# Independent policy gradient-based reinforcement learning for economic and reliable energy management of multi-microgrid systems

Junkai Hu[*]    Li Xia[†]

## Abstract

Efficiency and reliability are both crucial for energy management, especially in multi-microgrid systems (MMSs) integrating intermittent and distributed renewable energy sources. This study investigates an economic and reliable energy management problem in MMSs under a distributed scheme, where each microgrid independently updates its energy management policy in a decentralized manner to optimize the long-term system performance collaboratively. We introduce the mean and variance of the exchange power between the MMS and the main grid as indicators for the economic performance and reliability of the system. Accordingly, we formulate the energy management problem as a mean-variance team stochastic game (MV-TSG), where conventional methods based on the maximization of expected cumulative rewards are unsuitable for variance metrics. To solve MV-TSGs, we propose a fully distributed independent policy gradient algorithm, with rigorous convergence analysis, for scenarios with known model parameters. For large-scale scenarios with unknown model parameters, we further develop a deep reinforcement learning algorithm based on independent policy gradients, enabling data-driven policy optimization. Numerical experiments in two scenarios validate the effectiveness of the proposed methods. Our approaches fully leverage the distributed computational capabilities of MMSs and achieve a well-balanced trade-off between economic performance and operational reliability.

[*]J. Hu is with the School of Mechanical and Electrical Engineering, Shenzhen Polytechnic University, Shenzhen 518055, China.

[†]L. Xia is with the School of Business, Sun Yat-Sen University, Guangzhou 510275, China. (email: xiali5@sysu.edu.cn)

# 1    Introduction

Microgrids function as the foundational components of smart grids, providing an operational environment for the efficient control and utilization of distributed energy resources and consumer demand loads. To mitigate the impact of uncertainties in these resources and demand loads, microgrids are typically equipped with energy storage devices, which can shift energy between different time intervals by discharging or charging, facilitating energy management (Weitzel & Glock, 2018). However, the energy management capabilities of individual microgrids are often limited by their storage capacity. With advances in communication infrastructure, there is increasing attention on multi-microgrid systems (MMSs), where microgrids are interconnected through distribution buses and communication networks to better utilize distributed energy resources and mitigate power network uncertainties. Effective energy management methods can significantly improve both the economic efficiency and operational stability of MMSs (Alam et al., 2018).

In MMSs, each microgrid is equipped with its own energy management system (EMS) and implements a local policy for controlling its dispatchable resources. To facilitate collaboration among microgrids, various EMS topologies have been proposed, generally classified into four categories (Nawaz et al., 2022), as illustrated in Figure 1: Centralized EMS; Decentralized EMS; Hybrid EMS; Distributed EMS. The distributed scheme stands out for the absence of a centralized controller or coordinator. Instead, each microgrid communicates with neighboring microgrids via communication infrastructure and updates its energy management policy accordingly. This approach leverages the distributed computational capabilities of MMSs and avoids a single point of failure. In this work, we focus on the energy management problem in MMSs under the distributed EMS scheme.

A majority of the literature on energy management in MMSs focuses on achieving economic objectives within rigid operational constraints, such as minimizing operational costs
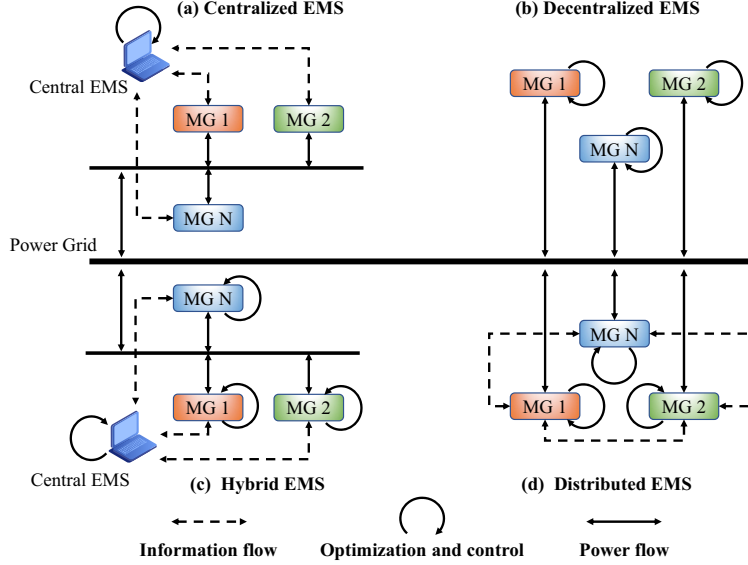
Figure 1: The illustration of EMS topologies in MMSs (MG denotes the microgrid).

(Chen et al., 2022) and increasing the penetration of renewable energy (Han et al., 2025). However, integrating a large number of uncertain units can negatively impact the reliability and stability of power network operations. The intermittent and unpredictable nature of renewable energy generation leads to significant power fluctuations, which may cause various power quality issues, including voltage instability and flicker effects (Yang et al., 2020).

To mitigate the fluctuations in exchange power, a thread of work focuses on power smoothing by formulating optimization objectives based on power ramp rates, which quantify the rate of power change between two consecutive time intervals (Arévalo et al., 2023; Wang et al., 2023). Another thread of work aims to achieve peak shaving and valley filling within the scheduling horizon, which is also a widely adopted practical objective (Tang et al., 2021; Manojkumar et al., 2022). It is evident that in MMS energy management, where renewable energy generation and loads are uncertain, reducing exchange power fluctuations is essential to ensuring the reliable operation of the overall system.

Traditional energy management methods, such as programming-based approaches, typically rely on accurate system modeling or precise predictions of uncertain variables. However, the increasing integration of stochastic components and the expanding scale of power networks have substantially heightened system complexity, posing significant challenges to con-

ventional model-based optimization and programming techniques. In response, data-driven approaches have emerged as a promising direction and have attracted widespread attention in the field of energy management (Li et al., 2023).

Reinforcement learning (RL) has recently gained growing interest in energy management for power networks, due to its notable success in complex decision-making tasks, ranging from mastering the game of Go (Silver et al., 2017) to empowering large language models. RL is a data-driven approach for solving stochastic sequential decision-making problems, typically modeled as Markov decision processes (MDPs). It enables an agent to repeatedly interact with an environment with unknown dynamics and rewards, collecting data on state transitions and rewards to iteratively update its policy and maximize expected cumulative rewards. Among various RL approaches, policy gradient methods are the most widely used, where the policy is iteratively updated along the gradient direction of the objective performance function (Agarwal et al., 2021). Additionally, compared to traditional tabular RL methods, deep reinforcement learning (DRL) incorporates deep neural networks to approximate value and policy functions, alleviating the curse of dimensionality and enhancing generalization and scalability.

However, standard RL or DRL methods are primarily designed to maximize the expected cumulative rewards. In many real-world applications, particularly in power networks, it is essential to consider reliability and risk-sensitive performance metrics, which are usually variance-related and do not exhibit the Markovian or additive characteristics of reward functions. These features make the dynamic programming principle fail (Filar et al., 1989; Xia, 2020) and limit the direct applicability of conventional RL algorithms to economic and reliability-aware energy management problems.

In this work, we investigate long-term economic and reliable energy management in MMSs under the distributed EMS scheme, aiming at updating the energy management policies of individual microgrids to optimize the exchange power between the MMS and the main grid. The average exchange power over the scheduling horizon corresponds to electricity sold (positive) or purchased (negative), reflecting the economic aspect of MMS operation. We introduce the variance metric as an additional optimization objective to capture the

4

volatility of exchange power, which reflects the reliability and stability of MMS operation. Since the problem requires multiple microgrids to jointly optimize a variance-based metric, standard policy gradient and RL methods, which only maximize the expected cumulative rewards, are not directly applicable. To address this challenge, we formulate the problem as a mean-variance team stochastic game (MV-TSG) and propose corresponding algorithms to solve it. The main contributions of this paper are as follows:

- Based on the characteristics of the distributed EMS scheme, we propose a mean-variance independent projected gradient ascent (MV-IPGA) method for MV-TSGs, as illustrated in Algorithm 1. In this approach, each microgrid independently and simultaneously updates its energy management policy to improve the overall economic and reliability performance of MMS operation. We also establish the global convergence of the proposed algorithm. Compared to methods that require centralized coordination, the proposed approach utilizes the distributed computing resources in MMSs more efficiently.

- Building upon Algorithm 1, we further develop an independent policy gradient-based reinforcement learning algorithm, referred to as mean-variance independent proximal policy optimization (MV-IPPO), as presented in Algorithm 2. This algorithm enables approximate solutions to large-scale MV-TSGs in a data-driven manner when environmental parameters are unknown.

The effectiveness of the proposed algorithms is validated through field-data-based experiments under two distinct scenarios: one with fully known model parameters and another with unknown parameters. The results demonstrate that our methods successfully achieve a balance between the operational economy and reliability of MMSs within a distributed EMS framework. Compared with existing works Qiu et al. (2021) and Shen et al. (2023), this study, to the best of our knowledge, is the first to develop a *fully distributed, risk-sensitive, cooperative RL algorithm with theoretical guarantees in heterogeneous multi-agent settings.*

The remainder of this paper is organized as follows. Section 2 reviews the relevant literature on reducing power fluctuations in microgrids and risk-sensitive RL. Section 3 presents

5

the problem formulation and the MV-TSG model. The MV-IPGA method and its convergence analysis are introduced in Section 4. The MV-IPPO algorithm is further proposed in Section 5. Section 6 provides the numerical experiments, and Section 7 concludes the paper.

# 2 Literature review

In this section, we first review the literature on reducing power fluctuations in microgrids. Subsequently, we briefly introduce the risk-sensitive RL, which optimizes some specific risk metrics instead of the expected cumulative reward.

## 2.1 Reducing power fluctuations in microgrids

Reducing power fluctuations is essential for the reliable operation of power networks. Some studies focus on power smoothing by optimizing objectives related to power ramp rates. For example, Arévalo et al. (2023) investigate photovoltaic power smoothing based on failure detection. They propose a method combining moving averages and ramp rate control, with energy storage systems (ESSs) to mitigate power fluctuations, and validate its effectiveness in a microgrid laboratory at the University of Cuenca. Similarly, Abdalla et al. (2023) study power smoothing by incorporating cloud information in a two-layer framework: the first layer predicts and classifies cloud types, while the second dynamically adjusts the filter time constant based on power ramp rates. However, the effectiveness of such approaches heavily depends on the accuracy of the predictions, which becomes increasingly difficult as the time horizon extends. Wang et al. (2023) address power smoothing for multiple wind turbines with energy storage using a multi-agent reinforcement learning (MARL) algorithm. They reshape the reward function by incorporating the power ramp rate. Nonetheless, the aforementioned approaches, which optimize power ramp rates, are limited to mitigating fluctuations between two successive time intervals, and are insufficient for achieving peak shaving and valley filling over the entire scheduling horizon.

Peak shaving is another critical objective in energy management for reducing power fluc-

tuations. Guo et al. (2021) study a cooperation problem involving distributed generators, ESSs, and voltage regulating devices, aiming to minimize daily operational costs while constraining peak load demand. They propose a two-stage programming method, with the peak demand limit determined via trial and error. Ghafoori et al. (2023) focus on optimal scheduling for electric vehicle charging and discharging to minimize peak power demand in commercial buildings. Their approach reduces the variance between actual and forecasted minimum demand across all time intervals, using machine learning for demand forecasting and binary linear programming for optimization. Similarly, Chapaloglou et al. (2019) propose a predict-then-optimize framework for load smoothing and peak shaving in microgrids. However, these methods are primarily designed for short-term, intra-day scheduling and are less applicable to long-term energy management due to the growing difficulty of accurate predictions over extended time horizons.

Markov model is a powerful tool for characterizing stochastic dynamic systems. Yang et al. (2020) investigate long-term energy management for microgrids to reduce fluctuations in the exchange power with the main grid. They formulate the problem as an MDP with a variance objective and propose a policy iteration type method to solve it, effectively addressing both power smoothing and peak shaving in long-term scheduling. Peirelinck et al. (2024) and Rostmnezhad & Dessaint (2023) investigate demand response optimization and building energy management problems, respectively, both aiming to reduce peak power. They formulate these problems as MDPs and solve them using DRL approaches. Hu et al. (2025) study an MMS energy management problem focusing on economic and reliable long-term operation. They incorporate a variance objective to reflect reliability and propose a MARL algorithm based on a centralized training with decentralized execution (CTDE) framework. However, all aforementioned studies rely on centralized control or coordination, which limits their applicability to the distributed MMS energy management problem studied in this work.

## 2.2 Risk-sensitive reinforcement learning

RL is a subfield of machine learning that has been successfully applied across various domains to address complex sequential decision-making problems, including energy management in

power networks. Policy gradient methods, such as proximal policy optimization (PPO) (Schulman et al., 2017), are the most popular type of RL algorithms, which utilize the gradient information of the performance function to guide policy updates. However, in power networks, the complexity and inherent uncertainty of real-world scenarios suggest that maximizing cumulative rewards (e.g., economic gains) alone is insufficient. To ensure practical and safety-critical operations, it is crucial to incorporate risk and system reliability metrics into the decision-making process (Blancas-Rivera & Jasso-Fuentes, 2024).

Risk-sensitive RL is a subfield of safe RL and has long been a prominent research direction (García & Fernández, 2015). It focuses on optimizing performance criteria that account for specific risks. Common risk metrics include variance, conditional value-at-risk (CVaR), absolute deviation, and semi-variance. However, many of these measures lack additivity and may not satisfy the Markovian property, which causes the Bellman equation—the foundation of standard RL methods—to no longer hold (Bäuerle & Jaśkiewicz, 2024; Ma et al., 2023). Consequently, conventional RL algorithms often struggle with optimization problems involving such risk metrics, requiring the development of specialized approaches to address these challenges.

