





Towards cost-optimal joint electricity-computation management: A novel predict-then-optimize framework

Yibo Ding^{a,b} , Xudong Li^a , Yuhong Zhao^{a,d}, Wenzhuo Shi^{a,f}, Cheng Lyu^e , Jiaqi Ruan^g, Zhao Xu^{a,b,c,*} 

^a Department of Electrical and Electronics Engineering, The Hong Kong Polytechnic University, Hong Kong Special Administrative Region of China

^b Research Institute for Smart Energy, The Hong Kong Polytechnic University, Hong Kong Special Administrative Region of China

^c Shenzhen Research Institute, The Hong Kong Polytechnic University, Shenzhen, 518129, China

^d School of Electrical Engineering, Xi'an Jiaotong University, Xi'an, 710049, China

^e Wu Jieh Yee School of Interdisciplinary Studies, Lingnan University, Hong Kong Special Administrative Region of China

^f School of Automation, Northwestern Polytechnical University, Xi'an, 710129, China

^g College of Electrical Engineering, Sichuan University, Chengdu, 610044, China

HIGHLIGHTS

- A cost-oriented predict-then-optimize decision-making framework is developed for the joint electricity-computation dispatch problem.
- A novel privacy-preserving iterative algorithm with guaranteed faster convergence performance is proposed.
- Heterogeneous uncertainties arising from power outputs of RES and workloads of DC are comprehensively considered.
- Extreme learning machine is employed as prediction model.

ARTICLE INFO

Keywords:

Energy management
Data center
Joint dispatch
Iterative algorithm

ABSTRACT

The escalating computing demand due to the flourishing of artificial intelligence is catalyzing more comprehensive and intricate interactions between modern power systems and data centers (DCs), necessitating joint electricity-computation management towards cost-optimal operation. The power system operator (SO) dispatches the generators, and the DC operator (DCO) optimizes the server dispatch strategies, where coupled information interactions exist. In practical, SO and DCO would encounter uncertainties arising from power outputs of renewable energy sources (RES) and computing workload requests submitted by end-users, respectively. Conventional accuracy-oriented predict-then-optimize (PTO) framework may lead to sub-optimal solutions due to the asymmetric relationship between prediction error and decision error. To achieve cost-optimal dispatch strategies, developing a cost-oriented PTO decision-making framework for the joint management is essential. Specially, the prediction models are trained by minimizing the decision regret. In addition, a privacy-preserving dual-boundary feedback-embedded adaptive iterative algorithm is specially proposed to solve the joint dispatch problem, realizing guaranteed and faster convergence. Simulation results on a modified IEEE-30 bus system over extensive scenarios demonstrate that the cost-oriented PTO framework saves about 1.4% of the total operational cost compared to conventional accuracy-oriented decision framework on average. Moreover, the proposed iterative algorithm averagely reduces 20% of iteration times than the existing binary search method.

* Corresponding author at: Department of Electrical and Electronics Engineering, The Hong Kong Polytechnic University, Hong Kong Special Administrative Region of China.

Email addresses: yibo0712.ding@connect.polyu.hk (Y. Ding), xudong.li@connect.polyu.hk (X. Li), yuhong.zhao@connect.polyu.hk (Y. Zhao), wenzhuo.shi@connect.polyu.hk (W. Shi), cheng.lyu@connect.polyu.hk (C. Lyu), jiaqiruan@scu.edu.cn (J. Ruan), eezhaoxu@polyu.edu.hk (Z. Xu).

<https://doi.org/10.1016/j.apenergy.2026.127734>

Received 15 December 2025; Received in revised form 10 February 2026; Accepted 14 March 2026

Available online 21 March 2026

0306-2619/© 2026 The Authors. Published by Elsevier Ltd. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

1. Introduction

1.1. Background

In recent years, the unprecedented development and ubiquitous deployment of artificial intelligence (AI) technologies have led to a tremendous increase in end-user computing demands. As a result, data centers (DCs), which integrate an extensive number of servers, have witnessed widespread deployment. This trend has significantly escalated their electricity consumption. According to a report by the International Energy Agency [1], the total electricity consumption of DCs has reached 415 TWh in 2024, accounting for approximately 1.5% of the total global electricity usage. Moreover, the growth rate is accelerating, as the electricity consumption of DCs has increased by 10% over the past decade, compared to only 3% during 2005–2015. This suggests that DCs will engage in more extensive and comprehensive interactions with power systems, making the cost-optimal joint management a critical concern for both DC operators (DCO) and power system operators (SO). SO optimizes power outputs and reserve capacities of generators to guarantee power supply-demand balance. DCO optimizes the processing strategies of computing workloads and schedules the servers.

In this work, we specifically focus on power grid-connected internet data centers (IDCs), which are hyperscale facilities integrating massive amounts of servers, storage, and network devices to provide on-demand computing services [2]. IDCs are characterized by high energy density and significant electricity consumption, which is primarily composed of IT equipment load and cooling system load. The energy efficiency of an IDC is typically measured by power usage effectiveness (PUE), defined as the ratio of total facility energy to IT equipment energy.

According to the real-world settings of large-scale cloud service providers such as Google, computing tasks submitted by end-users are generally categorized into two types: interactive workload, which requires immediate response (e.g., web search queries), and batch workload, which is delay-tolerant (e.g., data storage and machine learning) [3]. Unlike electricity loads that require real-time supply-demand balance, batch workload can be flexibly scheduled within deadlines defined by end-users. This allows the DCO to delay execution during peak-price hours and process them when electricity prices drop. Furthermore, due to the geo-distributed nature of DC infrastructures [4], a DCO typically operates multiple DCs located at different buses in the power system, each associated with a locational marginal price (LMP). Assuming that there is no long-term electricity wholesale contract between the DCO and the SO, and batch workload can be parallelized between the DCs using techniques such as MapReduce [5], the DCO can flexibly dispatch servers to process these workloads in both spatial and temporal dimensions to minimize its power consumption cost [6].

1.2. Motivation, related work and research gaps

Due to the increasing penetration rate of renewable energy sources (RESs) under the progressive low-carbon transition [7,8], SO must cope with uncertainties from RESs generation. Despite the reserve capacities of thermal generators, flexible loads are becoming critical assets for balancing renewable fluctuations. Since DCs with spatial and temporal flexibilities are flourishing, the integration of DCs into grid operation is essential. Also, DCO would spontaneously utilize servers when the LMP is lower, unlocking the flexible role of DCs [9–12].

The LMP at day-ahead stage is typically determined by the SO through solving an economic dispatch problem [13], which requires the accessibility of the total power consumption of DCs for SO. Meanwhile, the server dispatch strategies of DCO also depend on the LMP declared by the SO. This creates a coupled information exchange between the DCO and the SO, highlighting the necessity of formulating joint electricity-computation dispatch problem for these two stakeholders [14,15]. Such a joint dispatch problem faces two main challenges: how to handle the uncertainties of RESs and workloads and how to develop

an effective iterative algorithm. Literature review regarding these two critical aspects is demonstrated as follows.

As aforementioned, the SO needs to cope with uncertainties in RESs. Meanwhile, the DCO faces uncertainties in both batch and interactive workload arrivals [16]. These uncertainties necessitate the training of prediction models that serve for better dispatch strategies. However, most existing studies about the joint dispatch merely delve into the optimization stage [5,13,17]. Therefore, it is essential to develop a predict-then-optimize (PTO) decision framework tailored for the joint dispatch problem against uncertainties.

Most existing prediction approaches are evaluated solely by norm-based accuracy metrics, such as mean absolute error (MAE). Regrettably, for some real-world decision-making problems, more accurate predictions do not necessarily contribute to better decisions [18]. This is due to the asymmetric and even non-monotonic relationships between prediction error and decision error [19,20]. Asymmetry implies that over-prediction (OP) and under-prediction (UP) incur unequal decision error [21], while non-monotonicity indicates that reducing prediction error does not necessarily translate into lower decision error. Traditional accuracy-oriented prediction prioritizes optimal data fitting rather than decision performance in the downstream optimization [22]. Such approaches are optimal only when the relationships between prediction error and decision error are symmetric and monotonic. To achieve better decision performances for problems exhibiting asymmetric and non-monotonic relationships, developing cost-oriented PTO framework becomes essential, where the prediction models are optimized under the goal of minimizing decision regret.

The idea of cost-oriented prediction originates from quantitative finance [23], and has recently demonstrated better performance over accuracy-oriented PTO framework in various practical power systems problems. Notable examples include day-ahead market operation [24, 25] and energy storage arbitrage [26] based on predicted electricity price, unit commitment [27–30], energy dispatch [21,31], voltage regulation [32], inertia resources deployment [33] and distributed energy management [18] given predicted power outputs of RES. Such methodology has also been extended to interval forecasting [34]. These studies collectively confirm that whenever prediction errors translate into decision outcomes in a non-symmetric or non-monotonic manner, the cost-oriented framework offers practical advantages. For the joint electricity-computation dispatch problem, there exhibit asymmetric relationships between prediction error and decision error, necessitating the deployment of cost-oriented PTO framework. The physical interpretation behind such relationships would be discussed later in Section 5.2. Nevertheless, no existing work has so far proposed a cost-oriented PTO framework for the joint electricity-computation dispatch problem.

