\partial \Sigma_n}{\partial \theta_{\Sigma_n}} \quad (31b)$$

where, each gradient factor, \mathbf{g}_{μ_n} , \mathbf{b}_{m, μ_n} , \mathbf{g}_{Σ_n} , and \mathbf{b}_{m, Σ_n} , consists of four elements. All the elements in (30) and (31) can be obtained as follows.

1) $\partial J_{R_n} / \partial \mathbf{a}_n$ and $\partial J_{C_m} / \partial \mathbf{a}_n$: The gradients of the expected reward J_{R_n} and the expected constraint return J_{C_m} w.r.t. control actions \mathbf{a}_n can be obtained using a proposed four-step process, that leverages the current injection-based AC power flow equations. The details of this process are shown in Appendix A.

2) $\partial \mathbf{a}_n / \partial \pi_n$: Using the latest values for parameters μ_n , Σ_n , and actions \mathbf{a}_n , the gradient of control actions w.r.t. π_n is obtained from (17), as shown in (32):

$$\frac{\partial \mathbf{a}_n}{\partial \pi_n} = - \left(\frac{\Sigma_n^{-1} (\mathbf{a}_n - \mu_n)}{\sqrt{|\Sigma_n|} (2\pi)^{D_n}} e^{-\frac{1}{2}A} \right)^{-1} \quad (32)$$

where, $A = (\mathbf{a}_n - \mu_n)^T \Sigma_n^{-1} (\mathbf{a}_n - \mu_n)$. The detailed derivation of (32) can be found in Appendix B.

3) $\partial \pi_n / \partial \mu_n$ and $\partial \pi_n / \partial \Sigma_n$: Using the latest values for parameters μ_n , Σ_n and actions \mathbf{a}_n , the gradients of control policies, w.r.t. μ_n and Σ_n are determined using (17), as shown in (33) and (34):

$$\frac{\partial \pi_n}{\partial \mu_n} = \frac{\Sigma_n^{-1} (\mathbf{a}_n - \mu_n)}{\sqrt{|\Sigma_n|} (2\pi)^{D_n}} e^{-\frac{1}{2}A} \quad (33)$$

$$\frac{\partial \pi_n}{\partial \Sigma_n} = - \frac{1}{2} \frac{(\Sigma_n^{-1} - \Sigma_n^{-1} (\mathbf{a}_n - \mu_n) (\mathbf{a}_n - \mu_n)^T \Sigma_n^{-1})}{\sqrt{|\Sigma_n|} (2\pi)^{D_n}} e^{-\frac{1}{2}A} \quad (34)$$

where, the detailed derivations of (33) and (34) are shown in Appendix B.

4) $\partial \mu_n / \partial \theta_{\mu_n}$ and $\partial \Sigma_n / \partial \theta_{\Sigma_n}$: A *back-propagation process* [25] is performed on the two DNNs within each MG agent's control policy function, (18) and (19), to determine the gradients of DNNs' outputs w.r.t. their parameters. In each iteration, the latest values of state variables are employed as inputs of the DNNs. The back-propagation process exploits chain rule for stage-by-stage spreading of gradient information through layers of the DNNs, starting from the output layer and moving towards the input [25]. To enhance the stability of the back-propagation process, a sample batch approach is adopted, where the gradients obtained from several sampled actions are averaged to ensure robustness against outliers.

IV. MULTI-AGENT CONSENSUS-BASED SAFE POLICY LEARNING

A. Offline Policy Training

Using the gradient factors (30) and (31), the QCLP, (26)–(28), is fully specified and can be solved at each policy update iteration for training the agents' PL frameworks. However, we have identified two challenges in this problem: (i) the size of the DNN parameters θ can be extremely large,

579 which results in high computational costs during training; (ii)
580 the control policy privacy of the MG agents needs to be
581 preserved during training, which implies that the agents might
582 not have access to each other's control policies, cost func-
583 tions, and local constraints on assets. Centralized solvers can
584 be both time-consuming and lack guarantees for maintaining
585 data ownership boundaries.

586 In order to address these two challenges, we have developed
587 a *multi-agent consensus-based constrained training algo-*
588 *rithm* [18]. Due to its distributed nature this method is both
589 scalable and does not require sharing control policy param-
590 eters among agents. Thus, the proposed algorithm is able
591 to efficiently solve the QCLP (26)–(28), while relying only
592 on local inter-MG communication. The purpose of inter-MG
593 interactions is to satisfy global constraints, (3)–(4). To do this,
594 the agents repeatedly estimate and communicate dual vari-
595 able λ_n , corresponding to the Lagrangian multiplier of global
596 constraints. Furthermore, a local primal-dual gradient step is
597 included in the algorithm to move the primal and dual param-
598 eters towards their global optimum. The proposed distributed
599 algorithm consists of four stages that are performed iteratively,
600 as follows:

601 *Stage I [Initialize ($k \leftarrow 1$)]:* Gradient factors \mathbf{g}_n and \mathbf{b}_m are
602 obtained from Section II-D. The previous values of learning
603 parameters are input to the QCLP, $\theta_n^t(0) \leftarrow \theta_n^{t-1}$. Lagrangian
604 multipliers are initialized as zero for each MG agent.

605 *Stage II (Weighted Averaging Operation):* MG agent n
606 receives the Lagrangian multiplier $\lambda_{n'}$, for global con-
607 straints (3)–(4), from its neighbouring MG agents $n' \in$
608 $\{1, \dots, N_n\}$ and combines the received estimates using
609 weighted averaging:

$$\bar{\lambda}_n(k) = \sum_{n'=1}^{N_n} w_n(n') \lambda_{n'}(k) \quad (35)$$

611 where, $w_n(n')$ is the weight that MG agent n assigns to the
612 incoming message of the neighbouring MG agent n' . To guar-
613 antee convergence to consensus, the weight matrix, composed
614 of the agents' weight parameters is selected as a doubly
615 stochastic matrix [18], i.e., $w_n(n') = \frac{1}{N_n}$. This weight selection
616 strategy implies that the MG agents assign equal importance
617 to the information received from their neighboring agents.

618 *Stage III (Primal Gradient Update):* The n 'th MG agent
619 updates its primal parameters θ_n^t employing a gradient descent
620 operation, using the gradients of the agent's reward and the
621 global constraint returns, $m' \in M_C^G$, and step size ρ_1 :

$$\bar{\theta}_n(k) = \theta_n^t(k) - \rho_1 (\mathbf{g}_n(\theta_n^t(k)) + \mathbf{b}_{m'}(\theta_n^t(k)) \bar{\lambda}_n(k)). \quad (36)$$

623 *Stage IV (Projection on Local Constraints):* The agent
624 projects the local learning parameters to the feasible region
625 defined by the gradients of the local constraints (5)–(16):

$$\theta_n^t(k+1) = \arg \min_{\theta} \|\bar{\theta}_n(k) - \theta\| \quad (37)$$

$$s.t. \quad J_{c_m}(\theta_n^t(0)) + \mathbf{b}_m^T(\theta_n^t(0) - \theta) \leq d_m, \quad \forall m \in M_C^L \quad (38)$$

$$\frac{1}{2}(\theta_n^t(0) - \theta)^T H_n(\theta_n^t(0) - \theta) \leq \delta, \quad \forall n. \quad (39)$$