In the context of variance optimization in RL or MDPs, existing literature can be classified into two groups based on the variance definition. The first group concerns the variance of cumulative rewards, i.e., $\text{Var}(\sum_{t=0}^{\infty} \gamma^t r_t)$, where $\gamma$ is the discount factor and $r_t$ is the feedback reward at time step $t$ (Prashanth & Fu, 2022; Huang, 2018). This formulation is typically used to quantify fluctuations in total profits or costs. The second group focuses on the long-run variance or steady-state variance (Xia, 2018; Bisi et al., 2021; Filar et al., 1989), defined as $\lim_{T \to \infty} \frac{1}{T} \mathbb{E}_\mu [\sum_{t=0}^{T-1} (r(s_t, a_t) - \eta^\mu)^2]$, where $\eta^\mu$ is the long-run average reward under policy $\mu$. This metric captures long-term reward volatility and can serve as an indicator of system operational stability.

Due to the non-additivity of variance, standard RL algorithms are generally inadequate for solving either type of variance-based optimization problem. As a result, specialized policy gradient algorithms are often designed to address these challenges (Prashanth & Fu, 2022). However, in the energy management problem of MMSs under the distributed EMS

scheme, each microgrid operates independently, making decisions and updating its policy in a decentralized manner. When single-agent RL algorithms are used independently by each microgrid, and policies are updated simultaneously, the environment perceived by each microgrid becomes non-stationary. This violates the stationary environment assumption of single-agent RL, and leads to a lack of guarantees for improvement in the common objective function (Zhong et al., 2024). Therefore, directly applying single-agent RL methods in this setting is usually inappropriate, which highlights the complexity and challenges of cooperative risk-sensitive optimization in multi-agent settings.

A related research direction is cooperative MARL, which aims to learn behavior policies for multiple independent agents that collectively optimize a shared objective. Given the prevalence of cooperative systems in real-world applications (Rosa et al., 2024), this remains an active research area in MARL. However, few studies have addressed the collective optimization of risk-sensitive metrics. Qiu et al. (2021) and Shen et al. (2023) investigate cooperative MARL with objective functions such as CVaR and weighted quantile. Their approaches focus on the design of neural network architectures to approximately satisfy certain optimization conditions, thus lacking theoretical analysis. Moreover, both of their algorithms follow the CTDE framework, which requires centralized coordination during the training phase and suffers from high computational burdens (Chen & Cassandras, 2025). To the best of our knowledge, no effective algorithms with theoretical guarantees exist for fully distributed cooperative MARL under risk-sensitive objectives.

# 3 Problem description and modeling

In this section, we first briefly introduce the MMS under the distributed EMS scheme. Subsequently, we model the economic and reliable energy management of MMSs as an MV-TSG, and provide details for variable definitions and corresponding constraints.

## 3.1 Multi-microgrid system

In this study, we investigate an MMS under the distributed EMS scheme, as illustrated in Figure 2. The MMS consists of a set $\mathcal{N} = \{1, 2, \ldots, N\}$ of $N$ microgrids. Each microgrid may include renewable energy generators, demand loads, and ESSs. The renewable energy generator, such as a wind turbine or a solar panel, generates power by harvesting energy from renewable resources. Demand loads represent power consumers whose consumption cannot be curtailed. Consequently, microgrids are both energy producers and consumers, with bidirectional power flow. The ESS facilitates energy management by shifting energy across time steps through real-time charging or discharging. For example, when renewable energy generation exceeds the demand load, the ESS stores the surplus energy; conversely, when the demand load surpasses generation, the ESS discharges stored energy to meet the power shortfall. In the distributed EMS scheme, each microgrid operates an individual EMS and can communicate with neighboring microgrids through the communication infrastructure.
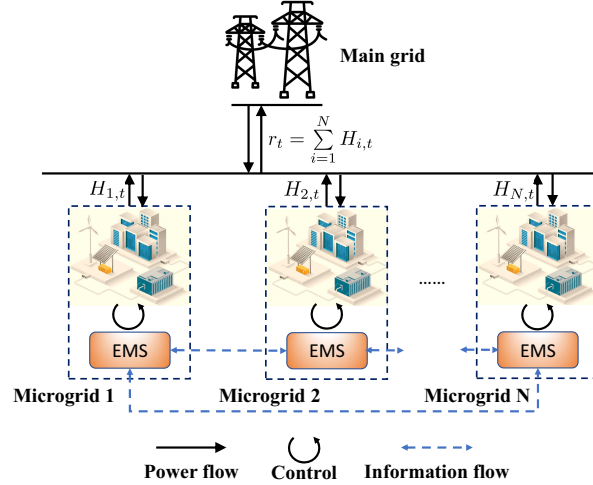


Figure 2: Grid-connected MMS under the distributed EMS scheme

For microgrid $i \in \mathcal{N}$, the maximum output power of the renewable energy generator is denoted by $G_i^{\max}$, and the minimum and maximum consumed power of the demand load is denoted by $D_i^{\min}$ and $D_i^{\max}$, respectively. At time step $t = 1, 2, \ldots$, the renewable energy generated power and demand load of microgrid $i$ are denoted by $G_{i,t}$ and $D_{i,t}$, respectively. Since both depend on random factors such as climate conditions, $G_{i,t}$ and $D_{i,t}$ are non-negative random variables. Each microgrid $i$ is also equipped with an ESS subject to minimum and

maximum energy level constraints, denoted by $B_i^{\min}$ and $B_i^{\max}$, respectively. At time step $t$, the energy level of the ESS is denoted by $B_{i,t}$, and the ESS is discharging with power $b_{i,t}$ (where $b_{i,t} < 0$ indicates charging). Due to power loss during the charging and discharging processes, we introduce an efficiency coefficient $\nu$ for both operations. The actual output or absorbed power of the ESS is given by $e_{i,t} = \nu[b_{i,t}]^+ + \frac{1}{\nu}[b_{i,t}]^-$, where $[b_{i,t}]^+ = \max\{0, b_{i,t}\}$ and $[b_{i,t}]^- = \min\{0, b_{i,t}\}$.

We define the discrete sets for the generated power, demand load power, and storage energy level of microgrid $i$ as $\mathcal{G}_i := \{0, \ldots, G_i^{\max}\}$, $\mathcal{D}_i := \{D_i^{\min}, \ldots, D_i^{\max}\}$, and $\mathcal{B}_i := \{B_i^{\min}, \ldots, B_i^{\max}\}$, respectively. For microgrids with renewable energy generators, power curtailment is permitted when necessary. The EMS of each microgrid controls the ESS discharging power $b_{i,t}$ and the curtailed power $v_{i,t}$ from the renewable energy generator. The output power of microgrid $i$ is then given by $H_{i,t} = G_{i,t} - D_{i,t} + e_{i,t} - v_{i,t}$.

Consequently, the total output power of the MMS or the exchange power between the MMS and the main grid is $H_t = \sum\limits_{i=1}^{N} H_{i,t}$. When $H_t > 0$, the MMS is outputting power to the main grid and getting profits. Conversely, when $H_t < 0$, the MMS purchases power from the main grid to meet demand or charge the ESSs, incurring additional costs. Moreover, fluctuations in exchange power may lead to power quality issues and serve as a reliable indicator for MMS operations.

In this work, we consider the economic and reliable energy management problem of MMSs, and use the variance to measure the fluctuations of exchange power. Given that the MMS consists of multiple microgrids with individual energy management policies, we model the MMS energy management problem as a team stochastic game (TSG) due to its advantages in dealing with stochastic dynamic systems with multiple decision makers.

## 3.2 Mean-variance team stochastic game model

A long-run TSG is defined by a tuple $\langle \mathcal{N}, \mathcal{S}, \mathcal{A}, P, r \rangle$. Here, $\mathcal{N} = \{1, \ldots, N\}$ is a finite set of agents, $\mathcal{S}$ is a finite set of system states, $\mathcal{A}$ is the joint action space of all agents, $\mathcal{A}_i$ is the action space of agent $i \in \mathcal{N}$, $P : \mathcal{S} \times \mathcal{A} \mapsto \Delta(\mathcal{S})$ is the transition probability function

($\Delta(\mathcal{S})$ is the probability distributions over $\mathcal{S}$), $r : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$ is a bounded common reward function. In the long-run TSG, the stationary policy for each agent $i \in \mathcal{N}$ is a mapping $\mu_i : \mathcal{S} \mapsto \Delta(\mathcal{A}_i)$ and the corresponding stochastic policy space is denoted as $\mathcal{U}_i$. The joint policy is denoted by $\boldsymbol{\mu} = (\mu_1, \ldots, \mu_N)$ and its corresponding policy space is $\mathcal{U} = \prod\limits_{i=1}^{N} \mathcal{U}_i$.

The dynamic decision-making process of multiple agents is stated as follows. At each time step $t$, each agent selects an action $a_{i,t} \in \mathcal{A}_i$ according to the system state $s_t$ and its randomized policy function $\mu_i(\cdot|s_t)$. Subsequently, the joint action $\boldsymbol{a_t} = (a_{1,t}, \ldots, a_{N,t})$ lead to a common reward $r_{t+1} = r(s_t, \boldsymbol{a_t})$ to all agents and induces the system's transition to a next state $s_{t+1}$ according to the state transition function $P(s_{t+1}|s_t, \boldsymbol{a_t})$. The entire decision-making process generates a trajectory $\tau = \{s_0, \boldsymbol{a_0}, r_1, s_1, \boldsymbol{a_1}, r_2, s_2, \ldots\}$.

The variable definitions for the TSG in the context of specific MMS energy management problems are described below.

- Number of agents $N$: Each microgrid corresponds to an agent, and $N$ denotes the total number of microgrids in the MMS.

- System state $s_t$: In MMSs under the distributed EMS scheme, each microgrid can communicate with each other, and the state variables of MMSs are fully observable for each microgrid. The MMS state should include the state variables of all microgrids, including the output power of the renewable energy generator, the demand load power, and the storage energy level. Then, the system state is given as $s_t = \{G_{1,t}, D_{1,t}, B_{1,t}, \ldots, G_{N,t}, D_{N,t}, B_{N,t}\}$.

- Action $a_{i,t}$: At each time step, the EMS of each microgrid $i$ can conduct energy management by controlling the output power $b_{i,t}$ of ESSs and the curtailed power $v_{i,t}$ of renewable energy generators. Then, the action is given as $a_{i,t} = (b_{i,t}, v_{i,t})$.

- State transition $P$: We assume that the variables $G_{i,t}$ and $D_{i,t}$ of each microgrid $i$ follow independent and identically distributed random processes, which can be described by Markov chains. This assumption generally holds when the microgrids are located relatively far from each other (Etesami et al., 2018). The state transition of ESSs is

given by $B_{i,t+1} = B_{i,t} - b_{i,t}$, where the next storage energy level depends on the action taken at the current time step. Then, the state transition $P$ of the MMS is determined by the randomness in renewable energy generation and demand loads, as well as the behavior policy.

- Reward $r(s_t, \boldsymbol{a}_t)$: In this study, we define the common reward as the exchange power between the MMS and the main grid, i.e., $r(s_t, \boldsymbol{a}_t) = H_t = \sum_{i=1}^{N}(G_{i,t} - D_{i,t} + e_{i,t} - v_{i,t})$. A positive reward, $r(s_t, \boldsymbol{a}_t) > 0$, indicates that the MMS is supplying power to the main grid and generating profit.

During the MMS operation, certain constraints must be satisfied in the dynamic operation of MMSs. Specifically, for $t = 1, 2, \ldots$, the ESS should meet the capacity constraints: $B_i^{\min} \leq B_{i,t} \leq B_i^{\max}$ and $B_{i,t} - B_i^{\max} \leq b_{i,t} \leq B_{i,t} - B_i^{\min}$. Moreover, the storage output power should also satisfy the power constraint of the interface inverters of the ESS: $-C_i^{\mathrm{ch}} \leq b_{i,t} \leq C_i^{\mathrm{dis}}$, where $C_i^{\mathrm{ch}}$ and $C_i^{\mathrm{dis}}$ denote the maximum charging and discharging power, respectively. Besides, the curtailed power of renewable energy generators should be no more than the power that can be generated, $0 \leq v_{i,t} \leq G_{i,t}$.

In this study, we consider the economic and reliable energy management problem of MMSs. For the economic aspect, the long-run average reward is regarded as the objective function,

$$\eta^{\boldsymbol{\mu}} := \lim_{T \to \infty} \frac{1}{T} \mathbb{E}_{\boldsymbol{\mu}} \Big[ \sum_{t=0}^{T-1} r(s_t, \boldsymbol{a}_t) \Big],$$

which reflects the electricity sales revenue or procurement cost during the MMS long-term operation. With $\pi^{\boldsymbol{\mu}}$ denoting the steady state distribution under the joint policy $\boldsymbol{\mu}$, it is convenient to rephrase the long-run average reward as

$$\eta^{\boldsymbol{\mu}} = \mathbb{E}_{s \sim \pi^{\boldsymbol{\mu}}, \boldsymbol{a} \sim \boldsymbol{\mu}}[r(s, \boldsymbol{a})]$$
$$= \sum_{s \in \mathcal{S}} \pi^{\boldsymbol{\mu}}(s) \sum_{\boldsymbol{a} \in \mathcal{A}} r(s, \boldsymbol{a}) \boldsymbol{\mu}(\boldsymbol{a}|s). \tag{1}$$

For the system operational reliability, we introduce the long-run variance of exchange power,

defined as

$$\zeta^{\boldsymbol{\mu}} := \lim_{T \to \infty} \frac{1}{T} \mathbb{E}_{\boldsymbol{\mu}} \left[ \sum_{t=0}^{T-1} (r(s_t, \boldsymbol{a}_t) - \eta^{\boldsymbol{\mu}})^2 \right],$$

which describes the fluctuation of exchange power in a long-term perspective. Since a power grid needs to maintain a strict balance between power supply and demand, large fluctuations in exchange power will degrade the stability and safety of the power system.