The second challenge lies in designing an efficient iterative algorithm. In some existing studies, the joint electricity-computation dispatch problem is formulated as a bi-level optimization problem. This structure allows for reformulation into a single-level problem using Karush–Kuhn–Tucker conditions [35]. However, linearizing the complementary slackness conditions introduces extensive binary variables, which significantly increase computational complexity. Moreover, since SO and DCO are separate stakeholders, transforming the problem into a single-level formulation requires that one party should access the other's private information, which is often impractical due to privacy concerns. Another classical approach to solve such bi-level problem is to utilize Benders decomposition [6,36], which typically assumes a leader-follower structure where the SO acts as a leader and the DCO responds accordingly as a follower. This formulation requires the SO to pass primal variables to the DCO, while in our concerned problem, the exchanged information is LMP, intrinsically the dual variables of power balance constraints. This mismatch in the information flow potentially limits the practical applicability of Benders decomposition based methods in the problem concerned.

Another commonly used approach is iterative optimization [37], in which the SO and the DCO update their strategies in an iterative manner. SO and DCO negotiate until a consensus is reached [38]. However, iterative optimization methods may suffer from the failure to converge due to the oscillation between successive iterations [13]. Although recent work has improved convergence performance by narrowing the oscillation interval using binary search technique [3], the convergence speed remains relatively slow. Thus, designing an iterative optimization method towards faster convergence emerges as an essential problem.

1.3. Contributions

To address these research gaps, this paper offers a novel solution towards cost-optimal joint electricity-computation management. Driven by the booming demand for AI and cloud computing, the joint management is becoming imperative. Our methodology would serve as theoretical guidance for future development. The major contributions are summarized as follows.

1. First, a novel cost-oriented PTO decision framework is developed for joint electricity-computation management. Compared to conventional accuracy-oriented approaches, the proposed framework realizes statistically lower operational costs in extensive scenarios. Specially, the prediction models are optimized through solving regret minimization problems, which distinguishes our framework from existing approaches. Also, heterogeneous uncertainties arising from power outputs of RES and workloads of DC have been comprehensively considered.
2. In addition, a dual-boundary feedback-embedded adaptive iterative algorithm is proposed to guarantee faster convergence. In contrast to existing methods that rely on repeated binary search to approximate equilibrium, the proposed algorithm achieves provably faster convergence by adaptively tightening the boundaries of oscillation interval. Moreover, the proposed algorithm only requires limited information disclosure for SO and DCO, which preserves the privacy.

Section 2 introduces the decision framework for the joint dispatch problem. The optimization problems and prediction models are formulated in Section 3. The proposed iterative algorithm and convergence analysis are presented in Section 4. Section 5 offers an analysis of the simulation results and Section 6 draws the conclusion.

2. Joint electricity-computation management framework

As stated in Section 1, more comprehensive interactions between SO and DCO are emerging. Therefore, in modern power systems containing RES and DC, it is essential to optimize generation and reserve plans to maintain the supply-demand balance of active power. From the perspective of DCO, the goal of server dispatch is to fulfill computing requests from end-users under the quality-of-service (QoS) requirements. However, SO and DCO face inherent uncertainties from power outputs of RES and workload arrivals, respectively. To improve decision quality under uncertainties, a cost-oriented PTO framework is adopted, as illustrated in Fig. 1.

In the day-ahead dispatch stage, prediction models are first trained over historical scenarios. In Step 1-1, given the realizations of power outputs of RES and workload, the SO and DCO jointly optimize their decisions with full information of the historical scenarios. The SO solves an economic dispatch problem (P1), and discloses the dual multipliers of the nodal power balance constraints, i.e., LMP, to the DCO. The DCO then optimizes the server dispatch strategies accordingly in (D1). Specially, the spatial-temporal flexibility of computing workloads could be employed to minimize operational costs. The resulting total power consumption P^{dc} is then fed back to the SO to solve (P1). The two entities iteratively update their decisions until convergence. The operation cost obtained from Step 1-1 serves as the labels for the training of

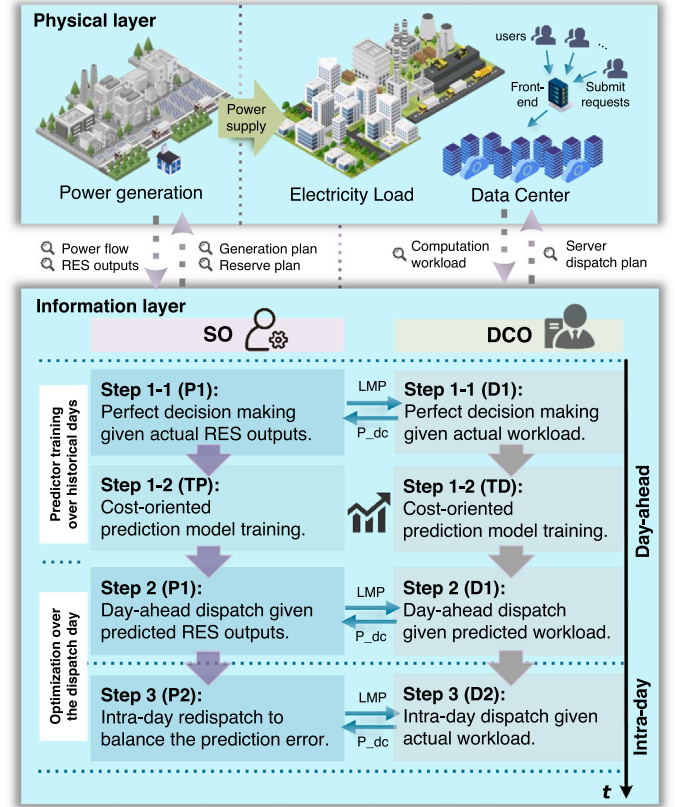


Fig. 1. Joint electricity-computation decision-making framework.

cost-oriented prediction models. As stated in Step 1-2, the training is performed by solving (TP) and (TD) for SO and DCO, respectively. After the training of the prediction models, the SO and DCO perform joint day-ahead dispatch over the testing scenario in Step 2, where (P1) and (D1) are solved iteratively until convergence. Due to inevitable prediction errors, the SO must perform redispatch in the intraday stage in Step 3, where generation strategies are adjusted based on the actual power outputs of RES by solving (P2). The DCO responds by solving (D2), and both entities continue to iterate until convergence.

Fig. 1 illustrates the decision-making process for one testing scenario. In practice, this framework operates in a rolling horizon manner [27]. For the k -th test day, historical data from the past h days is used for training, and the time window advances forward, i.e., using data from day $k - h + 1$ to day k to prepare for day $k + 1$, and so on.

3. Problem formulation

Let Ω_B and Ω_L denote the set of buses and transmission lines within the power system, respectively. $\Omega_G, \Omega_R, \Omega_D$ are respectively the sets of TG, RES and DC. \mathcal{T} is the set of time intervals of the dispatch day.

3.1. Optimization model of SO

At the day-ahead stage, SO aims to minimize the total operation cost C_P^{DA} given predicted power output of RES $\hat{P}_{i,t}^r$. The active power supply-demand balance constraint is modeled as follows.

$$P_{i,t}^G + \hat{P}_{i,t}^r + \sum_{(i,j) \in \Omega_L} P_{ij,t}^f + P_{i,t}^{sh} = P_{i,t}^L + P_{i,t}^{dc} + P_{i,t}^{cu} : \gamma_{i,t}, \quad \forall i, j \in \Omega_B, \forall t \quad (1)$$

where $P_{i,t}^G$ is the scheduled power generation strategy for each TG. $P_{ij,t}^f$ is the power flow from bus j to i . $P_{i,t}^L$ is the fixed load of each bus. $P_{i,t}^{dc}$ is the power consumption of DC. $P_{i,t}^{sh}$ and $P_{i,t}^{cu}$ are respectively the load

shedding and RES curtailment of the i -th bus, which will cause penalties. $\gamma_{i,t}$ represents the LMP of the i -th bus, which is intrinsically the dual variable of the power balance constraint (1).

The power output limits (2), ramping constraints of TGs (3) and the reserve dispatch requirements (4)-(6) could be modeled as follows:

$$R_{i,t}^- \leq P_{i,t}^G \leq \bar{P}_{i,t}^G - R_{i,t}^+, \quad \forall i \in \Omega_G, \forall t \quad (2)$$

$$-\Lambda_i^- \leq P_{i,t}^G - P_{i,t-1}^G \leq \Lambda_i^+, \quad \forall i \in \Omega_G, \forall t \quad (3)$$

$$0 \leq R_{i,t}^+ \leq v_{i,t}^+ \bar{P}_{i,t}^G, \quad \forall i \in \Omega_G, \forall t \quad (4)$$

$$0 \leq R_{i,t}^- \leq v_{i,t}^- \bar{P}_{i,t}^G, \quad \forall i \in \Omega_G, \forall t \quad (5)$$

$$\sum_{i \in \Omega_G} R_{i,t}^+ \geq \Pi_t^+, \quad \sum_{i \in \Omega_G} R_{i,t}^- \geq \Pi_t^-, \quad \forall t \quad (6)$$

where $R_{i,t}^+$ and $R_{i,t}^-$ are the scheduled upward and downward spinning reserve capacity, respectively. $\bar{P}_{i,t}^G$ is the rated power of the i -th TG. Λ_i^+ and Λ_i^- are the upward and downward ramping limits, respectively. $v_{i,t}^+$ and $v_{i,t}^-$ are the corresponding upper bounds of the upward and downward reserved capacity. Π_t^+ and Π_t^- are respectively the upward and downward reserve capacity requirements for the t -th time interval.