Algorithm 1 SMAS-PL Training

```

1: Select  $t^{max}, T, \delta, k^{max}, w_n(n'), \rho_1, \rho_2, \Delta\theta_n$ 
2: Initialize  $\theta_n^t$ 
3: for  $t \leftarrow 1$  to  $t^{max}$  do
4:    $\mathbf{S}_n \leftarrow [\mathbf{S}_n(t), \dots, \mathbf{S}_n(t+T)]$ 
5:    $\boldsymbol{\mu}_n \leftarrow (18)$  [Parameter insertion]
6:    $\boldsymbol{\Sigma}_n \leftarrow (19)$  [Parameter insertion]
7:    $\mathbf{a}_n \sim \pi_n(\mathbf{S}_n | \theta_n) \leftarrow (17)$  [Action selection]
8:    $\partial J_{R_n} / \partial \mathbf{a}_n \leftarrow (55)$ –(56)
9:    $\partial J_{C_m} / \partial \mathbf{a}_n \leftarrow (59), (57)$ –(58)
10:   $\partial \mathbf{a}_n / \partial \pi_n \leftarrow (32)$ 
11:   $\partial \pi_n / \partial \boldsymbol{\mu}_n \leftarrow (33)$ 
12:   $\partial \pi_n / \partial \boldsymbol{\Sigma}_n \leftarrow (34)$ 
13:   $\partial \boldsymbol{\mu}_n / \partial \theta_{\boldsymbol{\mu}_n} \leftarrow DNN_{\boldsymbol{\mu}_n}$  [Back-propagation]
14:   $\partial \boldsymbol{\Sigma}_n / \partial \theta_{\boldsymbol{\Sigma}_n} \leftarrow DNN_{\boldsymbol{\Sigma}_n}$  [Back-propagation]
15:   $\mathbf{g}_{\boldsymbol{\mu}_n}, \mathbf{b}_{m, \boldsymbol{\mu}_n} \leftarrow (30)$  [Chain rule]
16:   $\mathbf{g}_{\boldsymbol{\Sigma}_n}, \mathbf{b}_{m, \boldsymbol{\Sigma}_n} \leftarrow (31)$  [Chain rule]
17:   $H_n \leftarrow (29)$  [FIM Construction]
18:  Initialize  $\lambda_n(k_0)$ 
19:  for  $k \leftarrow 1$  to  $k^{max}$  do
20:     $\bar{\lambda}_n(k) \leftarrow (35)$  [Averaging operation]
21:     $\bar{\theta}_n(k) \leftarrow (36)$  [Primal gradient update]
22:     $\theta_n^t(k+1) \leftarrow (37)$ –(39) [Projection on  $M^L$ ]
23:     $\lambda_n(k+1) \leftarrow (40)$  [Dual gradient update]
24:    if  $\|\theta_n^t(k+1) - \theta_n^t(k)\| \leq \Delta\theta_n$  then
25:       $\theta_n^{t+1} \leftarrow \theta_n^t(k+1)$ ; Break;
26:    end if
27:  end for
28:  if  $\|\theta_n^{t+1} - \theta_n^t\| \leq \Delta\theta_n$  then
29:    Output  $\theta_n^* \leftarrow \theta_n^{t+1}$ ; Break;
30:  end if
31: end for
32: Output well-trained parameterized policy  $\pi_n(\theta_n^*)$ 

```

629 *Stage V (Dual Gradient Update):* Each agent's estimations
630 of dual variables λ_n for the global constraints, (3) and (4), will
631 be updated using a gradient ascent process over $\bar{\lambda}_n$:

$$\lambda_n(k+1) = [(\bar{\lambda}_n(k) + \rho_2(\mathbf{b}_{m'} \theta_n^t(k+1) - d_{m'}))]^+, \quad \forall m' \in M_C^G \quad (40)$$

634 where, ρ_2 is a penalty factor for global constraints violation,
635 and the operator $[\cdot]^+$ returns the non-negative part of its input.

636 *Stage VI (Stopping Criteria):* Check algorithm convergence
637 using the changes of $\theta_n^t(k)$; stop when the changes in param-
638 eters falls below the threshold value $\Delta\theta_n$; otherwise, go back
639 to Stage II.

640 The overall flowchart of the SMAS-PL training process
641 using the proposed distributed training technique is shown in
642 Algorithm 1. The calculations of Steps 8 and 9 can be found
643 in Appendix A.

B. Online Action Selection

645 The trained policy functions are used by the MG agents for
646 online action selection. This process can be simply represented
647 as sampling from the learned Gaussian policy functions (17).

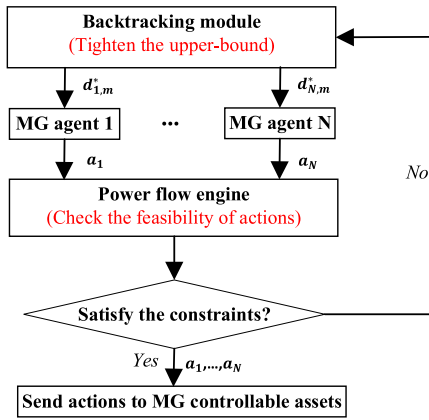
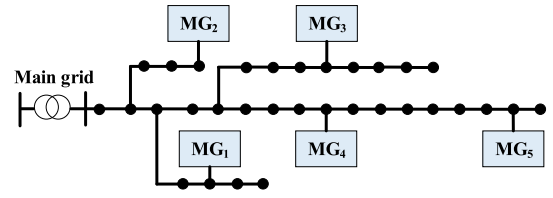
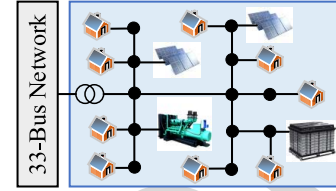


Fig. 2. Flowchart of the backtracking strategy.



(a) 33-bus system for distribution network



(b) 13-bus system for MGs

Fig. 3. Test system under study.

TABLE I
SELECTED COST FUNCTION PARAMETERS

Description	Notion	Value
Average electricity price (\$/kWh)	λ^R	0.046
Average DG fuel price (\$/L)	λ^f	0.57
Fuel cons. quadratic function parameter (L/kW ²)	a^f	0.0001773
Fuel cons. quadratic function parameter (L/kW)	b^f	0.1709
Fuel cons. quadratic function parameter (L)	c^f	14.67

648 First, the agents receive the latest values of the states, includ-
 649 ing the predicted solar irradiance and aggregate internal load
 650 power of MGs. These values are inserted into the trained
 651 DNNs (18) and (19) to obtain the mean and covariance matri-
 652 ces of the policy functions. Finally, samples are generated
 653 from the multivariate Gaussian distributions. These samples
 654 are averaged and passed to the local controllers of each
 655 controllable asset as a reference signal.

656 C. Backtracking Strategy

657 Due to convex approximations in the formulations
 658 (26)–(28), it is possible for few global constraints to be
 659 marginally violated in practice. To ensure feasibility, we can
 660 add a backtracking strategy into the proposed solution. This
 661 closed-loop backtracking strategy consists of two components,
 662 as shown in Fig. 2.

663 *Component 1 [Power Flow Engine (PFE)]:* The PFE
 664 receives the control actions from MG agents and runs a sim-
 665 ple power flow program to obtain the status of all constraints.
 666 If no constraint is violated, the control signals are passed to
 667 controllable assets. If some constraints are violated, then the
 668 PFE will engage the backtracking process.

669 *Component 2 (Backtracking Module):* The backtracking
 670 module tightens the upper-bound limit (d_m) (only) for the con-
 671 straints that have been violated. The parameters of the trained
 672 DNNs will be re-updated according to update rules (35)–(40)
 673 and with the modified upper-bounds. The purpose of tight-
 674 ening the upper-bound is to provide a safety margin. In this
 675 article the tightening process is performed using a user-defined
 676 coefficient multiplier, $0 < \tau < 1$, as follows:

$$677 \quad d_m^* = \tau d_m. \quad (41)$$

678 V. SIMULATION RESULTS

679 The proposed method is tested on a modified 33-bus dis-
 680 tribution network [26], which consists of five MGs as shown
 681 in Fig. 3(a). Each MG is modeled as a modified IEEE 13-bus
 682 network [26] at a low voltage level as shown in Fig. 3(b).
 683 When calculating the gradient factors, a single-phase AC
 684 power flow model is used for the sake of brevity. In the case
 685 study, the base power value is 100 kVA and base voltage values

686 in the 33-bus distribution network and 13-bus MG networks
 687 are 12.66 kV and 4.16 kV, respectively.

688 The input data for load demands and PV generations have
 689 15-minute time resolution are obtained from smart meter
 690 database [27] to provide realistic numerical experiments. The
 691 assumption in this article is that smart meters are installed
 692 throughout the network and the agents have access to a diverse
 693 data. The training and testing datasets are selected through
 694 uniform randomization to ensure that the proposed solver func-
 695 tions reasonably. Here, 1-month of the randomly selected data
 696 is used for testing and 11-month of the data is used for training.
 697 The energy price for the transferred power at the MG PCCs
 698 and the fuel price for the local DGs are adopted from [28]
 699 and [29], respectively. The quadratic polynomial parameters of
 700 DG fuel consumption are adopted from [30]. Table I presents
 701 selected parameters for operational cost calculation in simula-
 702 tions. The average capacities for DGs in MGs are 60 kWh. The
 703 average capacities for ESSs in MGs are 20 kWh, the maximum
 704 charging/discharging rate is 4kW and the charging/discharging
 705 efficiencies are 95% and 90%, respectively.