Then, the MV-TSG aims to maximize the expected mean-variance combined metric

$$\begin{aligned} J^{\boldsymbol{\mu}} &= \lim_{T \to \infty} \frac{1}{T} \mathbb{E}_{\boldsymbol{\mu}} \left\{ \sum_{t=0}^{T-1} [r(s_t, \boldsymbol{a}_t) - \beta(r(s_t, \boldsymbol{a}_t) - \eta^{\boldsymbol{\mu}})^2] \right\} \\ &= \eta^{\boldsymbol{\mu}} - \beta \zeta^{\boldsymbol{\mu}}, \end{aligned} \tag{2}$$

where $\beta \geq 0$ is the parameter for the trade-off between mean and variance. Next, we propose an independent policy gradient method and a DRL approach to address MV-TSGs.

# 4   Independent policy gradient

In this section, we first present the value function definitions and corresponding policy gradient results for MV-TSGs. Subsequently, we give the exact policy gradient method for MV-TSGs and analyze its global convergence properties.

## 4.1   Value functions and properties

In this work, we investigate the independent policy gradient method in the energy management problem. We consider direct decentralized policy parameterization, where each agent's policy is parameterized by $\theta_i$,

$$\mu_i^{\theta_i}(a_i|s) = \theta_{i,s,a_i}, \quad i = 1, 2, \ldots, N.$$

For any $s \in \mathcal{S}$ and $a_i$, we have $\theta_{i,s,a_i} > 0$ and $\sum_{a_i \in \mathcal{A}_i} \theta_{i,s,a_i} = 1$. For notational simplicity, we abbreviate $\mu_i^{\theta_i}(a_i|s)$ as $\mu^{\theta_i}(a_i|s)$, $\theta_{i,s,a_i}$ as $\theta_{s,a_i}$. The joint policy $\boldsymbol{\mu}^{\theta}(\boldsymbol{a}|s) = \prod_{i=1}^{N} \mu^{\theta_i}(a_i|s) = $

$\prod\limits_{i=1}^{N} \theta_{s,a_i}$. We use $\mathcal{X}_i := \Delta(\mathcal{A}_i)^{|S|}$, $\mathcal{X} := \mathcal{X}_1 \times \cdots \times \mathcal{X}_N$ denote the feasible region of $\theta_i$ and $\theta = (\theta_1, \ldots, \theta_N)$, respectively. In the rest of the paper, we use the notation $-i$ to denote the set of all agents except agent $i$.

Inspired by Equation (2), we provide a surrogate reward function as

$$f^\theta(s, \boldsymbol{a}) = r(s, \boldsymbol{a}) - \beta(r(s, \boldsymbol{a}) - \eta^\theta)^2. \tag{3}$$

Subsequently, similar to Equation (1), the mean-variance performance can be calculated by

$$J^\theta = \sum_{s \in \mathcal{S}} \pi^\theta(s) \sum_{\boldsymbol{a} \in \mathcal{A}} f^\theta(s, \boldsymbol{a}) \boldsymbol{\mu}^\theta(\boldsymbol{a}|s).$$

Following the definitions of average-reward MDPs in Sutton & Barto (2018), the value function $V_f^\theta$, action-value function $Q_f^\theta$ and advantage function $A_f^\theta$, with respective to the joint policy and surrogate reward function $f^\theta$, are defined as

$$V_f^\theta(s) := \mathbb{E}_{\boldsymbol{\mu}^\theta}\left[\sum_{t=0}^{\infty}(f^\theta(s_t, \boldsymbol{a}_t) - J^\theta)|s_0 = s\right],$$

$$Q_f^\theta(s, \boldsymbol{a}) := \mathbb{E}_{\boldsymbol{\mu}^\theta}\left[\sum_{t=0}^{\infty}(f^\theta(s_t, \boldsymbol{a}_t) - J^\theta)|s_0 = s, \boldsymbol{a}_0 = \boldsymbol{a}\right]$$

$$= f^\theta(s, \boldsymbol{a}) - J^\theta + \sum_{s' \in \mathcal{S}} P(s'|s, \boldsymbol{a})V_f^\theta(s'),$$

$$A_f^\theta(s, \boldsymbol{a}) := Q_f^\theta(s, \boldsymbol{a}) - V_f^\theta(s).$$

We now present the performance difference lemma in MV-TSGs, which quantifies the mean-variance performance difference between any two joint policies. Lemma 1 is a direct extension of the mean-variance difference formula for MDPs, as presented in Xia (2020), to the TSG setting, with the proof omitted for brevity. For clarity, the notations $J(\boldsymbol{\mu})$, $J^{\boldsymbol{\mu}}$, $J(\theta)$ and $J^\theta$ are used interchangeably when necessary.

**Lemma 1** (Performance difference in MV-TSGs). *For any two joint policies $\boldsymbol{\mu}, \boldsymbol{\mu}' \in \mathcal{U}$, the*

*mean-variance performance difference is*

$$J(\boldsymbol{\mu}') - J(\boldsymbol{\mu}) = \mathbb{E}_{s \sim \pi^{\boldsymbol{\mu}'}, \boldsymbol{a} \sim \boldsymbol{\mu}'}[A_f^{\boldsymbol{\mu}}(s, \boldsymbol{a})] + \beta(\eta^{\boldsymbol{\mu}'} - \eta^{\boldsymbol{\mu}})^2.$$

Different from the standard discounted or averaged reward cases, where the reward function $r$ is independent of policy $\mu$, the surrogate reward function $f$ defined in (3) includes a term of $\eta$, which depends on the joint policy. Then, a new policy gradient lemma is derived for MV-TSGs based on Lemma 1 and stated as follows.

**Lemma 2** (Mean-variance policy gradient)**.** *In MV-TSGs, for a joint policy parameterized by $\theta$, we have*

$$\nabla_\theta J(\theta) = \mathbb{E}_{s \sim \pi^\theta, \boldsymbol{a} \sim \mu^\theta}[\nabla_\theta \log \mu^\theta(\boldsymbol{a}|s) Q_f^\theta(s, \boldsymbol{a})].$$

The proof is provided in Appendix A.1. This result is similar to the well-known policy gradient theorem in the average-reward setting (Sutton et al., 1999), except that the action-value function is computed with respect to the surrogate reward function $f$ instead of $r$. Based on Lemma 2, we further derive the partial derivative formula of the mean-variance performance function. The proof is provided in Appendix A.2.

**Lemma 3** (Performance partial derivative)**.** *For any agent $i$ and any directly parameterized policy $\theta_i \in \mathcal{X}_i$, we have*

$$\frac{\partial J(\theta)}{\partial \theta_{s,a_i}} = \overline{Q_{f,i}^\theta}(s, a_i) \pi^\theta(s),$$

*where* $\overline{Q_{f,i}^\theta}(s, a_i) := \sum\limits_{\boldsymbol{a}_{-i}} \boldsymbol{\mu}^{\theta_{-i}}(\boldsymbol{a}_{-i}|s) Q_f^\theta(s, a_i, \boldsymbol{a}_{-i}).$

Based on Lemma 2 and Lemma 3, we now present the independent policy gradient method and analyze its convergence properties.

## 4.2 Exact policy gradient ascent

Because for each agent $i$, the policy parameters should satisfy the constraint of $\sum\limits_{a_i \in \mathcal{A}_i} \theta_{s,a_i} = 1, \forall s \in \mathcal{S}$, we propose an independent projected gradient ascent to update policies in MV-

TSGs, i.e., MV-IPGA. Specifically, each agent update its policy along the gradient direction and use the operator $\mathrm{Proj}_{\mathcal{X}_i}(\theta_i) := \arg\min_{x \in \mathcal{X}_i} \|x - \theta_i\|$ to project the updated parameters onto the feasible region $\mathcal{X}_i$,

$$\theta_i' = \mathrm{Proj}_{\mathcal{X}_i}(\theta_i + \alpha \nabla_{\theta_i} J(\theta_i, \theta_{-i})), \quad \alpha > 0.$$

In MV-TSGs, we have the following lemma. Similar results have also been provided in Leonardos et al. (2022)(Proposition B.1) and Zhang et al. (2024)(Proposition 1).

**Lemma 4.** *Let $\theta = (\theta_1, \ldots, \theta_N)$, $\theta' = \theta + \alpha \nabla_\theta J(\theta)$, where $\alpha$ is the step size. Then, we have*

$$\mathrm{Proj}_{\mathcal{X}}(\theta') = (\mathrm{Proj}_{\mathcal{X}_1}(\theta_1'), \ldots, \mathrm{Proj}_{\mathcal{X}_N}(\theta_N')).$$

Lemma 4 indicates that when all agents adopt the same step size for policy updates, the independent projected gradient ascent is equivalent to performing projected gradient ascent on the joint policy. Based on this observation, we propose the MV-IPGA algorithm, as detailed in Algorithm 1.

---

**Algorithm 1** Mean-variance independent projected gradient ascent for MV-TSGs

---

1: **Input:** step size $\alpha > 0$, number of iterations $K$.

2: **Initialization:** for each agent $i$, randomly initialize the parameters $\theta_i^{(0)}$ of the policy $\mu_i^{(0)}$.

3: **for** $k = 0, 1, \ldots, K - 1$ **do**

4: $\quad \theta_i^{(k+1)} = \mathrm{Proj}_{\mathcal{X}_i}(\theta_i^{(k)} + \alpha \nabla_{\theta_i} J(\theta_i^{(k)}, \theta_{-i}^{(k)})), \quad \forall i.$

5: **end for**

---

However, as demonstrated by Zhong et al. (2024), independent and simultaneous policy updates by all agents may fail to guarantee performance improvement or algorithmic convergence, due to the issue of environmental non-stationarity in multi-agent settings. In light of this, we provide a convergence analysis of MV-IPGA.

## 4.3 Global convergence of policy gradient in MV-TSGs

Based on Lemma 4, we analyze the convergence of MV-IPGA by following the standard approach for gradient-based methods in nonconvex optimization. The analysis proceeds in two steps: (1) establishing the smoothness properties of the mean-variance performance function, and (2) proving the convergence of MV-IPGA. Without loss of generality, we assume throughout this section that the rewards are normalized such that $r \in [0, 1]$.

We first analyze the smoothness properties of the mean-variance performance function, as Lemma 5 demonstrates. The proof is provided in Appendix A.3.

**Lemma 5.** *Denote* $A_{\max} := \max_i |\mathcal{A}_i|$, $S = |\mathcal{S}|$, $\kappa_0$ *is the mixing coefficient defined in Cheng et al. (2024),* $L := 6\beta A_{\max}(1 + \frac{\kappa_0}{2}S)^2 + (1 + \beta)\kappa_0 A_{\max}\sqrt{S}(\kappa_0 S + 1)$ *and* $L_J := N\left(6\beta A_{\max}(1 + \frac{\kappa_0}{2}S)^2 + (1 + \beta)A_{\max}(\kappa_0 S + \kappa_0^2 S^{\frac{3}{2}} + \kappa_0\sqrt{S} + 1)\right)$, *we have*

- *For any agent* $i$ *and* $\mu_{-i} \in \mathcal{U}_{-i}$, *the mean-variance performance function* $J^\theta$ *is* $L$-*smooth with respect to policy* $\theta_i$, *i.e.,* $\|\nabla_{\theta_i} J(\theta_i, \theta_{-i}) - \nabla_{\theta_i} J(\theta_i', \theta_{-i})\|_2 \leq L\|\theta_i - \theta_i'\|$, $\forall i$ *and* $\theta_i, \theta_i' \in \mathcal{X}_i$.

- *The mean-variance performance function* $J^\theta$ *is* $L_J$-*smooth with respect to the joint policy* $\theta$, *i.e.,* $\|\nabla_\theta J(\theta) - \nabla_\theta J(\theta')\|_2 \leq L_J\|\theta - \theta'\|$.

Compared with the smoothness results of the average-reward objective $\eta$ in Cheng et al. (2024), it is evident that the incorporation of variance complicates the smoothness of the objective function, which is also influenced by the coefficient $\beta$.

Based on the results of the smoothness properties above, we now analyze the convergence of MV-IPGA. To quantify the convergence behavior of the algorithm, we introduce the function $\mathrm{ST}(\theta) = \max_i \max_{\theta_i'} (\theta_i' - \theta_i)^\top \nabla_{\theta_i} J(\theta)$, which measures the maximal directional improvement over all agents. Since both the long-run average reward $\eta^{\boldsymbol{\mu}}$ and long-run variance $\zeta^{\boldsymbol{\mu}}$ are bounded due to the boundedness of the reward function $r$, we denote the maximum and minimum of the mean-variance performance function by $J_{\max}$ and $J_{\min}$. The convergence result of MV-IPGA is stated in Theorem 1, and its proof is provided in Appendix A.4.

**Theorem 1.** *If all agents follow Algorithm 1 with step size $\alpha = \frac{1}{L_J}$, then the joint policy asymptotically converges to a first-order stationary point, i.e., $\lim\limits_{k \to \infty} ST(\theta^{(k)}) \leq 0$. Moreover, if all agents independently play gradient ascent for at least $K \geq \frac{4L_J(J_{\max} - J_{\min})}{\epsilon}$ iterations, then there exists a $k \in \{1, \ldots, K\}$ such that $\theta^{(k)}$ is $\epsilon$-stationary, i.e., $ST(\theta^{(k)}) \leq \epsilon$.*

Theorem 1 shows that MV-IPGA converges within a finite number of iterations when all agents adopt an identical and theoretically valid step size. However, such theoretically guaranteed step sizes are often overly conservative and may be impractical for specific problem instances. In practice, the step size is typically chosen empirically to balance convergence stability and speed.

# 5   Independent deep reinforcement learning for MV-TSGs

In the previous section, we introduced the MV-IPGA algorithm and analyzed its convergence properties. However, applying MV-IPGA to practical energy management in MMSs may encounter two limitations. First, policy gradient methods typically require accurate knowledge of environmental parameters, such as the reward function and state transition dynamics. Second, as the number of microgrids increases, the state space of the system grows exponentially, leading to significant computational and memory challenges for policy optimization based on direct policy parameterization.