Credit to its convex feature and acceptable levels of error [13], the DC power flow constraint (7) is utilized in this work, where $\theta_{i,t}$ is the phase angle of i -th bus, X_{ij} is the line reactance. (8) limits the line power flow, where \bar{P}_{ij}^f is the transmission capacity. (9)-(10) ensure the non-negativity of load shedding and RES curtailment.

$$P_{ij,t}^f = (\theta_{i,t} - \theta_{j,t})/X_{ij}, \quad \forall (i,j) \in \Omega_L, \forall t \quad (7)$$

$$-\bar{P}_{ij}^f \leq P_{ij,t}^f \leq \bar{P}_{ij}^f, \quad \forall (i,j) \in \Omega_L, \forall t \quad (8)$$

$$0 \leq P_{i,t}^{\text{sh}} \leq P_{i,t}^L, \quad \forall i \in \Omega_B, \forall t \quad (9)$$

$$0 \leq P_{i,t}^{\text{cu}} \leq \hat{P}_{i,t}^r, \quad \forall i \in \Omega_R, \forall t \quad (10)$$

In summary, the day-ahead energy dispatch problem of SO ($P1$) could be modeled as follows:

$$(P1) \quad \min_{X_P^{\text{DA}}} C_P^{\text{DA}} = \sum_{i \in \mathcal{T}} \sum_{i \in \Omega_G} \gamma_i^G P_{i,t}^G + \sum_{i \in \mathcal{T}} \sum_{i \in \Omega_B} \gamma_i^{\text{sh}} P_{i,t}^{\text{sh}} + \sum_{i \in \mathcal{T}} \sum_{i \in \Omega_R} \gamma_i^{\text{cu}} P_{i,t}^{\text{cu}}, \quad (11)$$

s.t. (1) – (10).

where γ_i^G , γ_i^{sh} and γ_i^{cu} are respectively the cost coefficients of thermal power generation, load shedding and RES curtailment. For mathematical simplicity, an approximated linear generation cost function is assumed. Let X_P^{DA} collect the decision variables of problem ($P1$). The optimal solution of ($P1$) is expressed as $C_P^{\text{DA},*}(X_P^{\text{DA},*})$.

In the intraday stage, the actual power outputs of RES $P_{i,t}^r$ are available to the SO. The optimal generation strategies $P_{i,t}^{G,*}$ and upward/downward reserve capacities $R_{i,t}^{-,*}/R_{i,t}^{+,*}$ determined in ($P1$) are then served as input parameters for the intraday redispatch problem ($P2$). Accordingly, several constraints are modified as follows:

$$P_{i,t}^{G'} + P_{i,t}^r + P_{i,t}^{\text{sh}'} + \sum_{(i,j) \in \Omega_L} P_{ij,t}^f = P_{i,t}^L + P_{i,t}^{\text{de}'} + P_{i,t}^{\text{cu}'}, \quad \forall i \in \Omega_B, \forall t \quad (12)$$

$$-R_{i,t}^{-,*} \leq P_{i,t}^{G'} - P_{i,t}^{G,*} \leq R_{i,t}^{+,*}, \quad \forall i \in \Omega_G, \forall t \quad (13)$$

$$0 \leq P_{i,t}^{G'} \leq \bar{P}_{i,t}^G, \quad \forall i \in \Omega_G, \forall t \quad (14)$$

Then, ($P2$) could be formulated as follows:

$$(P2) \quad \min_{X_P^{\text{ID}}} C_P^{\text{ID}} = \sum_{i \in \mathcal{T}} \sum_{i \in \Omega_G} \gamma_i^G P_{i,t}^{G'} + \sum_{i \in \mathcal{T}} \sum_{i \in \Omega_B} \gamma_i^{\text{sh}} P_{i,t}^{\text{sh}'} + \sum_{i \in \mathcal{T}} \sum_{i \in \Omega_R} \gamma_i^{\text{cu}} P_{i,t}^{\text{cu}'}, \quad (15)$$

s.t. (3), (7) – (10), (12) – (14).

Similarly, the decision variables and optimal solution of ($P2$) are expressed as X_P^{ID} and $C_P^{\text{ID},*}(X_P^{\text{ID},*})$, respectively. The superscript ' denotes the variables for the intraday stage.

3.2. Optimization model of DCO

Since these two types of workloads are determined by separate users' demands, DCO may need to predict both interactive and batch workloads. Similarly to active power balance, the workload balance constraint is formulated as follows given predicted workloads [39]:

$$w_{i,t} = \widehat{W}_{i,t}^i + w_{i,t}^t + w_{i,t}^{\text{tr}} + w_{i,t}^{\text{cl}}, \quad \forall i \in \Omega_D, \forall t \quad (16)$$

$$w_{i,t}^t \geq 0, \quad w_{i,t}^{\text{cl}} \geq 0, \quad \forall i \in \Omega_D, \forall t \quad (17)$$

where $w_{i,t}$ is the workload that will be processed in each DC during each time interval t . $\widehat{W}_{i,t}^i$ is the interactive workload predicted. $w_{i,t}^t$ is the delay-tolerant batch workload processed immediately during t . $w_{i,t}^{\text{tr}}$ is the workload transferred among the DCs, making the DCs spatially flexible in power consumption. $w_{i,t}^{\text{cl}}$ is the workload processed in the cloud when local computation resources are inadequate. (17) ensures the non-negativity of $w_{i,t}^{\text{tr}}$ and $w_{i,t}^{\text{cl}}$.

For delay-tolerant batch workloads, DCO could decide whether to process them immediately or store them in the waiting queue. When the electricity price is high, the DCO can defer batch workloads and process them during lower-price periods [40], demonstrating the temporal flexibility of DCs. The associated constraints can be modeled as follows.

$$\widehat{W}_{i,t}^i = w_{i,t}^t + w_{i,t}^q, \quad \forall i \in \Omega_D, \forall t \quad (18)$$

$$w_{i,t}^q \geq 0, \quad \forall i \in \Omega_D, \forall t \quad (19)$$

$$W_{i,t+1}^q = W_{i,t}^q + w_{i,t}^q - w_{i,t}^t, \quad \forall i \in \Omega_D, \forall t \quad (20)$$

$$W_{i,1}^q = 0, \quad W_{i,T}^q = 0, \quad \forall i \in \Omega_D \quad (21)$$

$$0 \leq W_{i,t}^q \leq \bar{W}_i^q, \quad \forall i \in \Omega_D, \forall t \quad (22)$$

where $\widehat{W}_{i,t}^i$ is the predicted batch workload. $w_{i,t}^q$ is the stored batch workload. (18) describes the decision of DCO on the processing of batch workload. (19) ensures the non-negativity of $w_{i,t}^q$. $W_{i,t}^q$ is the existing batch workload stored in the queue. (20) is the dynamic continuity constraint of stored workload. (21) assumes that the queue is empty at the start of the dispatch day and must be fully cleared at the end. (22) defines the range of workload storage in the queue, where \bar{W}_i^q is the upper limit of storage capacity.

The upper limit of workload transfer among DCs $\bar{w}_{i,t}^{\text{tr}}$ is bounded by the communication bandwidth, as defined in (23). $w_{i,t}^{\text{tr}} \geq 0/w_{i,t}^{\text{tr}} < 0$ indicates that the workload is transferred from/to the i -th DC to/from other DCs. Also, similar to peer-to-peer energy sharing [18], the transferred workload at each time interval should be zero-sum, as defined in (24).

$$-\bar{w}_{i,t}^{\text{tr}} \leq w_{i,t}^{\text{tr}} \leq \bar{w}_{i,t}^{\text{tr}}, \quad \forall i \in \Omega_D, \forall t \quad (23)$$

$$\sum_{i \in \Omega_D} w_{i,t}^{\text{tr}} = 0, \quad \forall t \quad (24)$$

Moreover, as a computing services provider, the DC should satisfy QoS requirements, which are modeled by the M/M/1 queuing method in (25)[3]. Such a method could reflect the average response time and ensure acceptable service latency on an hourly dispatch time scale without introducing heavy computational complexity.

$$s_{i,t} \geq \frac{w_{i,t}}{\mu_i - 1/D}, \quad \forall i \in \Omega_D, \forall t \quad (25)$$

$$0 \leq s_{i,t} \leq \bar{s}_i, \quad \forall i \in \Omega_D, \forall t \quad (26)$$

where $s_{i,t}$ is the number of servers utilized at time interval t , mainly determined by the processed workload $w_{i,t}$. μ_i is the service rate of servers in the i -th DC. D is the maximum responding time. (26) defines the range of servers utilized, where \bar{s}_i is the total number of servers in the i -th DC.

Then, the total power consumption of a DC can be calculated as follows.

$$P_{i,t}^{dc} = P_i^b + [p_i^0 + (\xi_i - 1)p_i^{\max}]s_{i,t} + (w_{i,t} - w_{i,t}^{cl})(p_i^{\max} - p_i^0)/\mu_i, \quad \forall i \in \Omega_D, \forall t \quad (27)$$

where P_i^b is the base power consumption of the i -th DC. p_i^0 and p_i^{\max} are respectively the idle and peak power consumption of server in the i -th DC. $\xi_i \geq 1$ is the PUE, usually treated as a constant.