706 All the case studies are simulated using a PC with Intel Core
 707 i7-4790 3.6 GHz CPU and 16 GB RAM hardware. The sim-
 708 ulations are performed in MATLAB [31] and OpenDSS [32]
 709 to obtain the gradient factors, update the learning parameters,
 710 solve the distributed training problem, and validate the results.
 711 In training, each episode is a learning update iteration based
 712 on the data that comes from one moving decision window. The
 713 length of the moving window is 4 samples with a 15-minute
 714 time step, which gives us a 1-hour window. The activation
 715 functions of each layer (including the output layer) of the

TABLE II
SELECTED DNN HYPERPARAMETERS AND USER-DEFINED COEFFICIENTS

Description	Notion	Value
Length of the decision window in episode	T	4
Discount factor	γ	0.99
Step size for updating θ	δ	1×10^{-3}
Maximum iteration	k^{max}	200
Weight assigned to received information	w_n	0.2
Step size for primal gradient update	ρ_1	0.01
Step size for dual gradient update	ρ_2	0.01
Threshold for parameter updating	$\Delta\theta$	1×10^{-4}
Tightening multiplier	τ	0.9
Number of hidden layer	-	3
Number of neurons per hidden layer	-	10
Size of minibatches	-	128
Activation function of DNNs	-	<i>tansig</i>

716 feedforward networks are hyperbolic tangent-sigmoid (*tansig*).
 717 After various numerical tests, the parameters θ_μ and θ_Σ of
 718 the neural networks are initialized using uniform distributions
 719 defined over the intervals $(0, 0.2)$ and $(-0.03, 0.03)$, respec-
 720 tively. In our simulations, we have observed that $\tau = 0.9$ is
 721 sufficient for ensuring feasibility for those constraints that have
 722 been marginally violated after one-to-two rounds of backtrack-
 723 ing. Table II summarizes selected DNN hyperparameters and
 724 other user-defined coefficients in simulations. The hyperpa-
 725 rameters were optimized using a randomly-selected validation
 726 set (2 months worth of data) and Bayesian optimization with
 727 uninformative priors in MATLAB environment.

728 Further, to demonstrate the effectiveness of SMAS-PL,
 729 three benchmark methods have been considered, including
 730 an optimization-based method, an on-policy method and
 731 an off-policy method. The first benchmark method is an
 732 optimization-based method, which leverages YALMIP tool-
 733 box to solve the optimal power management of networked
 734 MGs using IBM ILOF CPLEX 12.9. The second one is the
 735 unconstrained policy gradient learning (U-PL) method, which
 736 leverages the same algorithm as the proposed SMAS-PL, how-
 737 ever, certain constraints are removed during the training process
 738 of U-PL. By comparing the SMAS-PL and the U-PL, we can
 739 show the effectiveness of the SMAS-PL when handling different
 740 local and global constraints. The U-PL can be considered as an
 741 on-policy benchmark. We also consider an off-policy bench-
 742 mark method, namely the deep Q-network (DQN). In [23],
 743 [33], DQN uses deep neural networks (DNNs) to approximate
 744 the Q-function and provide Q-value estimation for discretized
 745 control actions. To include the constraints in DQN, we have
 746 followed the suggestion in [23], [34] and added penalty terms
 747 to the reward function of the benchmark DQN to discour-
 748 age constraint violation. The penalty coefficients for global
 749 and local constraints are manually tuned based on the DQN
 750 performance. However, since the benchmark DQN was not
 751 originally designed for continuous actions, we have first dis-
 752 cretized the agents' action space with a step size of 33% of
 753 the constraint upper limit. For example, if the upper limit of a
 754 diesel generation (DG) power output is 60 kW, then, the power
 755 output action of DG has been discretized as $(0, 20, 40, 60)$ kW.

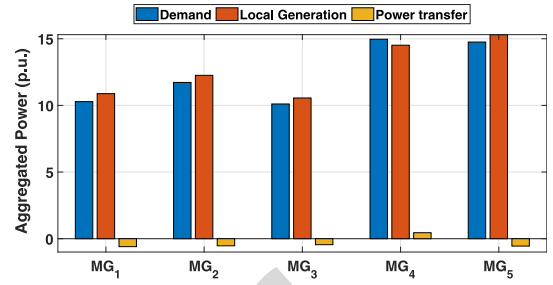


Fig. 4. Aggregated power of local demand, local generation and power transfer for MG_1 - MG_5 .

Similar discretization has been applied to the actions of PV
 inverters and ESSs. The inputs of the DNN are the system
 states, and the outputs of the DNN are estimations for the
 Q-value function for each discrete action. The DNN is param-
 eterized as a function approximator to represent the Q-value
 function. The temporal difference (TD) learning algorithm is
 used to train the DNN by minimizing the mean-squared TD
 error. The discount factor and learning rate in DQN are set to
 the same values as those of SMAS-PL. The exploration factor
 is set to 0.1 in the ϵ -greedy action selection of DQN. The
 structure of DNN in the benchmark DQN has been obtained
 using cross-validation. The dimensions of the input and output
 layers have been extended by the number of MG agents and the
 number of discrete actions. Note that the benchmark U-PL is
 implemented in a multi-agent framework, while the benchmark
 DQN is implemented in a centralized way.

A. System Operation Outcomes

In the case study, action selection is performed by sampling
 100 times from the trained policy functions (distributions).
 Then the dispatching action is obtained by averaging the
 selected samples. A trade-off is involved in choosing the
 number of action samples: if this number is too large, then
 the selected actions will converge to the policy mean, which
 implies that model uncertainties are ignored. This could result
 in erroneous and sub-optimal solutions in case the learned
 model is over-fitting (i.e., when the estimated mean has large
 errors). On the other hand, if the number of samples is too
 small, then the outcomes can deviate from the learned mean
 value, which can also result in low-quality outcomes. The
 average outcomes are shown in Fig. 4, Fig. 5 and Table III.
 The aggregate MG demand, aggregate MG generation, and
 aggregate power transfer through PCCs of MGs over a day
 are shown in Fig. 4. It can be seen that the main MG
 demands are supplied by the local generation within MGs
 due to low DG fuel prices and renewable outputs. While
 most MGs are exporting power to the upstream distribu-
 tion grid, MG_4 is importing power to satisfy the heavy local
 load that cannot be fully supplied internally. In all cases
 the power balance is maintained within the MGs. The ESS
 SOC for each MG are shown in Fig. 5, where can be seen
 that ESSs charge during off-peak period and discharge dur-
 ing peak time to provide optimal power balancing support
 for MGs. Table III presents comparisons between the bench-
 mark optimization-based method, the benchmark DQN and

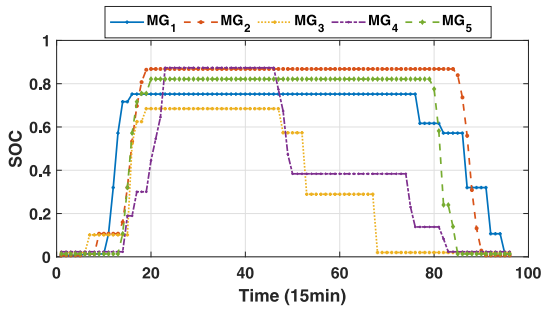
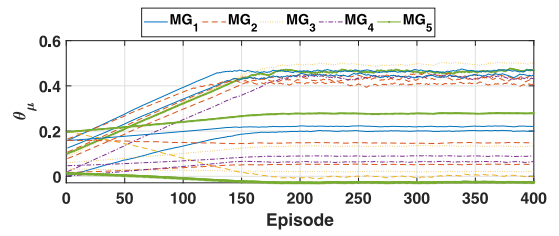
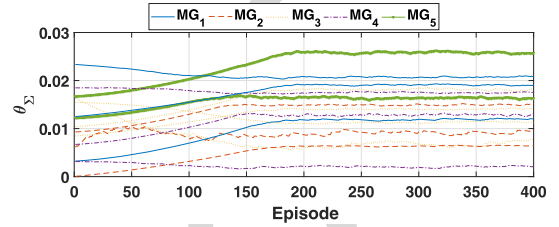
Fig. 5. ESS dispatching results for MG_1 - MG_5 .