To approximately solve the energy management problem in large-scale MMSs under scenarios with unknown model parameters, we further propose an independent DRL approach for MV-TSGs, built upon the independent policy gradient algorithm. In DRL, both the policy and value functions are parameterized via neural networks, which enables the framework to effectively handle problems with high-dimensional state and action spaces.

In the single-agent setting, PPO has demonstrated remarkable success in practical applications, including OpenAI Five (Berner et al., 2019) and generative models such as Chat-GPT. In practical applications, directly applying policy gradient methods may result in

excessively large policy updates due to potential parameter estimation errors, thereby desta-bilizing the learning process. Like other first-order policy gradient algorithms, PPO employs gradient information to guide optimization. The principal difference from methods such as A2C (Mnih et al., 2016), TD3 (Fujimoto et al., 2018), and SAC (Haarnoja et al., 2018) is its conservative update mechanism, which maximizes a clipped surrogate objective to con-strain policy updates within a local trust region, thereby enhancing stability. Empirical studies have shown that this approach achieves strong performance across a wide range of tasks, while offering several practical advantages such as ease of implementation, high sample efficiency, and robustness to hyper-parameter settings (Schulman et al., 2017).

Following the route of PPO, we propose a corresponding independent DRL algorithm based on Algorithm 1, referred to as MV-IPPO. Unlike standard PPO, which maximizes the discounted cumulative reward, in MV-IPPO, each agent performs policy optimization based on its surrogate action-value function $\overline{Q_{f,i}^{\theta}}(s, a_i)$ introduced in Lemma 3 or the advantage function $\overline{A_{f,i}^{\theta}}(s, a_i) := \sum_{\boldsymbol{a}_{-i}} \boldsymbol{\mu}^{\theta_{-i}}(\boldsymbol{a}_{-i}|s) A_f^{\theta}(s, a_i, \boldsymbol{a}_{-i})$. At the $k$th iteration, the policy of each agent and the common value function are parameterized by neural networks with parameters $\theta_i^{(k)}$ (actor network) and $\phi^{(k)}$ (critic network), respectively. Then, each agent optimizes its policy by maximizing the following surrogate objective function

$$\mathcal{L}_i^{\boldsymbol{\mu}}(\theta_i) := \mathbb{E}_{s \sim \pi^{\boldsymbol{\mu}}, a_i \sim \mu_i} \left[ \min \left( \omega_i(\theta_i) \hat{A}_{f,i}^{\theta^{(k)}}(s, a_i), \mathrm{clip}(\omega_i(\theta_i), 1 - \varepsilon, 1 + \varepsilon) \hat{A}_{f,i}^{\theta^{(k)}}(s, a_i) \right) \right],$$

where $\omega_i(\theta_i) = \frac{\mu^{\theta_i}(a_i|s)}{\mu^{\theta_i^{(k)}}(a_i|s)}$, and $\hat{A}_{f,i}^{\theta^{(k)}}$ is the estimation of $\overline{A_{f,i}^{\theta^{(k)}}}$.

For each agent, the advantage function $\hat{A}_{f,i}$ can be estimated using the generalized ad-vantage estimation method (Schulman et al., 2016). Specifically, we have

$$\hat{A}_{f,i}^{\theta^{(k)}}(s_n, a_{i,n}) = \sum_{t=n}^{T-1} \lambda_{t-n} (\hat{f}^{\theta^{(k)}}(s_t, a_{i,t}) - \hat{J}^{\theta^{(k)}} + V_f^{\phi^{(k)}}(s_t) - V_f^{\phi^{(k)}}(s_{t+1})),$$

where $\lambda$ is the hyper-parameter to trade-off bias and variance, and $T$ denotes the trajectory length.

Since we consider the long-run average performance in this study, we adopt the average value constraint (AVC) proposed by Ma et al. (2021) to assist in estimating the target value

function $\hat{V}_f^{\phi^{(k)}}$ and stabilizing the value learning. The value function network is updated using the loss function

$$\frac{1}{T}\sum_{t=0}^{T-1}\left(V_f^{\phi}(s_t) - \hat{V}_f^{\phi^{(k)}}(s_t)\right)^2.$$

The detailed pseudo-code of MV-IPPO is presented in Algorithm 2, where $(\cdot)_{t,m}$ denotes the corresponding variable value at time step $t$ in the $m$-th trajectory.

---

**Algorithm 2** Independent proximal policy optimization for MV-TSGs

---

1: **Input:** threshold hyper-parameter $\epsilon > 0$, step size $\alpha$, episode length $T$.
2: **Initialization:** for agent $i$, randomly initialize the parameters $\theta_i^{(0)}$ of the policy $\mu_i^{(0)}$ and the parameters $\phi^{(0)}$ of the value function, replay buffer $\mathcal{M}$. $M$ is the total number of trajectories collected in each iteration. Set $\hat{\eta} = 0, \hat{\zeta} = 0, \hat{J} = 0$.
3: **for** $k = 0, 1, \dots$ **do**
4:      Collect a set of trajectories by running policy $\mu^{\theta_i}$ in the environment.
5:      Push transitions $\{(s_{t,m}, a_{i,t,m}, s_{t+1,m}, r_{t,m}), t \in \{0, \dots, T-1\}, m \in \{1, \dots, M\}\}$ into $\mathcal{M}$.
6:      Update $\hat{\eta} \leftarrow (1-\alpha)\hat{\eta} + \alpha\frac{1}{MT}\sum\limits_{m=1}^{M}\sum\limits_{t=0}^{T-1} r_{t,m}$.
7:      Update $\hat{\zeta} \leftarrow (1-\alpha)\hat{\zeta} + \alpha\frac{1}{MT}\sum\limits_{m=1}^{M}\sum\limits_{t=0}^{T-1} (r_{t,m} - \hat{\eta})^2$.
8:      Compute the average mean-variance performance function $\hat{J}$.
9:      Compute the average $\hat{f}^{\theta^{(k)}}(s_t, a_{i,t})$ and $\hat{A}_f^{\theta^{(k)}}(s_t, a_{i,t})$ over all time steps and trajectories.
10:     Compute the average $\hat{V}_f^{\phi^{(k)}}(s_t)$ over all time steps and trajectories using AVC.
11:     Update the policy by maximizing the objective:

$$\theta_i^{k+1} = \arg\max_{\theta_i} \frac{1}{MT}\sum_{m=1}^{M}\sum_{t=0}^{T-1}\min\left(\frac{\mu^{\theta_i}(a_{i,t,m}|s_{t,m})}{\mu^{\theta_i^{(k)}}(a_{i,t,m}|s_{t,m})}\hat{A}_{f,i}^{\theta^k}(s_{t,m}, a_{i,t,m}),\right.$$
$$\left.\text{clip}\left(\frac{\mu^{\theta_i}(a_{i,t,m}|s_{t,m})}{\mu^{\theta_i^{(k)}}(a_{i,t,m}|s_{t,m})}, 1 \pm \epsilon\right)\hat{A}_{f,i}^{\theta^k}(s_{t,m}, a_{i,t,m})\right),$$

     using stochastic gradient ascent with Adam.
12:     Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg\min_{\phi} \frac{1}{MT}\sum_{m=1}^{M}\sum_{t=0}^{T-1}\left(V_f^{\phi}(s_{t,m}) - \hat{V}_f^{\phi^{(k)}}(s_{t,m})\right)^2,$$

     using stochastic gradient descent with Adam.
13: **end for**

---

# 6 Numerical experiments

This section first introduces the experimental settings in MMSs. We then evaluate the effectiveness of MV-IPGA under varying values of $\beta$ in a simplified scenario. Finally, MV-IPPO is applied to address the energy management problem in a larger-scale setting with unknown model parameters.

## 6.1 Experimental settings

We investigate an MMS under the distributed EMS scheme, specifying the configurations of renewable generators, demand loads, and ESSs based on typical equipment and power demand profiles (Yang et al., 2020). For simplicity, we assume that the equipment in different microgrids is identical and operates independently. All state and action variables are discretized using uniform quantization.

(1) Renewable energy generator. A wind turbine serves as a representative renewable energy generator in this study. The wind power $G(t)$ are calculated by

$$
G(t) = \begin{cases} W^{\text{cap}}, & V^{\text{rated}} \leq v_t < V^{\text{cutout}}, \\ W^{\text{cap}} \left( \frac{v_t}{V^{\text{rated}}} \right)^3, & V^{\text{cutin}} \leq v_t < V^{\text{rated}}, \\ 0, & \text{others}, \end{cases}
$$

where $v_t$ represents the wind speed at time step $t$, $V^{\text{rated}}$, $V^{\text{cutin}}$ and $V^{\text{cutout}}$ indicate the rated, cut-in and cut-out wind speed, respectively. $W^{\text{cap}}$ denotes the rated wind power. The parameter settings of the wind turbine are presented in Table 1.

Table 1: Wind turbine configurations.

| Symbol | Description | Setting |
|--------|-------------|---------|
| $V^{\text{cutin}}$ | Cut-in wind speed | 4 m/s |
| $V^{\text{cutout}}$ | Cut-out wind speed | 25 m/s |
| $V^{\text{rated}}$ | Rated wind speed | 15 m/s |
| $W^{\text{cap}}$ | Rated wind power | 3 MW |

Table 2: States of wind power.

| State | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| Wind power/MW | 0 | 0.6 | 1.2 | 1.8 | 2.4 | 3.0 |

Table 3: States of demand loads.

| State | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| Demand load power/MW | 0.6 | 1.2 | 1.8 | 2.4 | 3.0 | 3.6 |

$$\boldsymbol{P}_W = \begin{pmatrix} 0.94 & 0.05 & 0.01 & 0.00 & 0.00 & 0.00 \\ 0.40 & 0.44 & 0.10 & 0.03 & 0.01 & 0.02 \\ 0.16 & 0.37 & 0.26 & 0.11 & 0.05 & 0.05 \\ 0.08 & 0.24 & 0.25 & 0.19 & 0.10 & 0.14 \\ 0.08 & 0.14 & 0.18 & 0.19 & 0.14 & 0.27 \\ 0.04 & 0.07 & 0.08 & 0.10 & 0.10 & 0.61 \end{pmatrix}. \tag{4}$$

The output power of the wind turbine is discretized into six states, as illustrated in Table 2. The wind power and its state transition are determined by the random wind speed. The wind speed data we used were collected by the Measurement and Instrumentation Data Center at the National Renewable Energy Laboratory over the period from 2015 to 2024 (NREL, 2025). Accordingly, the probability transition matrix $\boldsymbol{P}_W$ in (4) is constructed. The $(k, l)^{\text{th}}$ element $p_{k,l}$ in $\boldsymbol{P}_W$ indicates the probability that the wind power from state $k$ to state $l$, and is calculated by $P_{k,l} = \frac{q_{k,l}}{q_k}$, where $q_{k,l}$ represents the observed transitions from state $k$ to state $l$, and $q_k$ represents the total occurrence of state $k$.

(2) Demand load. The demand load data is derived from the 2023 data in a publicly available database (IESO, 2023). The database is maintained by the independent electricity system operator (IESO), which is a non-profit corporate entity established in 1998 by the Electricity Act of Ontario. The dataset employed records the hourly demand load in Ontario across the entire year. With this demand load dataset, we divide the demand load power into six states and construct the probability transition matrix $\boldsymbol{P}_D$, as illustrated in Table 3 and (5), in the same way as in the wind power.

$$\boldsymbol{P}_D = \begin{pmatrix} 0.75 & 0.25 & 0.00 & 0.00 & 0.00 & 0.00 \\ 0.03 & 0.83 & 0.14 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.11 & 0.82 & 0.07 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.14 & 0.84 & 0.02 & 0.00 \\ 0.00 & 0.00 & 0.00 & 0.19 & 0.79 & 0.02 \\ 0.00 & 0.00 & 0.00 & 0.00 & 0.27 & 0.73 \end{pmatrix}. \tag{5}$$

(3) Energy storage system. The parameters of the ESS are presented in Table 4. The energy level is divided into six states, as illustrated in Table 5. The action space for discharging power is also discretized properly and illustrated in Table 6.

Table 4: Configurations of energy storage systems.

| Symbol | Description | Setting |
|---|---|---|
| $C^{\text{ch}}$ | Maximum charging power | 1.2 MW |
| $C^{\text{dis}}$ | Minimum discharging power | 1.2 MW |
| $B^{\text{cap}}$ | Rated capacity | 4.0 MWh |
| $B^{\text{min}}$ | Energy level lower bound threshold | 0.6 MWh |
| $B^{\text{max}}$ | Energy level upper bound threshold | 3.6 MWh |
| $\nu$ | Charging/discharging efficiency coefficient | 0.95 |

Table 5: States of storage energy levels.

| State | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| Energy level/MWh | 0.6 | 1.2 | 1.8 | 2.4 | 3.0 | 3.6 |

Table 6: Scheduling actions of energy storage systems.

| Action $(b_{i,t})$ | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| Discharging power/MW | -1.2 | 0.6 | 0.0 | 0.6 | 1.2 |

We conduct a series of experiments to validate the effectiveness of the MV-IPGA and MV-IPPO algorithms in Section 6.2 and Section 6.3, respectively. All experiments are conducted on a machine equipped with one AMD 3995WX CPU, 384GB of memory, and one Nvidia GeForce GTX4090 GPU. We note that the corresponding algorithms are evaluated under

24

different values of $\beta$, and compared against a baseline policy without energy management, due to the lack of existing methods that can address the economic and reliable MMS energy management under the distributed EMS scheme.