Finally, the day-ahead server dispatch problem of DCO (D1) could be formulated as follows:

$$(D1) \min_{X_D^{DA}} C_D^{DA} = \sum_{i \in \mathcal{I}} \sum_{t \in \Omega_D} \gamma_{i,t} P_{i,t}^{dc} + c^{tr} |w_{i,t}^{tr}| + c^{cl} w_{i,t}^{cl} \quad (28)$$

s.t. (16) – (27).

where X_D^{DA} collects decision variables for (D1). The first term in the objective is the power consumption cost. $\gamma_{i,t}$ is the LMP of the i -th bus. c^{tr} is the communication cost coefficient, and c^{cl} is the cost coefficient of using the cloud computing service.

Although DCO needs to reserve backup computing resources to deal with uncertainties, the response time from idle to full load is typically less than a few minutes [41]. Given the hourly time scale of dispatch in this work, it is reasonable to assume that servers can sufficiently respond within each time interval, and reserve computing capacity is not explicitly modeled in this work. Thus, unlike the SO, the DCO does not directly transmit decision variables from the day-ahead stage to intraday stage, which makes the day-ahead and intraday server dispatch problems share a similar formulation.

At the intraday stage, the workload balance constraints (16) and (18) could be modified as follows when the actual workload profiles are available.

$$w'_{i,t} = W_{i,t}^i + w'_{i,t} + w'_{i,t} + w'_{i,t}, \quad \forall i \in \Omega_D, \forall t \quad (29)$$

$$W_{i,t}^t = w'_{i,t} + w'_{i,t}, \quad \forall i \in \Omega_D, \forall t \quad (30)$$

Then, the intraday server dispatch problem of DCO (D2) could be formulated as follows. It is worth noting that while the primary objective of the DCO in this formulation is cost minimization. In our model, the service requirements are integrated as hard constraints.

$$(D2) \min_{X_D^{ID}} C_D^{ID} = \sum_{i \in \mathcal{I}} \sum_{t \in \Omega_D} \gamma'_{i,t} P_{i,t}^{dc'} + c^{tr} |w'_{i,t}| + c^{cl} w'_{i,t} \quad (31)$$

s.t. (17), (19) – (27), (29) – (30).

where X_D^{ID} collects decision variables for (D2) and the superscript $'$ denotes the variables for the intraday stage. Also, let $C_D^{ID,*}(X_D^{ID,*})$ denote the optimal solution of (D2).

3.3. Prediction model

Since the purpose of prediction is to better support decision-making [23], the primary objective of prediction model training is to fit the optimization oracle [22]. In addition, to reduce computational burden, the extreme learning machine (ELM) is adopted. ELM is a simple feed-forward neural network characterized by randomly pre-specified input weights and biases, eliminating the need for backward propagation [18]. Using ELM, the predicted power outputs of RES in vector form $\hat{P}_i^r \in \mathbb{R}^{|\mathcal{I}|}$ in (P1) are calculated as follows.

$$\hat{P}_i^r = (\Theta_i^r)^T \Phi_i^r, \quad \forall i \in \Omega_R, \quad (32)$$

where Θ_i^r is the learnable weights vector from hidden layer to output layer. H is the number of hidden neurons. Φ_i^r is the output of hidden

layer in ELM, which can be readily calculated given the input feature data, predefined input weight and bias vectors and the type of nonlinear activation function.

Similarly, the predicted interactive workload $\hat{W}_i^i \in \mathbb{R}^{|\mathcal{I}|}$ and batch workload $\hat{W}_i^t \in \mathbb{R}^{|\mathcal{I}|}$ for $\forall i \in \Omega_D$ in (D1) are calculated as follows:

$$\begin{cases} \hat{W}_i^i = (\Theta_i^{wi})^T \Phi_i^{wi}, & \forall i \in \Omega_D, \\ \hat{W}_i^t = (\Theta_i^{wt})^T \Phi_i^{wt}, & \forall i \in \Omega_D, \end{cases} \quad (33)$$

where Θ_i^{wi} and Θ_i^{wt} are the learnable weights vector of the ELMs for interactive workload and batch workload. Φ_i^{wi} and Φ_i^{wt} are the outputs of hidden layers of the ELMs for interactive workload and batch workload, respectively.

Traditionally, the goal of training an accuracy-oriented prediction model is to best fit the historical data. However, this does not necessarily lead to better decision quality [19]. To train a cost-oriented prediction model, the 'perfect decisions' of historical scenarios are first obtained to serve as labels, where the actual power outputs of RES and workloads are given.

For SO, the target is to optimize the learnable weights of ELM Θ_s^f under the objective of minimizing the decision error over $|\mathcal{S}|$ historical scenarios, i.e., the gap between the perfect operation costs obtained under actual power outputs of RES and the imperfect operation costs obtained under predicted power outputs. The training model can be formulated as follows.

$$(TP) \min_{X_{P,s}^{DA}, \Theta_s^f} \sum_{s \in \mathcal{S}} G_{P,s} + \lambda \|\Theta_s^f\|_1 \quad (34)$$

s.t. (1) – (10), (32),

$$G_{P,s} \geq C_{P,s}^{DA}(X_{P,s}^{DA}, \hat{P}_s^r) - C_{P,s}^{DA,*}(X_{P,s}^{DA,*}, P_s^r), \quad (35)$$

$$G_{P,s} \geq -[C_{P,s}^{DA}(X_{P,s}^{DA}, \hat{P}_s^r) - C_{P,s}^{DA,*}(X_{P,s}^{DA,*}, P_s^r)], \quad (36)$$

$$(1 - \kappa)P_s^r \leq \hat{P}_s^r \leq (1 + \kappa)P_s^r. \quad (37)$$

where $G_{P,s}$ is a slack variable representing the decision error of SO under scenario s . The second term in objective is to refrain from overfitting, where L1 regularization is commonly used [27]. λ is a hyperparameter that weights the importance of the regularization term. (35)–(36) are introduced to linearize the term of absolute value. P_s^r and \hat{P}_s^r are respectively the predicted and actual values of power outputs of RES under the s -th scenario in vector form. In addition, to ensure sufficient closeness of the predicted values to the actual values, (37) is introduced, where κ is the error range adjustment factor. This is because identical magnitudes of prediction error may lead to different decisions for certain problems. On the other hand, the same decision may correspond to different predicted values [22,26].

For the DCO, the corresponding training models can be similarly formulated as follows.

$$(TD) \min_{X_{D,s}^{DA}, \Theta_s^{wi}, \Theta_s^{wt}} \sum_{s \in \mathcal{S}} G_{D,s} + \lambda [\|\Theta_s^{wi}\|_1 + \|\Theta_s^{wt}\|_1] \quad (38)$$

s.t. (16) – (27), (33),

$$G_{D,s} \geq C_{D,s}^{DA}(X_{D,s}^{DA}, \hat{W}_s^i, \hat{W}_s^t) - C_{D,s}^{DA,*}(X_{D,s}^{DA,*}, W_s^i, W_s^t), \quad (39)$$

$$G_{D,s} \geq -[C_{D,s}^{DA}(X_{D,s}^{DA}, \hat{W}_s^i, \hat{W}_s^t) - C_{D,s}^{DA,*}(X_{D,s}^{DA,*}, W_s^i, W_s^t)], \quad (40)$$

$$(1 - \kappa)W_s^i \leq \hat{W}_s^i \leq (1 + \kappa)W_s^i, \quad (41)$$

$$(1 - \kappa)W_s^t \leq \hat{W}_s^t \leq (1 + \kappa)W_s^t. \quad (42)$$

where Θ_s^{wi} and Θ_s^{wt} are respectively the learnable weights of the ELM for interactive workload and batch workload under the s -th scenario. W_s^i

and \widehat{W}_s^i are respectively the predicted and actual values of interactive workload under the s -th scenario in vector form. W_s^t and \widehat{W}_s^t are respectively the predicted and actual values of batch workload under the s -th scenario in vector form.

By solving problems (TP) and (TD) separately, the cost-oriented prediction models could be obtained by SO and DCO. After that, the predicted power outputs of RES, interactive workload, and batch workload for the testing scenario could be calculated and taken into (P1) and (D1) for the day-ahead joint dispatch, followed by the intraday joint dispatch modeled in (P2) and (D2).

4. Algorithm design

In the context of joint management, privacy preservation is critical because these operators are distinct entities with separate commercial interests.

The sensitive information of the SO includes the detailed topology of the power grid, line parameters, and specific parameters of thermal generators and renewable units in the power system. The LMPs are shared with the DCO as public information, without revealing the internal physical parameters of the SO, which is consistent with the practice in the PJM market [42]. The private information of DCO includes computing workload arrival rates, server configuration details, QoS requirements, and specific power consumption characteristics. Keeping these details local protects the DCO's business strategies and operational status. The shared information of DCO is the total power consumption, which does not disclose the internal details of DCs.

In our designed algorithm, since only the LMP and the total power consumption of DCs are transmitted between SO and DCO during the iteration, the other private operational parameters are kept locally.

4.1. Iteration process

The detailed iterative algorithm is summarized in Algorithm 1. After initialization, SO and DCO optimize their respective dispatch problems sequentially. SO is assumed to be a leader in the iteration and DCO responds to the LMP declared by SO as a follower [43,44]. The SO then optimizes its strategies based on the total power consumption of DCs.