TABLE III
COMPARISON BETWEEN CENTRALIZED SOLVER, DQN AND SMAS-PL METHOD

	Cen. solver	DQN	SMAS-PL
Average daily cost (\$)	1356.60	1928.4	1372.11
Average time (second)	145.50	10.30	1.40 (per agent)
MG privacy maintenance	No	No	Yes

800 the proposed SMAS-PL, including the average daily cost of
801 operation over numerous scenarios, average online decision
802 time, and MG privacy maintenance.

803 In general, the SMAS-PL method has three fundamental
804 advantages over centralized optimization method: 1) Even
805 though the offline training process in our method takes a
806 long time (around 35 minutes per agent), the average online
807 decision time for the proposed SMAS-PL is about only
808 1.4 seconds per agent, which is much shorter than the average
809 time 145.5 seconds for the centralized optimization solver.
810 Thus, the real-time response of the trained policy function is
811 almost 100 times faster than that of the OPF solver. The reason
812 for this is that the OPF solver needs to find the optimal
813 solution of a complex optimization problem in real-time, while
814 our approach simply samples from multivariate Gaussian distributions
815 that embody optimal control policies. Furthermore,
816 we have observed that the computational cost of the centralized
817 OPF solver rises almost quadratically with the size of the
818 system; beyond a certain point the commercial solver is not
819 able to provide solutions in a reasonable time. On the other
820 hand, our SMAS-PL retains an almost constant online decision
821 time, while the cost of offline training increases almost
822 linearly. 2) The proposed PL method takes advantage of a
823 multi-agent (distributed) framework to train the policy function
824 of each MG agent; in practice, this distributed framework
825 can be implemented using parallel computation techniques,
826 which also enhances the scalability of the proposed SMAS-PL
827 method compared to centralized solvers. 3) Due to its distributed
828 nature, the proposed SMAS-PL method maintains the
829 privacy and data ownership boundaries of individual MGs.
830 During the training process, the MG agents do not need to
831 share control policy parameters, policy functions, cost functions,
832 and local asset constraints with each other. The only
833 variables that are shared among MG agents are the Lagrangian
834 multipliers corresponding to global network constraints. These
835 multipliers do not have a physical meaning and thus, do not
836 contain sensitive information.

(a) Selected θ_μ 's during training process(b) Selected θ_Σ 's during training processFig. 6. Convergence of learning parameters θ_μ and θ_Σ for MG_1 - MG_5 .

837 Based on the comparison between the centralized solver and
838 our proposed method, there is still a 1.14% difference between
839 the solutions from the centralized solver and the SMAS-PL
840 method, which might be caused by the following reasons:
841 (i) Unlike the centralized solver, which has access to the full
842 systemic model information, and thus, can guarantee at least a
843 local optimal solution, the proposed SMAS-PL method lacks a
844 guarantee of optimality. Also, in order to obtain a high-quality
845 solution, the SMAS-PL needs to first approximate the original
846 problem with a convex surrogate, which despite enhancing the
847 problem tractability, comes at the expense of loss of accuracy
848 and a reduction in performance. (ii) The proposed backtracking
849 mechanism is a heuristic strategy, which is aimed at obtaining
850 a feasible solution that might come at the expense of a loss
851 in the reward. (iii) To obtain a consensus-based solution, the
852 SMAS-PL needs a reliable inter-agent communication infrastruc-
853 ture, which could be costly. (iv) In case of changes in
854 system structure, the SMAS-PL will need an offline re-training
855 phase to adapt to new system conditions. This could take some
856 time, during which the agents will experience a temporary
857 decline in their payoffs. The comparison between DQN and
858 SMAS-PL is discussed in Section V-B.

B. Algorithm Performance

859 Fig. 6(a) and Fig. 6(b) show the convergence of a selected
860 group of learning parameters, θ_μ and θ_Σ during the training
861 process, for each MG agent. As can be seen, the changes in
862 θ_μ are relatively larger than that of θ_Σ . This is due to the
863 higher levels of sensitivity of MG agents' objective functions
864 to the mean values of the control actions compared with their
865 variance levels.
866

867 In Fig. 7, the average hourly rewards under SMAS-PL,
868 U-PL, and DQN are compared with each other. Note that
869 here, the moving average rewards and the episodic rewards
870 of different methods are depicted by dark and light curves.
871 It can be observed that SMAS-PL and U-PL both outper-
872 form DQN in term of the total reward. The reason for this is

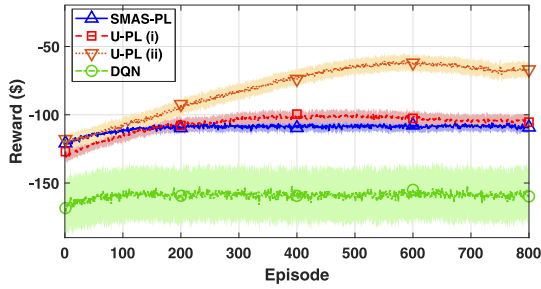


Fig. 7. Comparison of the average hourly rewards with different methods.

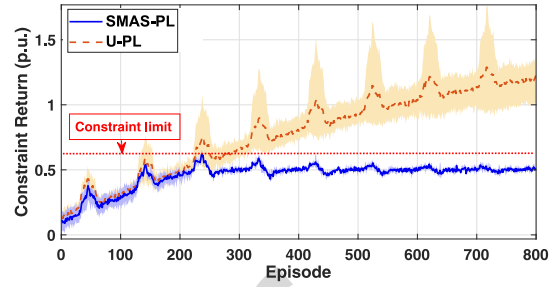
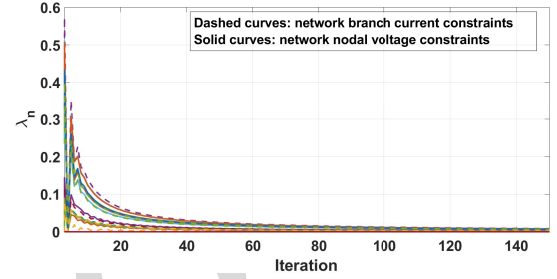

 Fig. 8. Comparison of constraint values w/ and w/o DG capacity constraints in MG_1 .


Fig. 9. The performance of the iterative distributed training method in one episode (no binding global constraints).

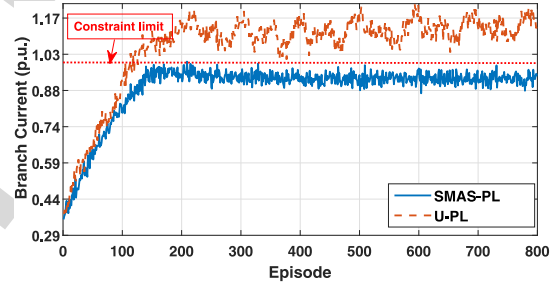


Fig. 10. Selected global branch current constraint return values for MG agents.

873 that the SMAS-PL and U-PL leverage the proposed iterative
 874 and distributed technique to adaptively tune the Lagrangian
 875 multipliers through information exchange between MG agents;
 876 on the other hand, the DQN needs to manually design penalty
 877 coefficients for constraint violations, which either offers inad-
 878 equate penalization of the constraint violations or excessive
 879 punishment for the constraints. Also, SMAS-PL and U-PL
 880 have continuous action spaces, while DQN employs action dis-
 881 cretization, which hinders accurate exploration of action space.
 882 After the NNs are fully trained, the SMAS-PL samples the
 883 actions from the learned multivariate Gaussian distributions
 884 that embody optimal control policies, while the benchmark
 885 DQN selects the control actions that have the highest estimated
 886 Q-values for the given state according to the trained DNN.
 887 Based on the results in Table III, the decision time for the
 888 SMAS-PL is around 1.4 seconds per agent, while the decision
 889 time for the benchmark DQN is approximately 10.3 seconds.
 890 Thus, the decision time for the proposed SMAS-PL is faster
 891 than the benchmark DQN, because the multi-agent framework
 892 enables the SMAS-PL to sample decision actions in parallel
 893 for each MG agent, while the benchmark DQN selects the
 894 control actions for all the MGs together in a centralized way.
 895 Two cases are considered in implementing U-PL: (i) no DG
 896 capacity constraints for MG_1 and MG_2 ; (ii) no DG capacity
 897 constraints for MG_1 - MG_5 . In cases (i) and (ii) of the U-PL,
 898 the agents obtain a higher reward compared to the SMAS-
 899 PL due to the constraint omission; however, this comes at the
 900 expense of decision infeasibility. In case of the SMAS-PL,
 901 these operational constraints are satisfied, which also leads
 902 to a drop in total reward, as expected. This shows that our
 903 proposed constrained PL decision model can ensure the feasi-
 904 bility of the control actions w.r.t. the constraints of the power
 905 management problem. Note that even though the proposed
 906 SMAS-PL framework is similar to the TRPO [15], the TRPO
 907 has theoretical guarantees for monotonic increase in return,
 908 while such guarantees do not exist for the approximate QCLP
 909 formulation in the proposed SMAS-PL. However, compared
 910 to TRPO our solution offers a simpler, more efficient, and
 911 tractable alternative, with fewer learning parameters.