## 6.2 Policy gradient in a small-scale case

In this section, we use the MV-IPGA method to deal with the economic and reliable energy management problem when the environmental parameters are known exactly. Due to the state space increasing exponentially with the number of microgrids, we focus on a small-scale MMS consisting of two microgrids. Each microgrid is equipped with an ESS, while only Microgrid 1 is configured with a renewable energy generator and demand loads. The step size $\alpha$ in the policy gradient method is set to 0.5.

Figure 3 illustrates the convergence process of the mean, variance, and mean-variance of exchange power under MV-IPGA. In the first two subplots, the blue dashed lines indicate the corresponding values under a baseline policy without energy management, with a mean of -1.52 and a variance of 0.57. Since energy management involves ESS scheduling (with efficiency coefficient $\nu < 1$) or power curtailment, the blue dashed line in the first subplot represents the theoretical upper bound of the mean exchange power. To evaluate the effectiveness of MV-IPGA, we conduct experiments under different values of $\beta$: 0.0, 0.3, and 1.0. As shown in the figure, when $\beta = 0.0$, the mean exchange power gradually converges to -1.52. In this case, the algorithm does not account for power fluctuation minimization, and the variance remains above the baseline level. We note that the variance decreases because the mean optimization avoids unnecessary energy management actions.

When $\beta = 0.3$, the power fluctuation is taken into account, and the converged variance of exchange power reduces significantly to 0.38, though the converged mean value incurs a slight decrease due to energy losses associated with storage scheduling and curtailment. When $\beta = 1.0$, more considerations are placed on reducing power fluctuations, further adjusting both the variance and mean.

The convergence curves of the mean-variance objective function in the third subplot

25

clearly show that the proposed objective is progressively optimized as the number of algorithm iterations increases, demonstrating the effectiveness of MV-IPGA.
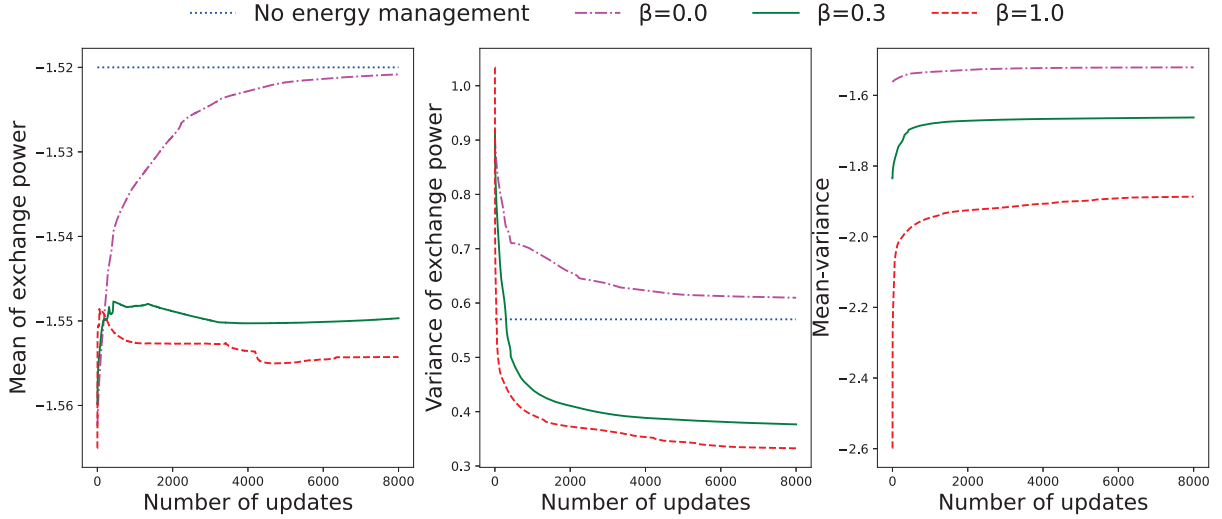


Figure 3: Convergence procedure of MV-PG with different coefficients.

To provide a more intuitive demonstration of the superiority of the energy management policy learned by the algorithm, a sample trajectory of exchange power over 72 time steps is plotted in Figure 4. It is illustrated that the exchange power exhibits significant fluctuations when no energy management is applied (blue curve) or when power fluctuation is not considered in the optimization (purple curve). In contrast, when $\beta$ is set to 0.3 and 1.0, the fluctuations in exchange power are effectively mitigated to varying degrees.

Furthermore, Figure 5 presents the state and action details of a representative trajectory under the setting $\beta = 1.0$. The net generated power is defined as the difference between the generated wind power and the demand load, i.e., $G_{1,t} - D_{1,t}$. It can be observed that the net generated power remains negative for most time steps, indicating that the demand load generally exceeds the available wind power in MMS operations. Microgrids only take action in states where significant fluctuations are likely to occur, as the discharging and charging of storage, as well as power curtailment, can result in energy losses.

Specifically, when the demand significantly surpasses the generated power, the microgrids compensate for power by discharging ESSs, such as during time steps 31 to 36. Conversely, microgrids tend to charge their ESSs when the net power is only slightly negative or turns
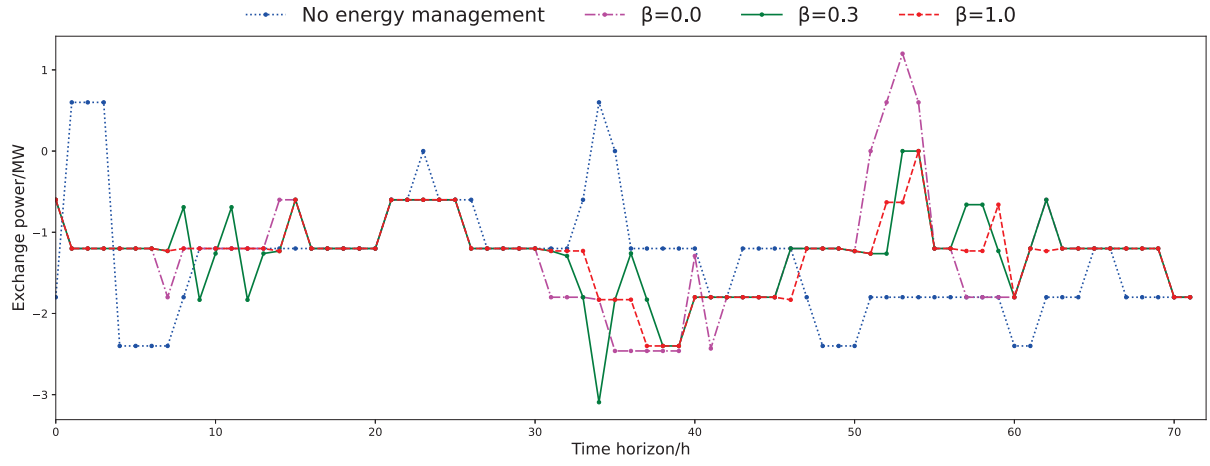
Figure 4: An episode of exchange power over 72 time steps under the MV-IPGA policy.
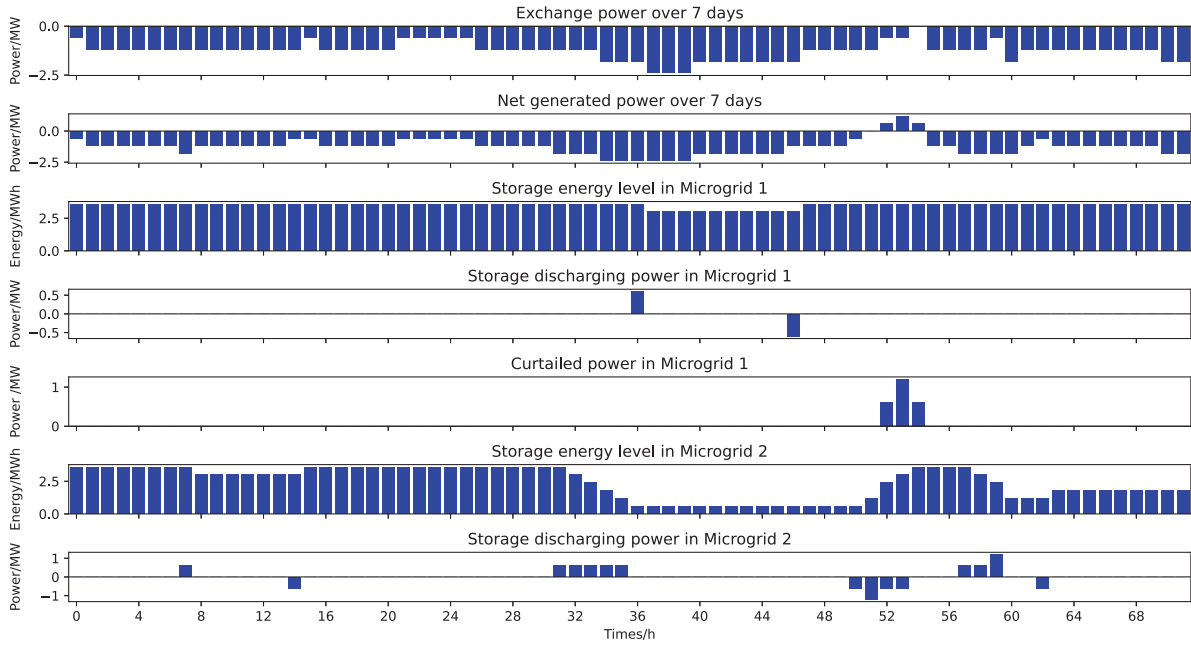


Figure 5: Trajectory details of states and actions when $\beta = 1.0$.

27

positive. Notably, when the net power becomes positive, Microgrid 1 implements power curtailment to prevent reverse power flow and to mitigate power fluctuations. Such behavior can be observed between time steps 50 and 54.

In summary, the results presented in the above figures jointly demonstrate the effectiveness of the proposed MV-IPGA in addressing the trade-off between economic performance and reliability in energy management for MMSs.

## 6.3  Reinforcement learning in a large-scale case

In this section, we use MV-IPPO to tackle the energy management problem when environmental parameters are unknown. The MMS considered consists of three microgrids, all are equipped with renewable energy generators, demand loads, and ESSs.

A simulated environment is constructed to train the energy management policy of microgrids. At each time step, the microgrid takes an action $a_{i,t}$ to interact with the environment and get a reward $r_{i,t}$ feedback. The system state $s_t$ is transitioned to $s_{t+1}$ according to the transition matrix $P_W$, $P_D$, and the joint action $\boldsymbol{a}$. The microgrids iteratively update their energy management policies based on the interaction data with the environment. The primary hyper-parameters of MV-IPPO are presented in Table 7. Both the actor and critic networks comprise two hidden layers, each with 64 units and rectified linear unit (ReLU) activation functions. The network architecture and most hyperparameter settings follow those of the state-of-the-art algorithm in Yu et al. (2022).

The performance of MV-IPPO is evaluated under different values of $\beta$, with the corresponding training curves shown in Figure 6. The blue curves represent the long-term mean and variance of exchange power in the absence of energy management, which are $-4.55$ and $1.70$, respectively. As illustrated in the figure, when $\beta = 0.0$, MV-IPPO focuses solely on optimizing the mean exchange power without accounting for power fluctuations. In this case, the average converged mean across multiple random seeds reaches $-4.59$, which is very close to the optimal value of $-4.55$, while the resulting variance exceeds $1.70$.

When $\beta = 0.3$, the variance of exchange power is explicitly considered. As shown in Figure 6, the variance is significantly reduced upon convergence, indicating effective coordi-

Table 7: Hyper-parameter settings in MV-IPPO

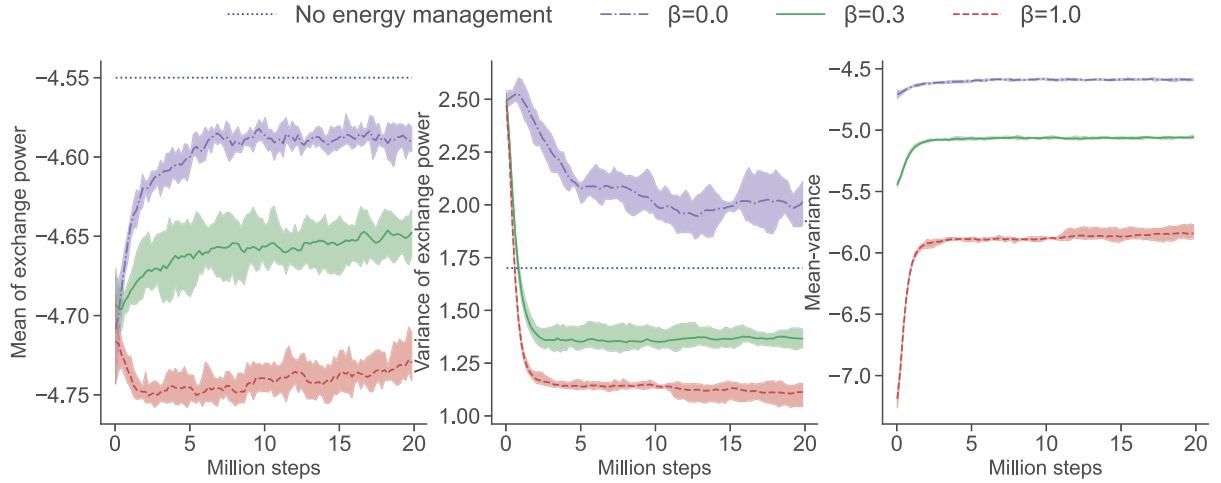| Hyper-parameter | Value |
| --- | --- |
| Number of training steps | 20 million |
| Number of environments collecting data in parallel | 20 |
| Length of time horizon $T$ | 2000 |
| Decay-rate parameter $\lambda$ for eligibility traces | 0.95 |
| Number of mini-bath | 40 |
| Clipping coefficient $\epsilon$ | 0.2 |
| Training epochs | 5 |
| Optimizer for gradient descent/ascent | Adam |
| Learning rate for optimizer | 0.0005 |
| Average value constraint coefficient in AVC | 0.01 |