The residual is defined as $\delta^k = \|P^{dc,k} - P^{dc,k-1}\|$ in each iteration, namely the differences in the total power consumption of DCs between successive iterations. If the residual δ^k remains unchanged after two or more iterations, an oscillation is detected. Then, in Step 7, the m -th oscillation interval and its length are identified as $I_m = [P_{i,t}^m, \overline{P}_{i,t}^m]$ and $L_m = |I_m|$, respectively. The problem would fail to converge due to persistent oscillations without intervention. To address this, a privacy-preserving dual-boundary feedback-embedded adaptive iterative algorithm is specially proposed. Before convergence analysis, the following theoretical basis is necessary.

Definition 1. In each iteration, the response of SO/DCO upon receiving information from DCO/SO could be respectively expressed as the following mappings:

$$\gamma_{i,t}^{k+1} = \mathcal{R}_s(P_{i,t}^{dc,k}), \quad P_{i,t}^{dc,k+1} = \mathcal{R}_d(\gamma_{i,t}^{k+1}), \quad (43)$$

Definition 2. When $\{P_{i,t}^{dc,k}\}$ oscillates in closed, compact and convex interval I_m , $\exists K > 0$ that satisfies $P_{i,t}^{dc,k} \in I_m, \forall k \geq K$. Furthermore, the mapping between successive elements in the sequence could be expressed as:

$$P_{i,t}^{dc,k+1} = T_m(P_{i,t}^{dc,k}) = \mathcal{R}_d(\mathcal{R}_s(P_{i,t}^{dc,k})), \quad (44)$$

Proposition 1. There exists at least one equilibrium point in the oscillation interval, namely $\exists P_{i,t}^\circ \in I_m$.

Proof. Since the energy and server dispatch problems are linear programming, the optimal solution set is convex and the mapping $T_m : I_m \rightarrow I_m$ is continuous [45]. When an oscillation occurs, $\forall P_{i,t}^{dc,k} \in I_m$,

Algorithm 1 Privacy-preserving dual boundary feedback-embedded adaptive iterative algorithm.

- 1: Initialization: set iteration times $k = 1$, oscillation times $m = 0$, convergence tolerances $\epsilon > 0$, $P_{i,t}^{dc,0} = 0$;
- 2: **repeat**
- 3: Solve energy-reserve joint dispatch problem for SO;
- 4: Solve server dispatch problem for DCO;
- 5: **if** residual δ^k remains unchanged after two or more iterations **then**
- 6: Set $m = m + 1$;
- 7: Obtain m -th oscillation interval $I_m = [P_{i,t}^m, \overline{P}_{i,t}^m]$:

$$\underline{P}_{i,t}^m = \min(P_{i,t}^{dc,k}, P_{i,t}^{dc,k-1}),$$

$$\overline{P}_{i,t}^m = \max(P_{i,t}^{dc,k}, P_{i,t}^{dc,k-1});$$
- 8: **if** $m = 1$ **then**
- 9: $\tilde{P}_{i,t}^m = \tau \overline{P}_{i,t}^m + (1 - \tau) \underline{P}_{i,t}^m$;
- 10: **else**
- 11: Calculate the interval variation:

$$\Delta_{i,t}^{m,-} = |\underline{P}_{i,t}^m - \underline{P}_{i,t}^{m-1}|, \Delta_{i,t}^{m,+} = |\overline{P}_{i,t}^m - \overline{P}_{i,t}^{m-1}|,$$
- 12: Calculate the step length:

$$\rho_{i,t}^m = \Delta_{i,t}^{m,-} / (\Delta_{i,t}^{m,-} + \Delta_{i,t}^{m,+}),$$
- 13: Update the searching point:

$$\tilde{P}_{i,t}^m = \underline{P}_{i,t}^m + (\overline{P}_{i,t}^m - \underline{P}_{i,t}^m) \cdot \rho_{i,t}^m,$$
- 14: **end if**
- 15: Solve dispatch problem for SO based on $\tilde{P}_{i,t}^m$.
- 16: Solve dispatch problem for DCO to obtain an auxiliary DC power consumption $\tilde{P}_{i,t}^{dc}$;
- 17: Calculate $I'_m = [P_{i,t}^{m'}, \overline{P}_{i,t}^{m'}]$, where:

$$\underline{P}_{i,t}^{m'} = (\tilde{P}_{i,t}^{dc} > \tilde{P}_{i,t}^m) \tilde{P}_{i,t}^m + (\tilde{P}_{i,t}^{dc} \leq \tilde{P}_{i,t}^m) \underline{P}_{i,t}^m,$$

$$\overline{P}_{i,t}^{m'} = (\tilde{P}_{i,t}^{dc} < \tilde{P}_{i,t}^m) \tilde{P}_{i,t}^m + (\tilde{P}_{i,t}^{dc} \geq \tilde{P}_{i,t}^m) \overline{P}_{i,t}^m;$$
- 18: Add constraint for DCO: $\underline{P}_{i,t}^{m'} \leq P_{i,t}^{dc} \leq \overline{P}_{i,t}^{m'}$.
- 19: **end if**
- 20: Set $k = k + 1$;
- 21: **until** $\delta^k \leq \epsilon$ or $k \geq k_{\max}$.

$T(P_{i,t}^{dc,k}) \in I_m$. According to the Brouwer theorem [46], there exists at least one fixed-point $P_{i,t}^\circ \in I_m$ satisfying $T(P_{i,t}^\circ) = P_{i,t}^\circ$, which is the desired equilibrium point.

In our problem, the same $\gamma_{i,t}$ may correspond to multiple $P_{i,t}^{dc}$. As the LMP is not solely determined by $P_{i,t}^{dc}$, but also influenced by other parameters and constraints. In other words, different values of $P_{i,t}^{dc}$ potentially result in the same LMP. Therefore, the mapping $\mathcal{R}_s(\cdot)$ is not injective and there could be more than one equilibrium point. \square

To approach the equilibrium point, a searching point $\tilde{P}_{i,t}^m$ is selected within the oscillation interval I_m to explore which sub-interval contains the equilibrium point. For the first oscillation, the searching point can be selected randomly in the interval. Here, the midpoint (i.e., $\tau = 0.5$) is used as an example. For the following oscillations, $\tilde{P}_{i,t}^m$ is adaptively updated based on dual-boundary feedback information, as in Steps 11–13. Specifically, the step length is calculated using difference in the interval boundaries between successive oscillations.

The searching point is then passed to SO, and DCO updates its decision $\tilde{P}_{i,t}^{dc}$ based on the corresponding price signal issued by the SO. According to Proposition 1, if the updated $\tilde{P}_{i,t}^{dc}$ is smaller than $\tilde{P}_{i,t}^m$, the equilibrium lies in the lower half of the interval I_m . Conversely, the equilibrium is in the upper half. The DCO's optimization problem is then modified by adding a constraint that restricts $P_{i,t}^{dc}$ to the identified sub-interval. This constraint actually tightens the next oscillation interval I_{m+1} . This process is repeated until $L_m \rightarrow 0$, where the residual would be sufficiently small.

4.2. Convergence analysis

By Proposition 1, since the equilibrium is guaranteed to exist within I_m , convergence of the interval to a point implies convergence to the equilibrium, where the residual δ^k becomes sufficiently small, satisfying the criteria in Step 21.

Proposition 2. *Compared to the bisection-embedding method (Method #a) in [3], the sequence of interval length $\{L_m\}$ exhibits a faster convergence rate utilizing the proposed method (Method #b).*

Definition 3. Let $\eta_m = L_{m+1}/L_m$ denote the shrinking rate of interval length between successive oscillations.

Proof. Method #a has linear convergence rate as $\|L_{m+1} - L_m^*\| \leq \eta_a \|L_m - L_m^*\|$, where the shrinking rate $\eta_a \equiv 1/2$ and $L_m^* = 0$. As for Method #b, the length of updated interval L'_m is expressed as follows after Step 17:

$$L'_m = \begin{cases} \frac{\bar{P}_{i,t}^m - \tilde{P}_{i,t}^m}{\tilde{P}_{i,t}^m - \underline{P}_{i,t}^m} = (1 - \rho_{i,t}^m)L_m, & \tilde{P}_{i,t}^{dc} \geq \tilde{P}_{i,t}^m \\ \frac{\tilde{P}_{i,t}^m - \underline{P}_{i,t}^m}{\tilde{P}_{i,t}^m - \underline{P}_{i,t}^m} = \rho_{i,t}^m L_m, & \tilde{P}_{i,t}^{dc} < \tilde{P}_{i,t}^m \end{cases} \quad (45)$$

After adding the constraint, we have $L'_m = \sup\{L_{m+1}\}$. In other words, I_{m+1} should be no broader than I'_m . Therefore, the shrinking rate of Method #b satisfies:

$$\eta_b \leq \min(\rho_{i,t}^m, 1 - \rho_{i,t}^m) = \begin{cases} \rho_{i,t}^m, & \rho_{i,t}^m < 1/2 \\ 1 - \rho_{i,t}^m, & \rho_{i,t}^m > 1/2 \\ 1/2, & \rho_{i,t}^m = 1/2 \end{cases}, \quad (46)$$

(46) indicates that $\eta_b \leq \eta_a$ always holds. Thus, Method #b enjoys faster convergence performance $\|L_{m+1} - L_m^*\| \leq \eta_b \|L_m - L_m^*\|$. \square

Noted that, although there exist cases when $\rho_{i,t}^m = 1/2$, where two methods share the same convergence rate, the convergence speed of our proposed method still outperforms that of Method #a through extensive numerical simulations, as demonstrated in Fig. 9(a).