912 Furthermore, Fig. 8 shows the constraint values during the
 913 training iterations for a 1-hour time window, for the two cases
 914 with and without DG capacity constraints in MG_1 , where the
 915 dark blue and red curves represent averaged constraint val-
 916 ues, and the light blue and red areas represent the variations
 917 around the average curves for the SMAS-PL and U-PL, respec-
 918 tively. During the training process, the U-PL violates the upper

919 boundary for DG generation limit (i.e., local constraint case
 920 study); on the other hand, the SMAS-PL solver satisfies the
 921 DG generation capacity constraints, which implies that the
 922 local constraints can be safely maintained. Therefore, com-
 923 pared to U-PL, the proposed SMAS-PL has shown to be able
 924 to generate control actions that not only improve the reward
 925

926 One example of the distributed training convergence pro-
 927 cess is shown in Fig. 9 for a policy gradient update step. As
 928 can be seen, the Lagrangian multipliers λ_n reach zero over
 929 iterations of the proposed multi-agent algorithm, which indi-
 930 cates that all the global constraints, including nodal voltage
 931 and branch current limits, are satisfied and feasible solutions
 932 are obtained. This also means that the bus voltage and line
 933 current constraints are not binding for this case.

934 Another example is given to demonstrate the effectiveness
 935 of the SMAS-PL in handling binding global constraints. This
 936 case shows a line flow constraint in the grid under the proposed
 937 SMAS-PL and a U-PL baseline; as is observed in Fig. 10,
 938 the U-PL has generated infeasible decisions that violate the
 939 constraint, while our approach has prevented the flow to go

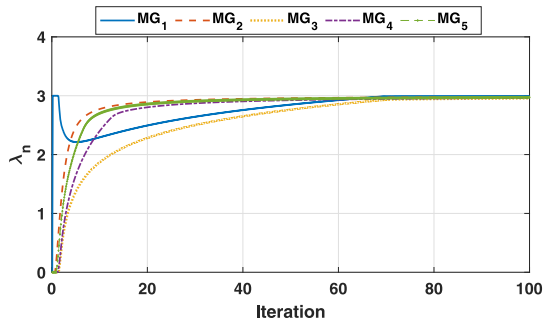


Fig. 11. MG agents' consensus on λ_n for the selected global constraint.

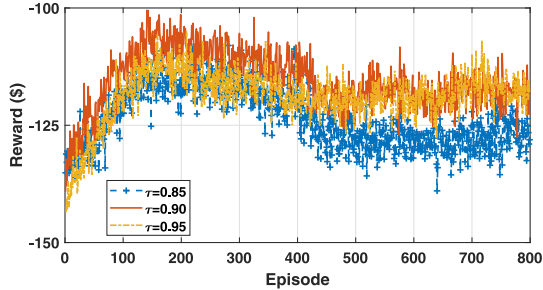


Fig. 12. Impact of backtracking on algorithm performance.

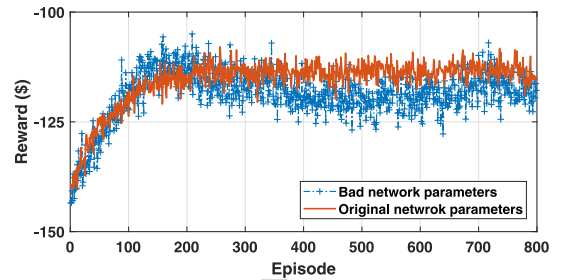


Fig. 13. Analysing the impact of bad network data on decision model outcomes.

convergence, the model still reaches reward values close to the ideal case. However, due to the information imperfection, a loss of reward is inevitable.

VI. CONCLUSION

Conventional model-based optimization methods suffer from high computational costs when solving large-scale multi-MG power management problems. On the other hand, the conventional model-free methods are black-box tools, which may fail to satisfy the operational constraints. Motivated by these challenges, in this article, a SMAS-PL method has been proposed for power management of networked MGs. Our proposed method exploits the gradients of the decision problem to learn control policies that achieve both optimality and feasibility. Furthermore, to enhance computational efficiency and maintain the policy privacy of the control agents, a distributed consensus-based training process is implemented to update the agents' policy functions over time using local communication.

Note that the current case study has been conducted over a balanced single-phase distribution system. However, our proposed SMAS-PL is not limited to single-phase distribution systems and can be potentially extended to unbalanced three-phase systems. One solution to this challenge could be using a single policy function for the resources connected to all phases (note that theoretically-speaking our method is not limited by the number of phases). However, this brute-force solution may lack scalability. Another solution extension to a multi-phase system cannot be fully addressed by having three separate policy functions per phase. A more efficient and scalable extension to unbalanced systems remains the subject of our future research.

APPENDIX A

CALCULATION OF $\partial J_{R_n}/\partial \mathbf{a}_n$ AND $\partial J_{C_m}/\partial \mathbf{a}_n$

The major difficulty in determining $\partial J_{R_n}/\partial \mathbf{a}_n$ and $\partial J_{C_m}/\partial \mathbf{a}_n$ pertains to the agents' reward functions and global constraint returns, (1)-(4), which are only implicitly related to the control actions. Since the reward and all the global constraint returns are functions of the observation variables, \mathbf{V} and \mathbf{I} , the gradients of these variables w.r.t. control actions are obtained and used to quantify $\partial J_{R_n}/\partial \mathbf{a}_n$ and $\partial J_{C_m}/\partial \mathbf{a}_n$. To do this, a four-step process is proposed that leverages the current injection-based AC power flow equations:

above its upper bound. Further, as can be seen in Fig. 11, the Lagrangian multipliers for this binding constraint reach a non-zero constant number over iterations. This also shows the agents' estimations of Lagrange multipliers for a global line flow constraint; as can be seen, using the proposed SMAS-PL the agents are capable of reaching consensus on the value of the multiplier without having any access to each other's policy functions, which corroborates the performance of our proposed method under incomplete information.

To validate the tightening parameter levels (τ), we have studied the impact of different τ values on the reward. Here, at episode 400, the value of τ is decreased from 1 to (0.95, 0.9, 0.85). The average rewards for different drops in τ are compared in Fig. 12. It can be observed that for values of τ close to 1 (i.e., $\tau = 0.95$ and $\tau = 0.9$) the reward values are very close to each other. However, as τ deviates from unity and reaches $\tau = 0.85$, the reward drops significantly. In our simulation, we have observed that $\tau = 0.9$ is sufficient for ensuring feasibility for those few constraints that have been marginally violated in certain operation scenarios after one-to-two rounds of backtracking. Note that this threshold needs to be fine-tuned for specific grids.