Figure 6: **Training curves of MV-IPPO with different coefficients.** Each training curve in the figure is averaged over six random seeds and shaded by standard deviations.
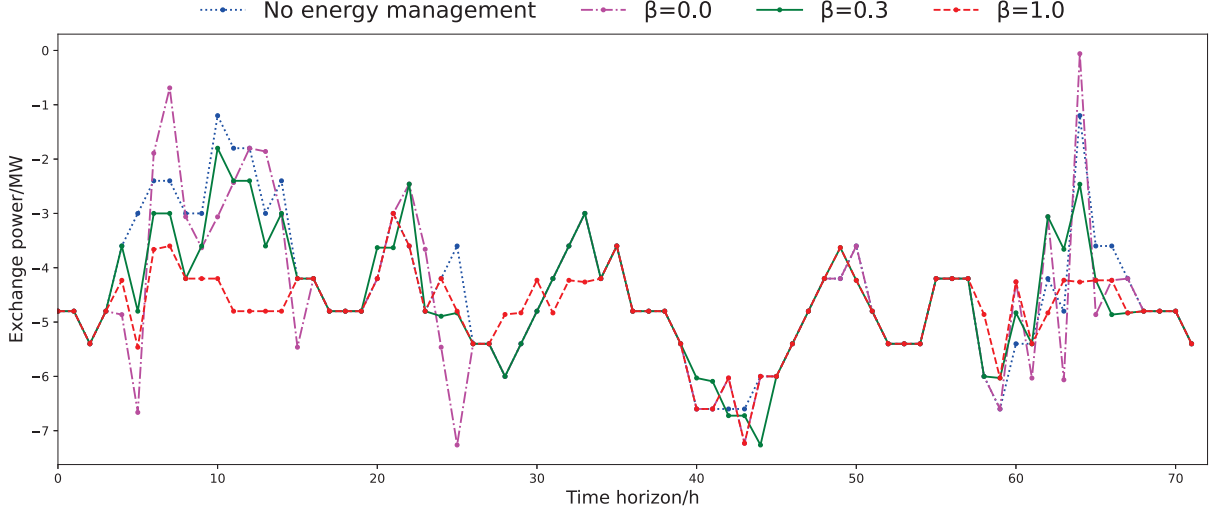
Figure 7: An episode of exchange power over 72 time steps under the MV-IPPO policy.

nation between energy storage and renewable generation resources. However, due to energy losses caused by storage charging/discharging and power curtailment, the mean value also decreases. As $\beta$ increases further to 1.0, both the converged mean and variance are further regulated. The convergence curves of the mean-variance objective confirm that the objective function is effectively optimized as training progresses.

Similarly, we select the best-performing policies across different random seeds for evaluation and generate a representative sample trajectory over 72 time steps, as depicted in Figure 7. It can be intuitively observed from the figure that the fluctuations in exchange power are effectively mitigated as the value of $\beta$ increases.

The results presented in Figure 6 and Figure 7 demonstrate that MV-IPPO effectively solves large-scale MV-TSG problems in a data-driven manner. The derived policies facilitate the joint optimization of economic performance and reliability in MMSs.

# 7 Conclusion and discussion

In this paper, we investigate an economic and reliable energy management problem for MMSs under the distributed EMS scheme. We introduce the mean and variance of exchange power between the MMS and the main grid as indicators of economic performance and reliability,

respectively. The problem is formulated as an MV-TSG. Given the absence of centralized coordination in the distributed EMS scheme, we propose an independent policy gradient method, termed MV-IPGA, to solve the MV-TSG when model parameters are fully known. A rigorous convergence analysis of MV-IPGA is also provided. Furthermore, we develop a DRL algorithm, called MV-IPPO, which extends MV-IPGA to scenarios with unknown model parameters. MV-IPPO enables the approximate solution of large-scale MV-TSGs in a data-driven manner. Experimental results demonstrate that the proposed algorithms effectively solve MV-TSGs and achieve a desirable trade-off between economic performance and reliability in the long-term operation of MMSs.

This study develops a distributed optimization framework for mean–variance objectives. Extending it to alternative risk measures, such as value-at-risk (VaR) and CVaR, is a promising direction for future research. However, these risk metrics often lack analytically tractable formulations for performance difference or gradient computation, which are essential to the proposed algorithm. In practice, challenges such as communication delays and missing data may be mitigated by leveraging recurrent neural networks or attention mechanisms to capture temporal dependencies. Moreover, addressing the sim-to-real gap through robustness- and generalization-oriented techniques is crucial to enhance the real-world applicability of DRL policies.

# Funding Declaration

# References

Abdalla, A. A., El Moursi, M. S., El-Fouly, T. H., & Al Hosani, K. H. (2023). A novel adaptive power smoothing approach for PV power plant with hybrid energy storage system. *IEEE Transactions on Sustainable Energy*, 14(3), 1457–1473.

Agarwal, A., Kakade, S. M., Lee, J. D., & Mahajan, G. (2021). On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 22(98), 1–76.

Alam, M. N., Chakrabarti, S., & Ghosh, A. (2018). Networked microgrids: State-of-the-art and future perspectives. *IEEE Transactions on Industrial Informatics*, 15(3), 1238–1250.

Arévalo, P., Benavides, D., Tostado-Véliz, M., Aguado, J. A., & Jurado, F. (2023). Smart monitoring method for photovoltaic systems and failure control based on power smoothing techniques. *Renewable Energy*, 205, 366–383.

Bäuerle, N. & Jaśkiewicz, A. (2024). Markov decision processes with risk-sensitive criteria: An overview. *Mathematical Methods of Operations Research*, 99(1), 141–178.

Berner, C., Brockman, G., Chan, B., Cheung, V., Dębiak, P., Dennison, C., Farhi, D., Fischer, Q., Hashme, S., Hesse, C., et al. (2019). Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*.

Bisi, L., Sabbioni, L., Vittori, E., Papini, M., & Restelli, M. (2021). Risk-averse trust region optimization for reward-volatility reduction. In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence* (pp. 4583–4589).

Blancas-Rivera, R. & Jasso-Fuentes, H. (2024). Discrete-time hybrid control with risk-sensitive discounted costs. *Discrete Event Dynamic Systems*, 34(4), 659–687.

Bubeck, S. et al. (2015). Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4), 231–357.

Chapaloglou, S., Nesiadis, A., Iliadis, P., Atsonios, K., Nikolopoulos, N., Grammelis, P., Yiakopoulos, C., Antoniadis, I., & Kakaras, E. (2019). Smart energy management algorithm for load smoothing and peak shaving based on load forecasting of an island's power system. *Applied energy*, 238, 627–642.

Chen, W., Wang, J., Yu, G., Chen, J., & Hu, Y. (2022). Research on day-ahead transactions between multi-microgrid based on cooperative game model. *Applied Energy*, 316, 119106.

Chen, Y. & Cassandras, C. G. (2025). Scalable adaptive traffic light control over a traffic network including turns, transit delays, and blocking. *Discrete Event Dynamic Systems*, (pp. 1–30).

Cheng, M., Zhou, R., Kumar, P., & Tian, C. (2024). Provable policy gradient methods for average-reward Markov potential games. In *International Conference on Artificial Intelligence and Statistics* (pp. 4699–4707).

Etesami, S. R., Saad, W., Mandayam, N. B., & Poor, H. V. (2018). Stochastic games for the smart grid energy management with prospect prosumers. *IEEE Transactions on Automatic Control*, 63(8), 2327–2342.

Filar, J. A., Kallenberg, L. C. M., & Lee, H. (1989). Variance-penalized Markov decision processes. *Mathematics of Operations Research*, 14(1), 147–161.

Fujimoto, S., Hoof, H., & Meger, D. (2018). Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning* (pp. 1587–1596).: PMLR.

García, J. & Fernández, F. (2015). A comprehensive survey on safe reinforcement learning. *Journal of Machine Learning Research*, 16(1), 1437–1480.

Ghafoori, M., Abdallah, M., & Kim, S. (2023). Electricity peak shaving for commercial buildings using machine learning and vehicle to building (V2B) system. *Applied Energy*, 340, 121052.

Guo, Y., Zhang, Q., & Wang, Z. (2021). Cooperative peak shaving and voltage regulation in unbalanced distribution feeders. *IEEE Transactions on Power Systems*, 36(6), 5235–5244.

Haarnoja, T., Zhou, A., Abbeel, P., & Levine, S. (2018). Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning* (pp. 1861–1870).: PMLR.

Han, J., Fang, Y., Li, Y., Du, E., & Zhang, N. (2025). Optimal planning of multi-microgrid system with shared energy storage based on capacity leasing and energy sharing. *IEEE Transactions on Smart Grid*, 16, 16–31.

Hu, J., Xia, L., Hu, J., & Wu, H. (2025). Economical and reliable energy management for networked microgrids in a multi-agent collaborative manner. *IEEE Transactions on Automation Science and Engineering*, 22, 8659 – 8669.

Huang, Y. (2018). Finite horizon continuous-time Markov decision processes with mean and variance criteria. *Discrete Event Dynamic Systems*, 28, 539–564.

IESO (2023). Ontario and Market Demand. https://ieso.ca/Power-Data/Data-Directory.

Leonardos, S., Overman, W., Panageas, I., & Piliouras, G. (2022). Global convergence of multi-agent policy gradient in Markov potential games. In *International Conference on Learning Representations*.

Li, J., Herdem, M. S., Nathwani, J., & Wen, J. Z. (2023). Methods and applications for artificial intelligence, big data, internet of things, and blockchain in smart energy management. *Energy and AI*, 11, 100208.

Ma, S., Ma, X., & Xia, L. (2023). A unified algorithm framework for mean-variance optimization in discounted Markov decision processes. *European Journal of Operational Research*, 311(3), 1057–1067.

Ma, X., Tang, X., Xia, L., Yang, J., & Zhao, Q. (2021). Average-reward reinforcement learning with trust region methods. In *International Joint Conference on Artificial Intelligence* (pp. 2797–2083).

Manojkumar, R., Kumar, C., Ganguly, S., Gooi, H. B., Mekhilef, S., & Catalão, J. P. (2022). Rule-based peak shaving using master-slave level optimization in a diesel generator supplied microgrid. *IEEE Transactions on Power Systems*, 38(3), 2177–2188.

Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., & Kavukcuoglu, K. (2016). Asynchronous methods for deep reinforcement learning. In *International Conference on Machine Learning* (pp. 1928–1937).: PMLR.

Nawaz, A., Zhou, M., Wu, J., & Long, C. (2022). A comprehensive review on energy management, demand response, and coordination schemes utilization in multi-microgrids network. *Applied Energy*, 323, 119596.

NREL (2025). National wind technology center. (online). https://midcdmz.nrel.gov/apps/sitehome.pl?site=NWTC. (accessed date January 1, 2025).

Peirelinck, T., Hermans, C., Spiessens, F., & Deconinck, G. (2024). Combined peak reduction and self-consumption using proximal policy optimisation. *Energy and AI*, 16, 100323.

Prashanth, L. & Fu, M. C. (2022). Risk-sensitive reinforcement learning via policy gradient search. *Foundations and Trends® in Machine Learning*, 15(5), 537–693.

Qiu, W., Wang, X., Yu, R., Wang, R., He, X., An, B., Obraztsova, S., & Rabinovich, Z. (2021). RMIX: Learning risk-sensitive policies for cooperative reinforcement learning agents. *Advances in Neural Information Processing Systems*, 34, 23049–23062.

Rosa, M., Cury, J. E., & Baldissera, F. L. (2024). A modular synthesis approach for the coordination of multi-agent systems: the multiple team case. *Discrete Event Dynamic Systems*, 34(1), 163–198.

Rostmnezhad, Z. & Dessaint, L. (2023). Power management in smart buildings using reinforcement learning. In *2023 IEEE Power & Energy Society Innovative Smart Grid Technologies Conference (ISGT)* (pp. 1–5).

Schulman, J., Moritz, P., Levine, S., Jordan, M., & Abbeel, P. (2016). High-dimensional continuous control using generalized advantage estimation. In *International Conference on Learning Representations*.

Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.

Shen, S., Ma, C., Li, C., Liu, W., Fu, Y., Mei, S., Liu, X., & Wang, C. (2023). RiskQ: Risk-sensitive multi-agent reinforcement learning value factorization. *Advances in Neural Information Processing Systems*, 36, 34791–34825.

Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al. (2017). Mastering the game of Go without human knowledge. *Nature*, 550(7676), 354–359.

Sutton, R. S. & Barto, A. G. (2018). *Reinforcement learning: An introduction.* MIT press.

Sutton, R. S., McAllester, D., Singh, S., & Mansour, Y. (1999). Policy gradient methods for reinforcement learning with function approximation. *Advances in Neural Information Processing Systems*, 12, 1057–1063.

Tang, H., Liu, C., Cao, Y., Lv, K., & Zhang, Q. (2021). Hierarchical scheduling learning optimisation of two-area active distribution system considering peak shaving demand of power grid. *Discrete Event Dynamic Systems*, 31(3), 439–468.

Wang, X., Zhou, J., Qin, B., & Guo, L. (2023). Coordinated power smoothing control strategy of multi-wind turbines and energy storage systems in wind farm based on MADRL. *IEEE Transactions on Sustainable Energy*, 15(1), 368–380.

Weitzel, T. & Glock, C. H. (2018). Energy management for stationary electric energy storage systems: A systematic literature review. *European Journal of Operational Research*, 264(2), 582–606.

Xia, L. (2018). Variance minimization of parameterized Markov decision processes. *Discrete Event Dynamic Systems*, 28, 63–81.

Xia, L. (2020). Risk-sensitive Markov decision processes with combined metrics of mean and variance. *Production and Operations Management*, 29(12), 2808–2827.

Yang, Z., Xia, L., & Guan, X. (2020). Fluctuation reduction of wind power and sizing of battery energy storage systems in microgrids. *IEEE Transactions on Automation Science and Engineering*, 17(3), 1195–1207.