5. Cases studies

5.1. Parameters setup

As depicted in Fig. 2, a modified IEEE 30-bus system containing PV stations and DCs is employed for case studies. Parameters of generator and power network are obtained from MATPOWER. Since the SO usually retains more $R_{i,t}^+$ to mitigate higher load shedding penalties [47], the requirements for upward/downward reserves are set as the sum of 60%/40% $\hat{P}_{i,t}^r$ and 15%/10% of total load, respectively. $v_{i,t}^+$ and $v_{i,t}^-$ are set at 40%/60%, with both Π_i^+ and Π_i^- capped at 30% of the installed generation capacity. γ_i^{sh} and γ_i^{cu} are set at 10,000\$/MW and 300\$/MW, respectively [21]. Real PV and load data from Ausgrid are utilized as datasets after properly scaling [48]. Besides, PV output data across different spatiotemporal contexts is selected as input feature dataset for the training of prediction model in this study [27]. $\bar{w}_{i,t}^{lr}$ is set as 4000 req. c^{lr} and c^{cl} are set as 0.01\$/req and 0.5\$/req, respectively. The other parameters about DC operation are obtained from [3]. The convergence threshold ϵ and the prediction error adjusting factor κ are set as 1e-20 and 0.5, respectively. All simulations are implemented via Gurobi on Matlab 2024b with a MIP gap tolerance of 1e-4. The hardware environment is an Intel i7-14700F CPU with 20 cores and an NVIDIA GeForce RTX 4060Ti GPU.

To present comprehensive and comparative results, three methods are implemented: *Method #1*: cost-oriented PTO using ELM; *Method #2*: cost-oriented PTO using linear regression (LR); *Method #3*: accuracy-oriented PTO using fully connected neural network (FCNN), where the FCNN is first trained under accuracy metric, followed by joint dispatch. MAE and root mean squared error (RMSE) are selected to evaluate the prediction accuracy.

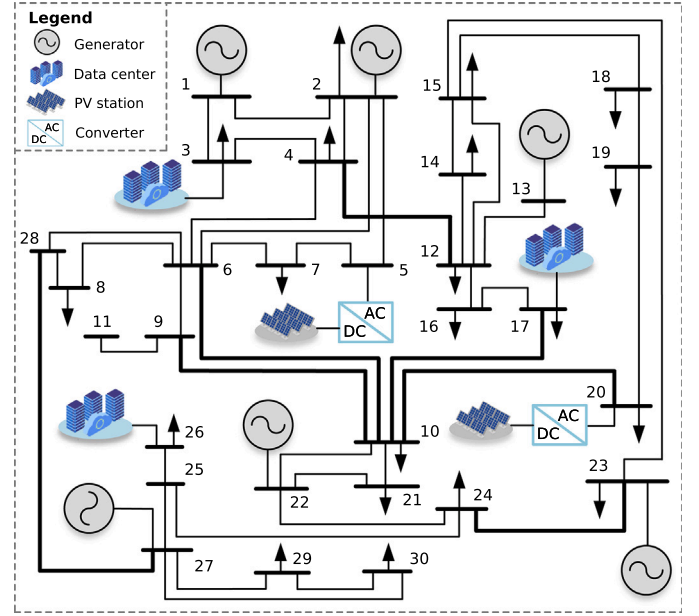


Fig. 2. Modified IEEE-30 bus system.

5.2. Decision and prediction performance

To comprehensively evaluate the cost-saving performance of the proposed framework for joint dispatch, total operational cost is defined as the sum of SO's and DCO's costs. Since the DCO's electricity consumption cost is actually the SO's revenue, this component is offset in the total cost calculation.

$$C^{ID} = C_P^{ID} + C_D^{ID} - \sum_{i \in T} \sum_{i \in \Omega_D} \gamma_{i,t}' P_{i,t}^{dc'} \quad (47)$$

Let $C^{ID,*}$ denote the perfect total operational cost of the testing scenario given actual power outputs of RES and workloads, the decision error could be defined as $[(C^{ID} - C^{ID,*})/C^{ID,*}] \times 100\%$.

Table 1 and Fig. 3 compare the decision error, prediction error, and average training time of the three methods. The number of scenarios is set as 426 due to data availability. As shown in Table 1, under the cost-oriented framework (Methods #1 and #2), the average decision error is lower than that of under the accuracy-oriented framework (Method #3), with reductions of approximately 1.38% and 1.52%, respectively. Although Method #2 yields a slightly lower median value of decision error than Method #1, it suffers from a higher standard deviation (Std), indicating a less stable decision performance.

In terms of prediction accuracy, Method #3 performs the best, achieving an MAE of 8.9784 MW, 0.6861×10^3 , and 0.4455×10^3 requests for power outputs of RES, batch workload, and interactive workload, respectively. Because the relationship between the RES prediction error and the total operation cost is asymmetric, a more accurate prediction does not necessarily lead to a more economic decision. The predicted values generated by Methods #1 and #2 tend to fall into a region that leads to smaller decision errors. Between the two cost-oriented methods, the ELM-based Method #1 outperforms the LR-based Method #2, benefiting from the superior nonlinear representation capability of the ELM. Additionally, Method #3 incurs a significantly higher training time due to the need for backward propagation, which requires approximately 19 s per scenario. In contrast, both Method #1 and Method #2 complete the training in less than 1 s, indicating a much lower computational burden. In summary, cost-oriented frameworks can achieve statistically better economic decision performance with significantly faster computation.

Table 1
Statistical results of decision error, prediction error and average training time.

Method	Decision error (%)			Prediction error (MW)		Prediction error (10^3 req)				Δt_r (s)
	Mean	Median	Std	P_r MAE	P_r RMSE	W_i MAE	W_i RMSE	W_i MAE	W_i RMSE	
#1	4.0171	1.9035	5.1528	10.081	11.798	3.5438	4.0748	2.3512	2.7058	0.8811
#2	3.8728	1.7652	6.2505	11.417	13.507	3.6659	4.1345	2.4803	2.7990	0.8794
#3	5.3997	1.1396	10.233	8.9784	10.862	0.6861	0.8508	0.4455	0.5545	18.925

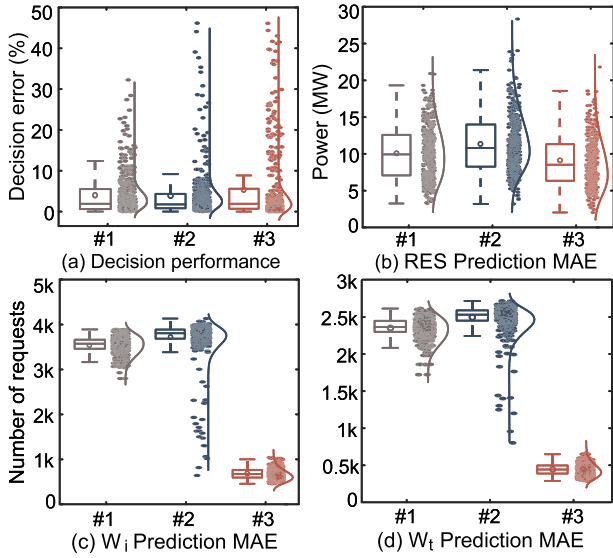


Fig. 3. Statistical results of decision and prediction error.

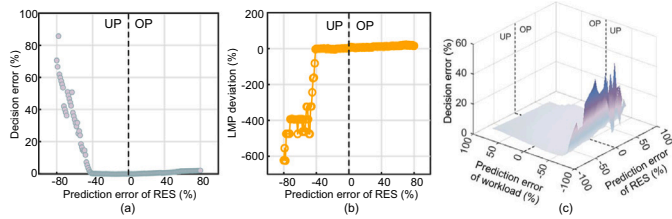


Fig. 4. Relationships between prediction error and decision error.

Fig. 4(a) reveals the asymmetric sensitivity of the SO’s decision error to RES prediction errors. Ideally, a perfect prediction yields zero decision error. However, under uncertainty, the error distribution is uneven.

In OP scenarios (Predicted > Actual): The SO underestimates day-ahead TG generation. To prevent high-penalty load shedding during the intraday stage, the SO typically schedules ample upward reserve capacities [21,47]. Consequently, the system is well-protected, keeping the realized decision error relatively low. In UP scenarios (Predicted < Actual): The SO schedules more day-ahead TG output. The primary risk here is renewable curtailment, which incurs a lower penalty compared to load shedding. Since the SO is less incentivized to reserve extensive downward capacities for low-penalty risks, curtailment occurs more frequently. As a result, UP scenarios statistically exhibit higher decision errors than OP scenarios, explaining the observed asymmetry.

Fig. 4(b) demonstrates the impact of RES prediction errors on intraday LMPs. When actual RES outputs significantly exceed forecasts (large UP errors), the surplus power drives LMPs to negative values, incentivizing users to consume excess generation. Notably, the LMP remains constant across certain ranges of prediction error. This step-like behavior confirms that a single LMP value can correspond to multiple distinct

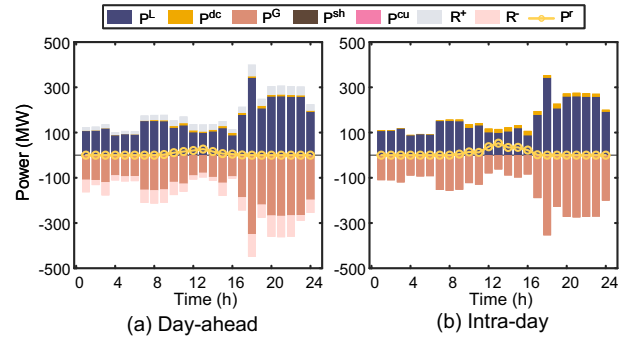


Fig. 5. Generation and reserve dispatch results.

levels of total DC power consumption, verifying the non-uniqueness of equilibrium points discussed in Section 4.1.