To simulate the impact of bad network parameter data on training model performance, we have added random errors (with a 10% variance) to the network resistance (R) and reactance (X) parameters during the training process. The bad network data will lead to errors in gradient factors (42)-(59) (see Appendix A). To validate the SMAS-PL under network data imperfection, we have compared the average reward obtained with perfect knowledge of network parameters and under bad network parameter information. It can be observed in Fig. 13, even though the learning process with bad network data shows more volatility and needs more time to reach

TABLE IV
PARTIAL DERIVATIONS OF I^{Re} AND I^{Im} W.R.T.
 $\mathbf{a}_n = [P_n^{DG}, P_n^{Ch}, P_n^{Dis}, Q_n^{DG}, Q_n^{PV}, Q_n^{ESS}]$

\mathbf{a}_n	P_n^{DG}	P_n^{Ch}	P_n^{Dis}	Q_n^{DG}	Q_n^{PV}	Q_n^{ESS}
$I_{i,t'}^{Re}$	$-\frac{V_{i,t'}^{Re}}{V_{i,t'}^2}$	$\frac{V_{i,t'}^{Re}}{V_{i,t'}^2}$	$-\frac{V_{i,t'}^{Re}}{V_{i,t'}^2}$	$\frac{V_{i,t'}^{Im}}{V_{i,t'}^2}$	$\frac{V_{i,t'}^{Im}}{V_{i,t'}^2}$	$-\frac{V_{i,t'}^{Im}}{V_{i,t'}^2}$
$I_{i,t'}^{Im}$	$-\frac{V_{i,t'}^{Im}}{V_{i,t'}^2}$	$\frac{V_{i,t'}^{Im}}{V_{i,t'}^2}$	$-\frac{V_{i,t'}^{Im}}{V_{i,t'}^2}$	$-\frac{V_{i,t'}^{Re}}{V_{i,t'}^2}$	$-\frac{V_{i,t'}^{Re}}{V_{i,t'}^2}$	$\frac{V_{i,t'}^{Re}}{V_{i,t'}^2}$

Step 1 - First, the gradients of real and imaginary parts of nodal current injection w.r.t. control actions are derived (denoted as $\partial I^{Re}/\partial \mathbf{a}_n$ and $\partial I^{Im}/\partial \mathbf{a}_n$, respectively.) To achieve this, the nodal power balance and nodal current injection relationships in the network are employed [35]:

$$I_{i,t'}^{Re} = \frac{P_{i,n,t'} V_{i,t'}^{Re} + q_{i,n,t'} V_{i,t'}^{Im}}{V_{i,t'}^2} \quad (42)$$

$$I_{i,t'}^{Im} = \frac{P_{i,n,t'} V_{i,t'}^{Im} - q_{i,n,t'} V_{i,t'}^{Re}}{V_{i,t'}^2} \quad (43)$$

$$P_{i,n,t'} = P_{i,n,t'}^D - P_{i,n,t'}^{DG} - P_{i,n,t'}^{PV} + P_{i,n,t'}^{Ch} - P_{i,n,t'}^{Dis} \quad (44)$$

$$q_{i,n,t'} = Q_{i,n,t'}^D - Q_{i,n,t'}^{DG} - Q_{i,n,t'}^{PV} + Q_{i,n,t'}^{ESS} \quad (45)$$

where, I_i^{Re} , I_i^{Im} and V_i^{Re} , V_i^{Im} denote the real and imaginary parts of nodal voltage and current injection at node i . Using these equations, $\partial I^{Re}/\partial \mathbf{a}_n$ and $\partial I^{Im}/\partial \mathbf{a}_n$ are derived and shown in Table IV. Note that the entries of this table can be calculated using the real and imaginary parts of nodal voltages, which in practice are either measured or estimated [35].

Step 2 - Using $\partial I^{Re}/\partial \mathbf{a}_n$ and $\partial I^{Im}/\partial \mathbf{a}_n$ from Step 1 (Table IV), $\partial V^{Re}/\partial \mathbf{a}$ and $\partial V^{Im}/\partial \mathbf{a}$ are obtained employing the network-wide relationship between nodal voltages and current injections:

$$\begin{bmatrix} \frac{\partial V^{Re}}{\partial \mathbf{a}_n} \\ \frac{\partial V^{Im}}{\partial \mathbf{a}_n} \end{bmatrix} = \begin{bmatrix} Y^{11} & Y^{12} \\ Y^{21} & Y^{22} \end{bmatrix}^{-1} \begin{bmatrix} \frac{\partial I^{Re}}{\partial \mathbf{a}_n} \\ \frac{\partial I^{Im}}{\partial \mathbf{a}_n} \end{bmatrix} \quad (46)$$

where, the modified network bus admittance sub-matrices are determined as follows:

$$Y^{11} = Y^{Re} - Y_D^{(Re,Re)}, Y^{12} = -Y^{Im} - Y_D^{(Re,Im)} \quad (47)$$

$$Y^{21} = Y^{Im} - Y_D^{(Im,Re)}, Y^{22} = Y^{Re} - Y_D^{(Im,Im)} \quad (48)$$

here, Y^{Re} and Y^{Im} are the real and imaginary parts of the original bus admittance matrix. The elements in diagonal matrices $Y_D^{(Re,Re)}$, $Y_D^{(Re,Im)}$, $Y_D^{(Im,Re)}$ and $Y_D^{(Im,Im)}$ are calculated using the following equations [35]:

$$Y_D^{(Re,Re)}(i, i) = \frac{P_{i,n,t'}}{V_{i,t'}^2} - \frac{2V_{i,t'}^{Re}(p_{i,n,t'} V_{i,t'}^{Re} + q_{i,n,t'} V_{i,t'}^{Im})}{V_{i,t'}^4} \quad (49)$$

$$Y_D^{(Re,Im)}(i, i) = \frac{q_{i,n,t'}}{V_{i,t'}^2} - \frac{2V_{i,t'}^{Im}(p_{i,n,t'} V_{i,t'}^{Re} + q_{i,n,t'} V_{i,t'}^{Im})}{V_{i,t'}^4} \quad (50)$$

$$Y_D^{(Im,Re)}(i, i) = -\frac{q_{i,n,t'}}{V_{i,t'}^2} - \frac{2V_{i,t'}^{Re}(p_{i,n,t'} V_{i,t'}^{Im} - q_{i,n,t'} V_{i,t'}^{Re})}{V_{i,t'}^4} \quad (51)$$

$$Y_D^{(Im,Im)}(i, i) = \frac{p_{i,n,t'}}{V_{i,t'}^2} - \frac{2V_{i,t'}^{Im}(p_{i,n,t'} V_{i,t'}^{Im} - q_{i,n,t'} V_{i,t'}^{Re})}{V_{i,t'}^4} \quad (52)$$

Step 3 - Noting that the current flow constraint returns and the rewards are also functions of branch current flows, the gradients of branch current flows are required to obtain $\partial J_{R_n}/\partial \mathbf{a}_n$ and $\partial J_{C_m}/\partial \mathbf{a}_n$. Using the branch current flow equations, these gradients are determined as a function of the derivatives of nodal voltages and current injections, as follows:

$$\frac{\partial I_{ij,t'}^{Re}}{\partial \mathbf{a}_{n,t'}} = y_{ij}^{Im} \left(\frac{\partial V_{i,t'}^{Im}}{\partial \mathbf{a}_{n,t'}} - \frac{\partial V_{j,t'}^{Im}}{\partial \mathbf{a}_{n,t'}} \right) - y_{ij}^{Re} \left(\frac{\partial V_{i,t'}^{Re}}{\partial \mathbf{a}_{n,t'}} - \frac{\partial V_{j,t'}^{Re}}{\partial \mathbf{a}_{n,t'}} \right) \quad (53)$$

$$\frac{\partial I_{ij,t'}^{Im}}{\partial \mathbf{a}_{n,t'}} = y_{ij}^{Im} \left(\frac{\partial V_{i,t'}^{Re}}{\partial \mathbf{a}_{n,t'}} - \frac{\partial V_{j,t'}^{Re}}{\partial \mathbf{a}_{n,t'}} \right) + y_{ij}^{Re} \left(\frac{\partial V_{i,t'}^{Im}}{\partial \mathbf{a}_{n,t'}} - \frac{\partial V_{j,t'}^{Im}}{\partial \mathbf{a}_{n,t'}} \right) \quad (54)$$

where, I_{ij}^{Re} and I_{ij}^{Im} are the real and imaginary parts of branch currents, y_{ij}^{Re} and y_{ij}^{Im} are the real and imaginary parts of branch admittance.