Yu, C., Velu, A., Vinitsky, E., Gao, J., Wang, Y., Bayen, A., & Wu, Y. (2022). The surprising effectiveness of ppo in cooperative multi-agent games. *Advances in Neural Information Processing Systems*, 35, 24611–24624.

Zhang, R., Ren, Z., & Li, N. (2024). Gradient play in stochastic games: Stationary points, convergence, and sample complexity. *IEEE Transactions on Automatic Control.*

Zhong, Y., Kuba, J. G., Feng, X., Hu, S., Ji, J., & Yang, Y. (2024). Heterogeneous-agent reinforcement learning. *Journal of Machine Learning Research*, 25(32), 1–67.

# Appendix

## A  Proofs

### A.1  Proof of Lemma 2

*Proof.* Consider two joint policies, $\boldsymbol{\mu}'$ and $\boldsymbol{\mu}$, parameterized by $\theta'$ and $\theta$, respectively. According to Lemma 1, we have

$$J(\boldsymbol{\mu}') - J(\boldsymbol{\mu}) = \mathbb{E}_{s\sim\pi^{\boldsymbol{\mu}'},\boldsymbol{a}\sim\boldsymbol{\mu}'}[A_f^{\boldsymbol{\mu}}(s,\boldsymbol{a})] + \beta(\eta^{\boldsymbol{\mu}'} - \eta^{\boldsymbol{\mu}})^2.$$

We denote $h_1(\boldsymbol{\mu}',\boldsymbol{\mu}) = \mathbb{E}_{s\sim\pi^{\boldsymbol{\mu}'},\boldsymbol{a}\sim\boldsymbol{\mu}'}[A_f^{\boldsymbol{\mu}}(s,\boldsymbol{a})]$ and $h_2(\boldsymbol{\mu}',\boldsymbol{\mu}) = (\eta^{\boldsymbol{\mu}'} - \eta^{\boldsymbol{\mu}})^2$. Considering the policy parameterization, we denote $\Delta\theta = \theta' - \theta$ and let $\theta' \to \theta$. For $h_1(\theta',\theta)$ we have

$$
\begin{aligned}
\nabla_\theta h_1 &= \lim_{\Delta\theta\to 0} \frac{1}{\Delta\theta} \sum_s \pi^{\theta'}(s) \sum_{\boldsymbol{a}} \left[ \boldsymbol{\mu}^{\theta'}(\boldsymbol{a}|s) A_f^\theta(s,\boldsymbol{a}) \right] \\
&\overset{(i)}{=} \lim_{\Delta\theta\to 0} \sum_s \pi^{\theta'}(s) \sum_{\boldsymbol{a}} \frac{\boldsymbol{\mu}^{\theta'}(\boldsymbol{a}|s) - \boldsymbol{\mu}^\theta(\boldsymbol{a}|s)}{\Delta\theta} A_f^\theta(s,\boldsymbol{a}) \\
&= \sum_s \pi^\theta(s) \sum_{\boldsymbol{a}} \nabla_\theta \boldsymbol{\mu}^\theta(\boldsymbol{a}|s) A_f^\theta(s,\boldsymbol{a}) \\
&\overset{(ii)}{=} \mathbb{E}_{s\sim\pi,\boldsymbol{a}\sim\boldsymbol{\mu}^\theta} \left[ \nabla_\theta \log \boldsymbol{\mu}^\theta(\boldsymbol{a}|s) A_f^\theta(s,\boldsymbol{a}) \right] \\
&= \mathbb{E}_{s\sim\pi,\boldsymbol{a}\sim\boldsymbol{\mu}^\theta} \left[ \nabla_\theta \log \boldsymbol{\mu}^\theta(\boldsymbol{a}|s) Q_f^\theta(s,\boldsymbol{a}) \right],
\end{aligned}
$$

where the Equality (i) holds is due to $\mathbb{E}_{\boldsymbol{a}\sim\boldsymbol{\mu}}[A_f^{\boldsymbol{\mu}}(s,\boldsymbol{a})] = 0$, Equality (ii) holds is due to $\nabla_\theta \log \boldsymbol{\mu}^\theta(\boldsymbol{a}|s) = \frac{\nabla_\theta \boldsymbol{\mu}^\theta(\boldsymbol{a}|s)}{\boldsymbol{\mu}^\theta(\boldsymbol{a}|s)}$. For $h_2(\theta',\theta)$ we have

$$
\begin{aligned}
\nabla_\theta h_2 &= \lim_{\Delta\theta\to 0} \frac{(\eta^{\theta'} - \eta^\theta)^2}{\Delta\theta} \\
&= \lim_{\Delta\theta\to 0} 2(\eta^{\theta'} - \eta^\theta)\nabla_\theta \eta^\theta \\
&= 0.
\end{aligned}
$$

Then, we arrive at

$$\nabla_\theta J(\theta) = \nabla_\theta h_1 - \beta \nabla_\theta h_2$$
$$= \mathbb{E}_{s \sim \pi^\theta, \boldsymbol{a} \sim \boldsymbol{\mu}^\theta} \left[ \nabla_\theta \log \boldsymbol{\mu}^\theta(\boldsymbol{a}|s) Q_f^\theta(s, \boldsymbol{a}) \right],$$

then the proof is finished. $\qquad \square$

## A.2   Proof of Lemma 3

*Proof.* According to the policy gradient Lemma 2, we have

$$\frac{\partial J(\theta)}{\partial \theta_{s,a_i}} = \sum_{\boldsymbol{a}} \pi^\theta(s) \boldsymbol{\mu}^\theta(\boldsymbol{a}|s) \frac{\partial \log \boldsymbol{\mu}^\theta(\boldsymbol{a}|s)}{\partial \theta_{s,a_i}} Q_f^\theta(s, \boldsymbol{a}).$$

For the direct policy parameterization, with $\mu^\theta(\boldsymbol{a}|s) = \prod_i^N \mu^{\theta_i}(a_i|s)$, we have

$$\frac{\partial \log \boldsymbol{\mu}^\theta(\boldsymbol{a}|s)}{\partial \theta_{s,a_i}} = \frac{\partial \log \mu^{\theta_i}(a_i|s)}{\partial \theta_{s,a_i}} = \frac{1}{\mu^{\theta_i}(a_i|s)}.$$

Then we arrive at

$$\frac{\partial J(\theta)}{\partial \theta_{s,a_i}} = \sum_{\boldsymbol{a}_{-i}} \pi^\theta(s) \mu^{\theta_{-i}}(\boldsymbol{a}_{-i}|s) \mu^{\theta_i}(a_i|s) \frac{1}{\mu^{\theta_i}(a_i|s)} Q_f^\theta(s, \boldsymbol{a})$$
$$= \sum_{\boldsymbol{a}_{-i}} \pi^\theta(s) \mu^{\theta_{-i}}(\boldsymbol{a}_{-i}|s) Q_f^\theta(s, \boldsymbol{a})$$
$$= \pi^\theta(s) \overline{Q_{f,i}^\theta}(s, a_i),$$

where $\overline{Q_{f,i}^\theta}(s, a_i) := \sum_{\boldsymbol{a}_{-i}} \boldsymbol{\mu}^{\theta_{-i}}(\boldsymbol{a}_{-i}|s) Q_f^\theta(s, a_i, \boldsymbol{a}_{-i})$. The proof is finished. $\qquad \square$

## A.3 Proof of Lemma 5

*Proof.* First, we provide the following definitions.

$$\theta_{i,\epsilon}(a_i|s) = \theta_{s,a_i} + \epsilon u_{s,a_i},$$

$$\theta_{j,\tau}(a_j|s) = \theta_{s,a_j} + \tau u_{s,a_j},$$

$$\theta_\epsilon(\boldsymbol{a}|s) = \theta_{i,\epsilon}(a_i|s)\theta_{-i}(a_{-i}|s),$$

$$\theta_\tau(\boldsymbol{a}|s) = \theta_{j,\tau}(a_j|s)\theta_{-j}(a_{-j}|s),$$

$$J_\epsilon = J(\theta_{i,\epsilon}, \theta_{-i}),$$

$$J_{\epsilon,\tau} = J(\theta_{i,\epsilon}, \theta_{j,\tau}, \theta_{-ij}),$$

$$\overline{r_i^\theta}(s, a_i) := \sum_{\boldsymbol{a}_{-i}} \boldsymbol{\mu}^{\theta_{-i}}(\boldsymbol{a}_{-i}|s)r(s, a_i, \boldsymbol{a}_{-i}),$$

$$\overline{f_i^\theta}(s, a_i) := \sum_{\boldsymbol{a}_{-i}} \boldsymbol{\mu}^{\theta_{-i}}(\boldsymbol{a}_{-i}|s)f^\theta(s, a_i, \boldsymbol{a}_{-i}).$$

According to these definitions above, we have

$$J_\epsilon = \sum_{s,a} \pi^{\theta_\epsilon}(s)\theta_\epsilon(\boldsymbol{a}|s)f^{\theta_\epsilon}(s, \boldsymbol{a}) = \sum_{s,a_i} \pi^{\theta_\epsilon}(s)\theta_{i,\epsilon}(a_i|s)\overline{f_i^{\theta_\epsilon}}(s, a_i),$$

where $\overline{f_i^{\theta_\epsilon}}(s, a_i) = \sum_{\boldsymbol{a}_{-i}} \theta_{-i}(\boldsymbol{a}_{-i}|s)f^{\theta_\epsilon}(s, a_i, \boldsymbol{a}_{-i})$. Next, we compute the first-order and second-order derivatives of $\overline{f_i^{\theta_\epsilon}}(s, a_i)$, respectively. For the first-order derivative, we have

$$\frac{\mathrm{d}\overline{f_i^{\theta_\epsilon}}(s, a_i)}{\mathrm{d}\epsilon} = 2\beta(\overline{r_i}(s, a_i) - \eta^{\theta_\epsilon})\frac{\mathrm{d}\eta^{\theta_\epsilon}}{\mathrm{d}\epsilon},$$

$$\frac{\mathrm{d}^2\overline{f_i^{\theta_\epsilon}}(s, a_i)}{\mathrm{d}\epsilon^2} = 2\beta\left[\left(\overline{r_i}(s, a_i) - \eta^{\theta_\epsilon}\right)\frac{\mathrm{d}^2\eta^{\theta_\epsilon}}{\mathrm{d}\epsilon^2} - \left(\frac{\mathrm{d}\eta^{\theta_\epsilon}}{\mathrm{d}\epsilon}\right)^2\right].$$

For TSGs, it holds that $\nabla_{\boldsymbol{\mu}}J^{\boldsymbol{\mu}} = (\frac{\partial J^{\boldsymbol{\mu}}}{\partial\mu_1}, \dots, \frac{\partial J^{\boldsymbol{\mu}}}{\partial\mu_N})^T$. To investigate the smooth property of

the mean-variance performance function, we have

$$\left\| \nabla_{\boldsymbol{\mu}} J^{\boldsymbol{\mu}} - \nabla_{\boldsymbol{\mu}} J^{\boldsymbol{\mu}'} \right\|_2^2 = \sum_{i=1}^N \left\| \nabla_{\mu_i} J^{\boldsymbol{\mu}} - \nabla_{\mu_i} J^{\boldsymbol{\mu}'} \right\|_2^2$$

$$\leq \sum_{i=1}^N \sum_{j=1}^N \left\| \nabla_{\mu_i} J^{(\mu_{1\sim j-1}, \mu'_{j\sim N})} - \nabla_{\mu_i} J^{(\mu_{1\sim j}, \mu'_{j+1\sim N})} \right\|_2^2,$$

where $J^{\boldsymbol{\mu}'} = J^{(\mu_{1\sim 0}, \mu'_{1\sim N})}$, $J^{\boldsymbol{\mu}} = J^{(\mu_{1\sim N}, \mu'_{N+1\sim N})}$.

Cheng et al. (2024) show that $\left\| \frac{\mathrm{d}\pi^{\theta_\epsilon}}{\mathrm{d}\epsilon} \right\|_2 \leq \frac{\kappa_0 \sqrt{SA_{\max}}}{2}$, $\left\| \frac{\mathrm{d}^2 \pi^{\theta_\epsilon}}{\mathrm{d}\epsilon^2} \right\|_2 \leq \kappa_0^2 SA_{\max}$, and for the average performance function $\eta$, it holds that

$$\frac{\mathrm{d}\eta^{\theta_\epsilon}}{\mathrm{d}\epsilon} = \sum_{s,a_i} \pi^{\mu_\epsilon}(s) u_{s,a_i} \overline{r_i}(s, a_i) + \sum_{s,a_i} \frac{\mathrm{d}\pi^{\mu_\epsilon}(s)}{\mathrm{d}\epsilon} \theta_{i,\epsilon}(a_i | s) \overline{r_i}(s, a_i),$$

$$\left| \frac{\mathrm{d}\eta^{\theta_\epsilon}}{\mathrm{d}\epsilon} \right| \leq \left\| \pi^{\theta_\epsilon} \right\|_2 \sqrt{\sum_s (\sum_{a_i} u_{s,a_i} \overline{r_i}(s, a_i))^2} + \left\| \frac{\mathrm{d}\pi^{\theta_\epsilon}}{\mathrm{d}\epsilon} \right\|_1 \leq \sqrt{A_{\max}}(1 + \frac{\kappa_0 S}{2}),$$

Next, we first derive the first-order derivative of the mean-variance performance function with respect to $\epsilon$,

$$\frac{\mathrm{d}J_\epsilon}{\mathrm{d}\epsilon} = \underbrace{\sum_{s,a_i} \pi^{\theta_\epsilon}(s) u_{s,a_i} \overline{f_i^{\theta_\epsilon}}(s, a_i)}_{\text{Part A}}$$

$$+ \underbrace{\sum_{s,a_i} \frac{\mathrm{d}\pi^{\theta_\epsilon}(s)}{\mathrm{d}\epsilon} \theta_{i,\epsilon}(a_i | s) \overline{f_i^{\theta_\epsilon}}(s, a_i)}_{\text{Part B}}$$

$$+ \underbrace{\sum_{s,a_i} \pi^{\theta_\epsilon}(s) \theta_{i,\epsilon}(a_i | s) \frac{\mathrm{d}\overline{f_i^{\theta_\epsilon}}(s, a_i)}{\mathrm{d}\epsilon}}_{\text{Part C}}.$$