Fig. 4(c) visualizes the joint impact of RES and workload uncertainties. The relationship is complex and non-monotonic. RES prediction errors primarily dictate the allocation of reserve capacities, while workload prediction errors determine the day-ahead power consumption baseline of DCs. These two factors are coupled: the error in workloads indirectly shifts the operating point of thermal generation and reserves. Their interaction jointly determines the intraday LMPs and the final decision performance, resulting in the composite error surface shown in the figure.

5.3. Dispatch results on a typical scenario

Fig. 5(a) and (b) present the joint dispatch results of generation and reserve for the day-ahead and intraday stages in a typical testing scenario. Each term in the stack bar chart is the system-wide sum of each variable. On this day, the system load peaks between 18:00 and 24:00. It can be observed that upward reserve capacities are more than downward reserve capacities. In addition, the predicted RES outputs is lower than the actual values. Moreover, load shedding or RES curtailment penalties are not triggered, indicating the effectiveness of generation and reserve capacities dispatch strategies.

Fig. 6(a)–(c) illustrate the intraday workload processing strategies of three DCs in the testing system given the actual workload realizations. Since interactive workloads require immediate processing, they do not reflect temporal flexibility and are omitted from detailed discussion. The focus is on the batch workloads. In conjunction with the LMP signals shown in Fig. 7, it is evident that DCs delay processing during high-price periods by storing batch workloads in the waiting queues and concentrate processing during low-price hours (12:00–17:00). After evening peak loads, a processing surge occurs during 21:00–23:00 to ensure that the waiting queues are cleared by the end of the day. This demonstrates the temporal flexibility. Notably, when the system electricity load peaks at 18:00, the electricity price at DC #3 spikes sharply. In response, part of its workload is transmitted to DC #1, illustrating spatial flexibility in server dispatch. Fig. 6(d) shows the server provisioning status of each DC, which closely follows the trend of workloads. At 18:00, DC #3 significantly reduces server utilization to avoid high operation costs.

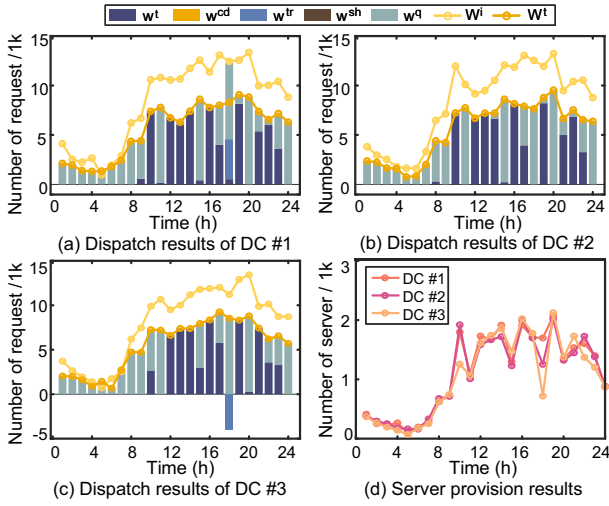


Fig. 6. Server dispatch results.

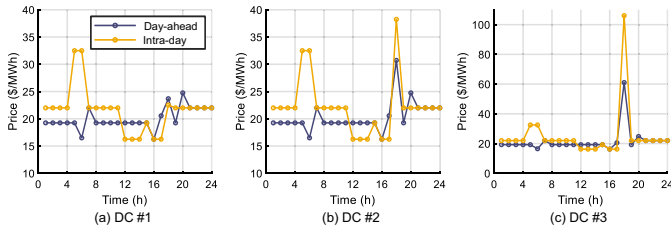


Fig. 7. LMP at day-ahead and intraday stage.

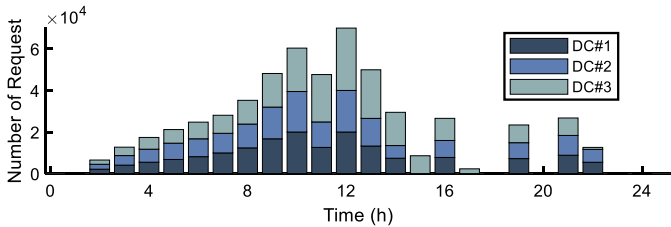


Fig. 8. Stored batch workload in the waiting queue.

Fig. 8 represents the accumulation of tasks in the waiting queues of the three DCs. First, all tasks in the queues are fully cleared before 24:00, ensuring sufficient buffer capacity for the operation of the next day. Moreover, by leveraging both electricity price and received workload information, the DCs gradually clear queued tasks during low-price periods, during which server utilization peaks. Conversely, during high-price hours, workloads are stored in the queues, keeping the power consumption low. These behaviors clearly demonstrate that the waiting queues serve as the key enabler of temporal flexibility.

5.4. Convergence performance

Fig. 9(a) compares the number of iterations required for convergence with approximately 1000 different combinations of prediction errors for the same testing scenario. The average, median, and standard deviation of the iteration times of the proposed dual-boundary feedback-embedded adaptive iterative algorithm (Method #b) are 16, 15, and 4.227, respectively. In contrast, for the existing bisection-embedded method (Method #a), the corresponding values are 20, 19, and 6.788. Although the maximum number of iterations for the proposed method is comparable to the existing method in some extreme cases, our method

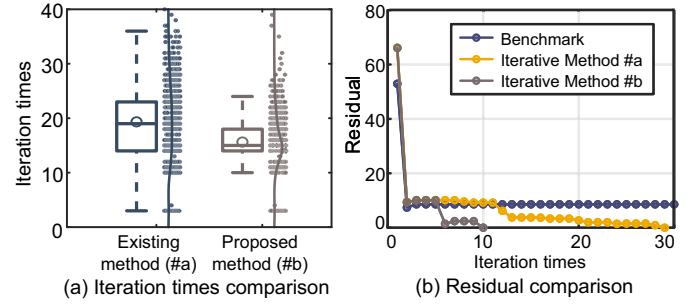


Fig. 9. Convergence performance comparison.

achieves a statistically reduction of approximately 20% in the average number of iterations. This strongly supports the effectiveness of the proposed approach and validates the theoretical analysis presented in Section 4.2. Furthermore, Fig. 9(b) shows the residual of a typical testing scenario. Without interval tightening, the solution keeps oscillating and fails to converge, whereas the proposed method achieves a remarkable faster convergence speed than the bisection-embedded method.

6. Conclusion

6.1. Major results and practical implications

To conclude, this paper aims to address the escalating interdependence between power systems and DCs, offering actionable insights for the sustainable integration of computing workloads into renewable-penetrated power systems. Specifically, to attain more economic joint electricity-computation dispatch strategies for both SO and DCO against heterogeneous uncertainties, a cost-oriented PTO framework is developed in this work. Moreover, a dual-boundary feedback-embedded adaptive iterative algorithm is proposed to ensure rigorous and faster convergence in iterative decision-making between SO and DCO. Major conclusions and practical implications drawn from extensive simulations are summarized as follows:

1. More accurate predictions do not necessarily yield lower operational costs due to the asymmetric risks of power shortages versus surpluses. The proposed cost-oriented PTO framework effectively captures this asymmetry, achieving an average cost reduction of 1.4%. For SO, this underscores the necessity of embedding downstream decision logic into upstream prediction models to avoid high-penalty events like load shedding.
2. The proposed dual-boundary feedback-embedded algorithm reduces the average iteration rounds by approximately 20% compared to standard binary search methods in [3]. It also proves that SO and DCO can reach a faster equilibrium without sharing sensitive data (e.g., grid topology or detailed workload profiles), thereby overcoming the primary barrier of commercial confidentiality.

6.2. Transferability and general applicability

While this study focuses on the coordination between power systems and DCs, the proposed methodologies possess significant transferability and general applicability to a broader range of energy management and optimization problems.

From the methodological perspective, the core contribution, the cost-oriented PTO framework, offers a novel paradigm for decision-making under uncertainty. Traditional accuracy-oriented forecasting often falls short in complex systems where the relationship between prediction errors and decision costs is asymmetric or even non-monotonic. Our framework effectively bridges this gap by embedding the downstream optimization cost directly into the training phase. Consequently, this

approach is not limited to the electricity-computation nexus but is highly applicable to other optimization problems in power systems, such as UC, optimal bidding strategies, etc. Furthermore, its potential extends to interdisciplinary fields like quantitative trading in finance and inventory management in logistics, where the economic penalty of over-forecasting differs from that of under-forecasting.

From the algorithmic perspective, the specific problem addressed in this paper can be viewed as a representative case of flexible loads participating in demand response. The iterative coordination algorithm developed herein relies on a generalized information exchange interface: the SO issues price signals (i.e., LMPs), and the flexible load operator responds with total power consumption. This decentralized mechanism does not require the disclosure of the internal physics of the load. Therefore, the proposed algorithm can be readily extended to coordinate other types of flexible resources, such as heating, ventilation, and air conditioning systems in smart buildings or electric vehicle charging stations [49]. As long as the flexible resources can adjust their consumption profile in response to price incentives, our privacy-preserving coordination algorithm provides a robust solution for achieving system-wide equilibrium.