Step 4 - Finally, using the derivatives obtained from Steps 1, 2, and 3, $\partial J_{R_n}/\partial \mathbf{a}_n$ and $\partial J_{C_m}/\partial \mathbf{a}_n$ are determined through straightforward algebraic manipulations. As an example, the gradient of reward function w.r.t. $P_{n,t'}^{DG}$ is calculated as:

$$\frac{\partial J_{R_n}}{\partial P_{n,t'}^{DG}} = \sum_{t'=t}^{t+T} \left(\lambda_{i,n}^F (2a_f + b_f) - \lambda_n^R \frac{\partial P_{n,t'}^{PCC}}{\partial P_{n,t'}^{DG}} \right) \quad (55)$$

where, $\partial P_{n,t'}^{PCC}/\partial P_{n,t'}^{DG}$ is obtained using the outcomes of Steps 2 and 3, as follows:

$$\begin{aligned} \frac{\partial P_{n,t'}^{PCC}}{\partial P_{n,t'}^{DG}} &= \frac{\partial V_{i,t'}^{Re}}{\partial P_{n,t'}^{DG}} I_{ij,t'}^{Re} + V_{i,t'}^{Re} \frac{\partial I_{ij,t'}^{Re}}{\partial P_{n,t'}^{DG}} \\ &+ \frac{\partial V_{i,t'}^{Im}}{\partial P_{n,t'}^{DG}} I_{ij,t'}^{Im} + V_{i,t'}^{Im} \frac{\partial I_{ij,t'}^{Im}}{\partial P_{n,t'}^{DG}} \end{aligned} \quad (56)$$

Furthermore, $\partial J_{C_m}/\partial \mathbf{a}_n$ for the global constraints (3) and (4) can be calculated using the outcomes of Steps 2 and 3:

$$\frac{\partial V_{i,t'}}{\partial \mathbf{a}_{n,t'}} = \frac{V_{i,t'}^{Re}}{V_{i,t'}} \frac{\partial V_{i,t'}^{Re}}{\partial \mathbf{a}_{n,t'}} + \frac{V_{i,t'}^{Im}}{V_{i,t'}} \frac{\partial V_{i,t'}^{Im}}{\partial \mathbf{a}_{n,t'}} \quad (57)$$

$$\frac{\partial I_{ij,t'}}{\partial \mathbf{a}_{n,t'}} = \frac{I_{ij,t'}^{Re}}{I_{ij,t'}} \frac{\partial I_{ij,t'}^{Re}}{\partial \mathbf{a}_{n,t'}} + \frac{I_{ij,t'}^{Im}}{I_{ij,t'}} \frac{\partial I_{ij,t'}^{Im}}{\partial \mathbf{a}_{n,t'}} \quad (58)$$

As can be seen in (5)-(16), the local constraint returns are trivial functions of the control actions. For example, the constraint return value for (5) is $J_{C_{5,t'}} = P_{n,t'}^{DG}$ which induces a simple gradient element w.r.t. control action $P_{n,t'}^{DG}$:

$$\frac{\partial J_{C_{5,t'}}}{\partial P_{n,t'}^{DG}} = 1 \quad (59)$$

The gradients of constraint returns w.r.t. control actions for the remaining local constraints, (6)-(16), can be obtained in a similar way.

APPENDIX B

DERIVATION OF $\partial \mathbf{a}_n / \partial \pi_n$, $\partial \pi_n / \partial \boldsymbol{\mu}_n$ AND $\partial \pi_n / \partial \Sigma_n$

$\partial \mathbf{a}_n / \partial \pi_n$, $\partial \pi_n / \partial \boldsymbol{\mu}_n$ and $\partial \pi_n / \partial \Sigma_n$ are obtained using the probability density function of (D-dimensional) multivariate Gaussian distribution [36], which has the following general formulation:

$$f(\mathbf{x}; \boldsymbol{\mu}, \Sigma) = \frac{1}{\sqrt{|\Sigma|(2\pi)^D}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})} \quad (60)$$

where \mathbf{x} is a random vector. To derive the gradients, first, the log-likelihood function of this multivariate Gaussian distribution (60) is obtained as follows:

$$L = \ln(f) = \ln \frac{1}{\sqrt{|\Sigma|(2\pi)^D}} - \frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}) \quad (61)$$

The derivative of L w.r.t. mean vector $\boldsymbol{\mu}$ and covariance matrix Σ can be written as follows:

$$\begin{aligned} \frac{\partial L}{\partial \boldsymbol{\mu}} &= -\frac{1}{2} \frac{\partial (\mathbf{x}-\boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})}{\partial \boldsymbol{\mu}} \\ &= -\frac{1}{2} \left(-2\Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}) \right) = \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}) \end{aligned} \quad (62)$$

$$\begin{aligned} \frac{\partial L}{\partial \Sigma} &= -\frac{1}{2} \left(\frac{\partial \ln(|\Sigma|)}{\partial \Sigma} + \frac{\partial (\mathbf{x}-\boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})}{\partial \Sigma} \right) \\ &= -\frac{1}{2} \left(\Sigma^{-1} - \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})(\mathbf{x}-\boldsymbol{\mu})^\top \Sigma^{-1} \right) \end{aligned} \quad (63)$$

Thus, using (62) and (63), the derivatives of the function f w.r.t. $\boldsymbol{\mu}$ and Σ can be shown in (64) and (65), respectively:

$$\frac{\partial f}{\partial \boldsymbol{\mu}} = \frac{\Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})}{\sqrt{|\Sigma|(2\pi)^D}} e^{-\frac{1}{2}A} \quad (64)$$

$$\frac{\partial f}{\partial \Sigma} = -\frac{1}{2} \frac{(\Sigma^{-1} - \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})(\mathbf{x}-\boldsymbol{\mu})^\top \Sigma^{-1})}{\sqrt{|\Sigma|(2\pi)^D}} e^{-\frac{1}{2}A} \quad (65)$$

where $A = (\mathbf{x}-\boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})$. Similarly, the derivative of the function f w.r.t. \mathbf{x} is shown as follows:

$$\frac{\partial f}{\partial \mathbf{x}} = -\frac{\Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})}{\sqrt{|\Sigma|(2\pi)^D}} e^{-\frac{1}{2}A}. \quad (66)$$

REFERENCES

- [1] Z. Li, M. Shahidehpour, F. Aminifar, A. Alabdulwahab, and Y. Al-Turki, "Networked microgrids for enhancing the power system resilience," *Proc. IEEE*, vol. 105, no. 7, pp. 1289–1310, Jul. 2017.
- [2] M. N. Alam, S. Chakrabarti, and A. Ghosh, "Networked microgrids: State-of-the-art and future perspectives," *IEEE Trans. Ind. Informat.*, vol. 15, no. 3, pp. 1238–1250, Mar. 2019.
- [3] M. Ban, M. Shahidehpour, J. Yu, and Z. Li, "A cyber-physical energy management system for optimal sizing and operation of networked nanogrids with battery swapping stations," *IEEE Trans. Sustain. Energy*, vol. 10, no. 1, pp. 491–502, Jan. 2019.
- [4] Z. Wang, B. Chen, J. Wang, M. 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.
- [5] L. Xiao, N. B. Mandayam, and H. V. Poor, "Prospect theoretic analysis of energy exchange among microgrids," *IEEE Trans. Smart Grid*, vol. 6, no. 1, pp. 63–72, Jan. 2015.
- [6] 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.
- [7] D. Shi *et al.*, "A distributed cooperative control framework for synchronized reconnection of a multi-bus microgrid," *IEEE Trans. Smart Grid*, vol. 9, no. 6, pp. 6646–6655, Nov. 2018.
- [8] B. Zhao *et al.*, "Energy management of multiple microgrids based on a system of systems architecture," *IEEE Trans. Power Syst.*, vol. 33, no. 6, pp. 6410–6421, Nov. 2018.
- [9] J. Duan, Z. Yi, D. Shi, C. Lin, X. Lu, and Z. Wang, "Reinforcement-learning-based optimal control of hybrid energy storage systems in hybrid AC–DC microgrids," *IEEE Trans. Ind. Informat.*, vol. 15, no. 9, pp. 5355–5364, Sep. 2019.
- [10] P. Zeng, H. Li, H. He, and S. Li, "Dynamic energy management of a microgrid using approximate dynamic programming and deep recurrent neural network learning," *IEEE Trans. Smart Grid*, vol. 10, no. 4, pp. 4435–4445, Jul. 2019.
- [11] B.-G. Kim, Y. Zhang, M. van der Schaar, and J.-W. Lee, "Dynamic pricing and energy consumption scheduling with reinforcement learning," *IEEE Trans. Smart Grid*, vol. 7, no. 5, pp. 2187–2198, Sep. 2016.
- [12] T. Yu, H. Z. Wang, B. Zhou, K. W. Chan, and J. Tang, "Multi-agent correlated equilibrium Q(λ) learning for coordinated smart generation control of interconnected power grids," *IEEE Trans. Power Syst.*, vol. 30, no. 4, pp. 1669–1679, Jul. 2015.
- [13] X. Lu, X. Xiao, L. Xiao, C. Dai, M. Peng, and H. V. Poor, "Reinforcement learning-based microgrid energy trading with a reduced power plant schedule," *IEEE Internet Things J.*, vol. 6, no. 6, pp. 10728–10737, Dec. 2019.
- [14] Q. Zhang, K. Dehghanpour, Z. Wang, and Q. Huang, "A learning-based power management method for networked microgrids under incomplete information," *IEEE Trans. Smart Grid*, vol. 11, no. 2, pp. 1193–1204, Mar. 2020.
- [15] J. Achiam, D. Held, A. Tamar, and P. Abbeel, "Constrained policy optimization," in *Proc. 34th Int. Conf. Mach. Learn.*, 2017, pp. 22–31.
- [16] M. Wen and U. Topcu, "Constrained cross-entropy method for safe reinforcement learning," in *Proc. 32nd Conf. NeurIPS*, 2018, pp. 7461–7471.
- [17] Y. Chow, O. Nachum, M. Ghavamzadeh, and E. Duenez-Guzman, "A Lyapunov-based approach to safe reinforcement learning," in *Proc. 32nd Conf. NeurIPS*, 2018, pp. 8103–8112.
- [18] T.-H. Chang, A. Nedic, and A. Scaglione, "Distributed constrained optimization by consensus-based primal-dual perturbation method," *IEEE Trans. Autom. Control*, vol. 59, no. 6, pp. 1524–1538, Jun. 2014.
- [19] G. Mokryani, "Active distribution networks planning with integration of demand response," *Solar Energy*, vol. 122, pp. 1362–1370, Dec. 2015.
- [20] M. Wahbah, T. H. M. El-Fouly, B. Zahawi, and S. Feng, "Hybrid Beta-KDE model for solar irradiance probability density estimation," *IEEE Trans. Sustain. Energy*, vol. 11, no. 2, pp. 1110–1113, Apr. 2020.
- [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] H. Li, Z. Wan, and H. He, "Constrained EV charging scheduling based on safe deep reinforcement learning," *IEEE Trans. Smart Grid*, vol. 11, no. 3, pp. 2427–2439, May 2020.
- [24] O. Besson and Y. I. Abramovich, "On the Fisher information matrix for multivariate elliptically contoured distributions," *IEEE Signal Process. Lett.*, vol. 22, no. 11, pp. 1130–1133, Nov. 2013.
- [25] J. Li, J. Cheng, J. Shi, and F. Huang, "Brief introduction of back propagation neural network algorithm and its improvement," in *Advances in Computer Science and Information Engineering*. Heidelberg, Germany: Springer, 2012, pp. 553–558.
- [26] W. H. Kersting, "Radial distribution test feeders," in *Proc. IEEE PES Winter Meeting*, vol. 2, 2001, pp. 908–912.
- [27] *Iowa Distribution Test Systems*. [Online]. Available: <http://wzy.ece.iastate.edu/Testsystem.html>
- [28] *Wholesale Electricity and Natural Gas Market Data*. [Online]. Available: <https://www.eia.gov/electricity/wholesale/#history>
- [29] *Gasoline and Diesel Fuel Update*. [Online]. Available: <https://www.eia.gov/petroleum/gasdiesel/>
- [30] 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.
- [31] *MATLAB*. [Online]. Available: <https://www.mathworks.com/help/matlab/index.html>
- [32] *OpenDSS Basic Tutorial*. [Online]. Available: <https://sourceforge.net/p/electricdss/discussion/beginners.html>