To compute the second-order derivative, we take the derivative of each of the three parts

41

in the above expression separately. The derivative of Part A is as follows:

$$\frac{\mathrm{d}\mathrm{Part\ A}}{\mathrm{d}\epsilon} = \sum_{s,a_i} \frac{\mathrm{d}\pi^{\theta_\epsilon}(s)}{\mathrm{d}\epsilon} u_{s,a_i} \overline{f_i^{\theta_\epsilon}}(s,a_i) + \sum_{s,a_i} \pi^{\theta_\epsilon}(s) u_{s,a_i} \frac{\mathrm{d}\overline{f_i^{\theta_\epsilon}}(s,a_i)}{\mathrm{d}\epsilon}.$$

The derivative of Part B is as follows:

$$\begin{aligned}
\frac{\mathrm{d}\mathrm{Part\ B}}{\mathrm{d}\epsilon} &= \sum_{s,a_i} \frac{\mathrm{d}\pi^{\theta_\epsilon}(s)}{\mathrm{d}\epsilon} u_{s,a_i} \overline{f_i^{\theta_\epsilon}}(s,a_i) \\
&+ \sum_{s,a_i} \frac{\mathrm{d}\pi^{\theta_\epsilon}(s)}{\mathrm{d}\epsilon} \theta_{i,\epsilon}(s,a_i) \frac{\mathrm{d}\overline{f_i^{\theta_\epsilon}}(s,a_i)}{\mathrm{d}\epsilon} \\
&+ \sum_{s,a_i} \frac{\mathrm{d}^2\pi^{\theta_\epsilon}(s)}{\mathrm{d}\epsilon^2} \theta_{i,\epsilon}(s,a_i) f_i^{\theta_\epsilon}.
\end{aligned}$$

The derivative of Part C is as follows:

$$\begin{aligned}
\frac{\mathrm{d}\mathrm{Part\ C}}{\mathrm{d}\epsilon} &= \sum_{s,a_i} \frac{\mathrm{d}\pi^{\theta_\epsilon}(s)}{\mathrm{d}\epsilon} \theta_{i,\epsilon}(s,a_i) \frac{\mathrm{d}\overline{f_i^{\theta_\epsilon}}}{\mathrm{d}\epsilon} \\
&+ \sum_{s,a_i} \pi^{\theta_\epsilon}(s) u_{s,a_i} \frac{\mathrm{d}\overline{f_i^{\theta_\epsilon}}(s,a_i)}{\mathrm{d}\epsilon} \\
&+ \sum_{s,a_i} \pi^{\theta_\epsilon}(s) \theta_{i,\epsilon}(s,a_i) \frac{\mathrm{d}^2\overline{f_i^{\theta_\epsilon}}(s,a_i)}{\mathrm{d}\epsilon^2}.
\end{aligned}$$

By arranging the results, we obtain that the second-order derivative of the mean-variance

performance function with respect to $\epsilon$ is:

$$
\begin{aligned}
\frac{\mathrm{d}^2 J}{\mathrm{d}\epsilon^2} =& \frac{\mathrm{d}\text{Part A}}{\mathrm{d}\epsilon} + \frac{\mathrm{d}\text{Part B}}{\mathrm{d}\epsilon} + \frac{\mathrm{d}\text{Part C}}{\mathrm{d}\epsilon} \\
=& \underbrace{2 \sum_{s,a_i} \frac{\mathrm{d}\pi^{\theta_\epsilon}(s)}{\mathrm{d}\epsilon} u_{s,a_i} \overline{f_i^{\theta_\epsilon}}(s, a_i)}_{①} \\
& + \underbrace{2 \sum_{s,a_i} \pi^{\theta_\epsilon}(s) u_{s,a_i} \frac{\mathrm{d}\overline{f_i^{\theta_\epsilon}}(s, a_i)}{\mathrm{d}\epsilon}}_{②} \\
& + \underbrace{2 \sum_{s,a_i} \frac{\mathrm{d}\pi^{\theta_\epsilon}(s)}{\mathrm{d}\epsilon} \theta_{i,\epsilon}(s, a_i) \frac{\mathrm{d}\overline{f_i^{\theta_\epsilon}}(s, a_i)}{\mathrm{d}\epsilon}}_{③} \\
& + \underbrace{\sum_{s,a_i} \frac{\mathrm{d}^2 \pi^{\theta_\epsilon}(s)}{\mathrm{d}\epsilon^2} \theta_{i,\epsilon}(s, a_i) \overline{f_i^{\theta_\epsilon}}(s, a_i)}_{④} \\
& + \underbrace{\sum_{s,a_i} \pi^{\theta_\epsilon}(s) \theta_{i,\epsilon}(s, a_i) \frac{\mathrm{d}^2 \overline{f_i^{\theta_\epsilon}}(s, a_i)}{\mathrm{d}\epsilon^2}}_{⑤} .
\end{aligned}
$$

Based on the above result, we derive an upper bound for each of the five terms separately. For the first term, we have:

$$
\begin{aligned}
|①| &\leq 2 \left\| \frac{\mathrm{d}\pi^{\theta_\epsilon}}{\mathrm{d}\epsilon} \right\|_2 \sqrt{\sum_s (\sum_{a_i} u_{s,a_i} \overline{f_i^{\theta_\epsilon}}(s, a_i))^2} \\
&\leq 2 \left\| \frac{\mathrm{d}\pi^{\theta_\epsilon}}{\mathrm{d}\epsilon} \right\|_2 \sqrt{A_i} f_{\max} \|u_i\|_2 \\
&\leq \kappa_0 f_{\max} A_{\max} \sqrt{S}.
\end{aligned}
$$

where $f_{\max}$ represents the upper bound of the magnitude of the surrogate reward function

$f$. For the second term, we have:

$$|②| = \left|4\beta \sum_{s,a_i} \pi^{\theta_\epsilon} u_{s,a_i}(\overline{r_i}(s,a_i) - \eta^{\theta_\epsilon})\frac{\mathrm{d}\eta^{\theta_\epsilon}}{\mathrm{d}\epsilon}\right|$$

$$\leq 4\beta \left\|\pi^{\theta_\epsilon}\right\|_2 \sqrt{\sum_s [\sum_{a_i} u_{s,a_i}(\overline{r_i}(s,a_i) - \eta^{\theta_\epsilon})]^2} \left|\frac{\mathrm{d}\eta^{\theta_\epsilon}}{\mathrm{d}\epsilon}\right|$$

$$\leq 4\beta \cdot \sqrt{A_i} \cdot (\sqrt{A_i} + \frac{\kappa_0}{2}S\sqrt{A_i})$$

$$\leq 4\beta A_{\max}(1 + \frac{\kappa_0}{2}S).$$

For the third term, we have:

$$|③| = \left|4\beta \sum_{s,a_i} \frac{\mathrm{d}\pi^{\theta_\epsilon}}{\mathrm{d}\epsilon}\theta_{i,\epsilon}(a_i|s)(\overline{r_i}(s,a_i) - \eta^{\theta_\epsilon})\frac{\mathrm{d}\eta^{\theta_\epsilon}}{\mathrm{d}\epsilon}\right|$$

$$\leq 4\beta \left\|\frac{\mathrm{d}\pi^{\theta_\epsilon}}{\mathrm{d}\epsilon}\right\|_1 \cdot \left|\frac{\mathrm{d}\eta^{\theta_\epsilon}}{\mathrm{d}\epsilon}\right|$$

$$\leq 4\beta\sqrt{S} \left\|\frac{\mathrm{d}\pi^{\theta_\epsilon}}{\mathrm{d}\epsilon}\right\|_2 \sqrt{A}(1 + \frac{\kappa_0 S}{2})$$

$$\leq \beta\kappa_0 S A_{\max}(2 + \kappa_0 S).$$

For the fourth term, we have:

$$|④| \leq \left\|\frac{\mathrm{d}^2\pi^{\theta_\epsilon}}{\mathrm{d}\epsilon^2}\right\|_1 \cdot f_{\max}$$

$$\leq \kappa_0^2 S^{\frac{3}{2}} A_{\max} f_{\max}.$$

For the fifth term, we have:

$$|⑤| = \left|\sum_{s,a_i} \pi^{\theta_\epsilon}\theta_{i,\epsilon}(a_i|s) \cdot 2\beta \left[(\overline{r_i}(s,a_i) - \eta^{\theta_\epsilon})\frac{\mathrm{d}^2\eta^{\theta_\epsilon}}{\mathrm{d}\epsilon^2} - (\frac{\mathrm{d}\eta^{\theta_\epsilon}}{\mathrm{d}\epsilon})^2\right]\right|$$

$$= \left|2\beta \left[(\eta^{\theta_\epsilon} - \eta^{\theta_\epsilon})\frac{\mathrm{d}^2\eta^{\theta_\epsilon}}{\mathrm{d}\epsilon^2} - \sum_{s,a_i} \pi^{\theta_\epsilon}\theta_{i,\epsilon}(a_i|s)(\frac{\mathrm{d}\eta^{\theta_\epsilon}}{\mathrm{d}\epsilon})^2\right]\right|$$

$$\leq 2\beta \left|\frac{\mathrm{d}\eta^{\theta_\epsilon}}{\mathrm{d}\epsilon}\right|^2$$

$$\leq 2\beta A_{\max} R_{\max}^2 (1 + \frac{\kappa_0 S}{2})^2,$$

where $R_{\max}$ represents the upper bound of the magnitude of the reward function $r$.

In this proof, we assume that the reward function is normalized, i.e., $r \in [0, 1]$, then we have $R_{\max} = 1$ and $f_{\max} = 1 + \beta$. Furthermore, we have

$$\left| \frac{\mathrm{d}^2 J}{\mathrm{d}\epsilon^2} \right| \leq |\text{①}| + |\text{②}| + |\text{③}| + |\text{④}| + |\text{⑤}|$$

$$\leq 6\beta A_{\max}(1 + \frac{\kappa_0}{2}S)^2 + (1 + \beta)\kappa_0 A_{\max}\sqrt{S}(\kappa_0 S + 1) = L.$$

In a similar manner, we consider the successive second-order differentiation of the performance function $J_{\epsilon,\tau}$ with respect to $\epsilon$ and $\tau$. As the derivation is similar to the above, it is omitted here. We obtain:

$$\frac{\partial^2 J_{\epsilon,\tau}}{\partial \epsilon \partial \tau} \leq 6\beta A_{\max}(1 + \frac{\kappa_0}{2}S)^2 + (1 + \beta)A_{\max}(\kappa_0 S + \kappa_0^2 S^{\frac{3}{2}} + \kappa_0\sqrt{S} + 1) = \frac{L_J}{N}.$$

Then, the proof is finished. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ □

## A.4   Proof of Theorem 1

*Proof.* Prior to the proof of Theorem 1, we introduce the auxiliary Lemma 6 and Lemma 7.

**Lemma 6** (Bubeck et al. (2015), Lemma 3.6)**.** *Assume that $J(\theta)$ is $L_J$-smooth with respect to $\theta \in \mathcal{X}$. Define the gradient mapping as follows:*

$$G^\alpha(\theta) := \frac{1}{\alpha}\left(Proj_{\mathcal{X}}(\theta + \alpha\nabla J(\theta)) - \theta\right).$$

*Let $\theta^+ = \theta + \alpha G^\alpha(\theta) = Proj_{\mathcal{X}}(\theta + \alpha\nabla J(\theta))$, then for $\alpha \leq \frac{1}{L_J}$ we have*

$$J(\theta^+) - J(\theta) \geq \frac{\alpha}{2}\|G^\alpha(\theta)\|_2^2.$$

**Lemma 7** (Agarwal et al. (2021), Proposition 37)**.** *Assume that $J(\theta)$ is $L_J$-smooth over*

$\theta \in \mathcal{X}$. The projected gradient update is defined as $\theta^+ = \theta + \alpha G^\alpha(\theta)$. If $|G^\alpha(\theta)|_2 \leq \epsilon$, then:

$$\max_{\theta \in \mathcal{X}} (\theta - \theta^+)^\top \nabla_\theta J(\theta^+) \leq \epsilon(1 + \alpha L_J).$$

In particular, when $\alpha \leq \frac{1}{L_J}$, we have

$$\max_{\theta \in \mathcal{X}} (\theta - \theta^+)^\top \nabla_\theta J(\theta^+) \leq 2\epsilon.$$

Next, we prove Theorem 1. According to Lemma 6, when $\alpha \leq \frac{1}{L_J}$, we have

$$
\begin{aligned}
\frac{1}{K} \sum_{k=0}^{K-1} \|G^\alpha(\theta)\|_2^2 &\leq \frac{2L_J}{K} \sum_{k=0}^{K-1} \left( J(\theta^{(k+1)}) - J(\theta^{(k)}) \right) \\
&= \frac{2L_J}{K} (J(\theta^{(K)})) - J(\theta^{(0)})) \\
&\leq \frac{2L_J}{K} (J_{\max} - J_{\min}).
\end{aligned}
$$

Let $\frac{2L_J}{K}(J_{\max} - J_{\min}) = \frac{\epsilon}{2}$, which implies that $K = \frac{4L_J(J_{\max} - J_{\min})}{\epsilon}$, then we have

$$\min_k \|G^\alpha(\theta^{(k)})\| \leq \frac{\epsilon}{2}.$$

Let $k^* = \arg\min_k |G(\theta^{(k)})|$. By Lemma 7, we can conclude that $\theta^{(k+1)}$ is $\epsilon$-stationary. This completes the proof.

$\square$