- 1201 [33] Y. Ye, D. Qiu, X. Wu, G. Strbac, and J. Ward, "Model-free real-
1202 time autonomous control for a residential multi-energy system using
1203 deep reinforcement learning," *IEEE Trans. Smart Grid*, vol. 11, no. 4,
1204 pp. 3068–3082, Jul. 2020.
- 1205 [34] W. Wang, N. Yu, Y. Gao, and J. Shi, "Safe off-policy deep reinforcement
1206 learning algorithm for Volt-VAR control in power distribution systems,"
1207 *IEEE Trans. Smart Grid*, vol. 11, no. 4, pp. 3008–3018, Jul. 2020.
- 1208 [35] R. A. Jabr and I. Džafić, "Sensitivity-based discrete coordinate-descent
1209 for Volt/VAr control in distribution networks," *IEEE Trans. Power Syst.*,
1210 vol. 31, no. 6, pp. 4670–4678, Nov. 2016.
- 1211 [36] P. K. Brandt and M. S. Pedersen, *The Matrix Cookbook*. Lyngby,
1212 Denmark: Tech. Univ. Denmark, 2008.



Feng Qiu (Senior Member, IEEE) received the Ph.D. 1253
degree from the School of Industrial and Systems 1254
Engineering, Georgia Institute of Technology in 1255
2013. He is a Principal Computational Scientist 1256
with the Energy Systems Division, Argonne National 1257
Laboratory, Argonne, IL, USA. His current research 1258
interests include optimization in power system oper- 1259
ations, electricity markets, and power grid resilience. 1260



Qianzhi Zhang (Graduate Student Member, IEEE) 1213
received the M.S. degree in electrical and com- 1214
puter engineering from Arizona State University in 1215
2015, and he is currently pursuing the Ph.D. degree 1216
with the Department of Electrical and Computer 1217
Engineering, Iowa State University, Ames, IA, 1218
USA. He has worked with Huadian Electric Power 1219
Research Institute as a Research Engineer, from 1220
2015 to 2016. His research interests include the 1221
applications of machine learning and optimization 1222
techniques in power system operation and control. 1223



Kaveh Dehghanpour (Member, IEEE) received the 1224
B.Sc. and M.S. degrees in electrical and computer 1225
engineering from the University of Tehran, in 2011 1226
and 2013, respectively, and the Ph.D. degree in elec- 1227
trical engineering from Montana State University in 1228
2017. He was a Postdoctoral Research Associate 1229
with Iowa State University. His research interests 1230
include the applications of machine learning and 1231
data-driven techniques in power system monitoring 1232
and control. 1233



Zhaoyu Wang (Senior Member, IEEE) received the 1234
B.S. and M.S. degrees in electrical engineering from 1235
Shanghai Jiaotong University in 2009 and 2012, 1236
respectively, and the M.S. and Ph.D. degrees in elec- 1237
trical and computer engineering from the Georgia 1238
Institute of Technology in 2012 and 2015, respec- 1239
tively. He is the Harpole-Pentair Assistant Professor 1240
with Iowa State University. His research interests 1241
include power distribution systems and microgrids, 1242
particularly on their data analytics and optimization. 1243
He was the recipient of the 2020 IEEE Power and 1244

Energy Society (PES) Outstanding Young Engineer Award. He is the Secretary 1245
of IEEE PES Award Subcommittee, a Co-Vice Chair of PES Distribution 1246
System Operation and Planning Subcommittee, and a Vice Chair of PES 1247
Task Force on Advances in Natural Disaster Mitigation Methods. He is 1248
an Associate Editor of IEEE TRANSACTIONS ON POWER SYSTEMS, IEEE 1249
TRANSACTIONS ON SMART GRID, IEEE POWER AND ENERGY SOCIETY 1250
LETTERS, IEEE OPEN ACCESS JOURNAL OF POWER AND ENERGY, and 1251
IET Smart Grid. 1252



Dongbo Zhao (Senior Member, IEEE) received the 1261
B.S. degree in electrical engineering from Tsinghua 1262
University, Beijing, China, the M.S. degree in elec- 1263
trical engineering from Texas A&M University, 1264
College Station, TX, USA, and the Ph.D. degree in 1265
electrical engineering from the Georgia Institute of 1266
Technology, Atlanta, GA, USA. 1267

From 2014 to 2016, he has worked with Eaton 1268
Corporation as a Lead Engineer in its Corporate 1269
Research and Technology Division, and with ABB 1270
in its U.S. Corporate Research Center, from 2010 1271
to 2011. He is Currently a Principal Energy System Scientist with Argonne 1272
National Laboratory, Lemont, IL, USA. He is also an Institute Fellow with 1273
Northwestern Argonne Institute of Science and Engineering, Northwestern 1274
University. His research interests include power system control, protection, 1275
reliability analysis, transmission and distribution automation, and electric 1276
market optimization. He is an Editor of IEEE TRANSACTIONS ON POWER 1277
DELIVERY, IEEE TRANSACTIONS ON SUSTAINABLE ENERGY, and IEEE 1278
POWER ENGINEERING LETTERS. He is the Subject Editor of subject "Power 1279
system operation and planning with renewable power generation" of *IET* 1280
Renewable Power Generation and an Associate Editor of IEEE ACCESS. 1281
He is a Member of IEEE PES, IAS, and IES Societies. 1282