6.3. Future work

Future work will focus on two directions to further bridge the gap between theory and practice: 1) Integrating end-to-end decision-focused learning architectures (e.g., Transformer-based forecasting models) to decision-making; and 2) Incorporating more granular computing models (e.g., G/G/1 queuing) to reflect diverse QoS requirements.

CRedit authorship contribution statement

Yibo Ding: Writing – original draft, Software, Methodology. **Xudong Li:** Writing – review & editing, Software. **Yuhong Zhao:** Visualization, Investigation. **Wenzhuo Shi:** Methodology, Investigation. **Cheng Lyu:** Software, Resources. **Jiaqi Ruan:** Data curation, Conceptualization. **Zhao Xu:** Writing – review & editing, Project administration, Funding acquisition.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

Data will be made available on request.

References

- [1] International Energy Agency. Energy and AI; 2025. <https://www.iea.org/reports/energy-and-ai>.
- [2] Armbrust M, Fox A, Griffith R, Joseph AD, Katz R, Konwinski A, Lee G, Patterson D, Rabkin A, Stoica I, et al. A view of cloud computing. *Commun ACM* 2010;53(4):50–8.
- [3] Yan D, Chow M-Y, Chen Y. Low-Carbon operation of data centers with joint workload sharing and carbon allowance trading. *IEEE Trans Cloud Comput* 2024;12(2):750–61. <https://doi.org/10.1109/TCC.2024.3396476>
- [4] Tran NH, Tran DH, Ren S, Han Z, Huh E-N, Hong CS. How geo-distributed data centers do demand response: a game-theoretic approach. *IEEE Trans Smart Grid* 2015;7(2):937–47.
- [5] Kwon S, Ntamo L, Gautam N. Demand response in data centers: integration of server provisioning and power procurement. *IEEE Trans Smart Grid* 2018;10(5):4928–38.
- [6] Liu S, Zhao T, Liu X, Li Y, Wang P. Proactive resilient day-ahead unit commitment with cloud computing data centers. *IEEE Trans Ind Appl* 2022;58(2):1675–84.
- [7] Ding Y, Liu Y, Ruan J, Sun X, Shi W, Xu Z. Carbon management for modern power system: an overview. *Smart Power Energy Security* 2025;1(1):12–24. <https://doi.org/10.1016/j.spes.2024.06.001>
- [8] Das P, Kayal P. An advantageous charging/discharging scheduling of electric vehicles in a PV energy enhanced power distribution grid. *Green Energy Intell Transp* 2024;3(2):100170. <https://doi.org/10.1016/j.geits.2024.100170>
- [9] Jin T, Bai L, Yan M, Chen X. Unlocking Spatio-Temporal flexibility of data centers in multiple regional Peer-to-Peer energy transaction markets. *IEEE Trans Power Syst* 2025;40(5):3914–27. <https://doi.org/10.1109/TPWRS.2025.3532208>
- [10] Ye Y, Ma D, Wu Y, Hu H, Zhang X, Liu C, Xu D. Harvesting spatial-temporal load migration flexibility of data centers: a chance-constrained bi-level optimization model with endogenously formed risk-reflective locational prices. *Appl Energy* 2026;402:126971. <https://doi.org/10.1016/j.apenergy.2025.126971>
- [11] Tsiligkaridis A, Andrianesis P, Coskun AK, Caramanis MC, Paschalidis IC. Distributed economic dispatch in power networks incorporating data center flexibility. *IEEE Trans Sustain Comput* 2025;10(4):768–83.
- [12] Zeng B, Zhou Y, Xu X, Cai D. Bi-level planning approach for incorporating the demand-side flexibility of cloud data centers under electricity-carbon markets. *Applied Energy* 2024;357:122406. <https://doi.org/10.1016/j.apenergy.2023.122406>
- [13] Ruan J, Zhu Y, Cao Y, Sun X, Lei S, Liang G, Qiu J, Xu Z. Privacy-preserving bi-level optimization of internet data centers for electricity-carbon collaborative demand response. *IEEE Internet Things J* 2024;11(14):24948–59. <https://doi.org/10.1109/JIOT.2024.3391762>
- [14] Yang L, Chen X, Fang X, Yang Q. Optimal coordinated management of integrated electricity-heat-computation systems in geographically distributed data centers. *Energy* 2025;332:137232. <https://doi.org/10.1016/j.energy.2025.137232>
- [15] Zhang Y, Zou B, Jin X, Luo Y, Song M, Ye Y, Hu Q, Chen Q, Zamboni AC. Mitigating power grid impact from proactive data center workload shifts: a coordinated scheduling strategy integrating synergistic traffic - data - power networks. *Appl Energy* 2025;377:124697. <https://doi.org/10.1016/j.apenergy.2024.124697>
- [16] Gmach D, Rolia J, Cherkasova L, Kemper A. Workload analysis and demand prediction of enterprise data center applications. In: 2007 IEEE 10th International Symposium on workload characterization. IEEE; 2007. p. 171–80.
- [17] Fan W, Fan Y, Liu P, Wang Y, Tong F, Yi B, Yao X. Distributionally robust optimization scheduling model for electric power and computing power coordination considering spatiotemporal response. *Appl Energy* 2025;402:126895. <https://doi.org/10.1016/j.apenergy.2025.126895>
- [18] Ding Y, Sun X, Shi W, Ruan J, Chen J, Zhao Y, Xu Z. Optimal distributed energy management for local energy community: a electricity regret oriented smart predict and optimize approach. *IEEE Trans Ind Informatics* 2025;21(12):9690–700.
- [19] Zhang H, Li R, Du Q, Tao J, Pineda S, Kariniotakis G, Camal S, Monroc CB, Sun M, Wan C, et al. Decision-focused learning for power system decision-making under uncertainty. *IEEE Trans Power Syst* 2026;41(1):307–23. <https://doi.org/10.1109/TPWRS.2025.3597806>
- [20] Wang Y, Wang J, Zhang H, Song J. Bridging prediction and decision: advances and challenges in data-driven optimization. *Nexus* 2025;2(1):100057. <https://doi.org/10.1016/j.nynex.2025.100057>
- [21] Ding Y, Sun X, Zhao Y, Lyu C, Chen J, Li X, Shi W, Ruan J, Xu Z. Multi-stage electricity-carbon joint management with decision-oriented predict-then-optimize method. *J Mod Power Syst Clean Energy* 2025. <https://doi.org/10.35833/MPCE.2025.000388>
- [22] Elmachtoub AN, Grigas P. Smart “predict, then optimize”. *Manag Sci* 2022;68(1):9–26.
- [23] Bengio Y. Using a financial training criterion rather than a prediction criterion. *Int J Neural Syst* 1997;8(4):433–43.
- [24] Alrasheedi AF, Alnowibet KA, Alshamrani AM. A smart predict-and-optimize framework for microgrid's bidding strategy in a day-ahead electricity market. *Electr Power Syst Res* 2024;228:110016.
- [25] Zhang Y, Wen H, Bian Y, Shi Y. Improving sequential market coordination via value-oriented renewable energy forecasting. [arXiv preprint] arXiv:2405.09004. 2024.
- [26] Sang L, Xu Y, Long H, Hu Q, Sun H. Electricity price prediction for energy storage system arbitrage: a decision-focused approach. *IEEE Trans Smart Grid* 2022;13(4):2822–32.
- [27] Chen X, Yang Y, Liu Y, Wu L. Feature-driven economic improvement for network-constrained unit commitment: a closed-loop predict-and-optimize framework. *IEEE Trans Power Syst* 2021;37(4):3104–18.
- [28] Chen X, Liu Y, Wu L. Towards improving unit commitment economics: an Add-On tailor for renewable energy and reserve predictions. *IEEE Trans Sustain Energy* 2024;15(4):2547–66. <https://doi.org/10.1109/TSTE.2024.3426337>
- [29] Wu HT, Ke DP, Song L, Liao SY, Xu J, Sun YZ, Fang K. A novel stochastic unit commitment characterized by closed-loop forecast-and-decision for wind integrated power systems. *IEEE Trans Power Syst* 2024;39(2):2570–86.
- [30] Ghazanfariharandi M, Mieth R. Value-oriented forecast combinations for unit commitment. *IEEE Control Syst Lett* 2025;9:1466–71. <https://doi.org/10.1109/LCSYS.2025.3579578>
- [31] Zhang Y, Jia M, Wen H, Bian Y, Shi Y. Toward value-oriented renewable energy forecasting: an iterative learning approach. *IEEE Trans Smart Grid* 2025;16(2):1962–74. <https://doi.org/10.1109/TSG.2024.3503554>
- [32] Sang L, Xu Y, Long H, Wu W. Safety-aware semi-end-to-end coordinated decision model for voltage regulation in active distribution network. *IEEE Transactions on Smart Grid* 2022;14(3):1814–26.
- [33] Zhang H, Li R, Chen Y, Chu Z, Sun M, Teng F. Risk-aware objective-based forecasting in inertia management. *IEEE Trans Power Syst* 2024;39(2):4612–23. <https://doi.org/10.1109/TPWRS.2023.3305452>
- [34] Zhang Y, Wen H, Wu Q. A contextual bandit approach for value-oriented prediction interval forecasting. *IEEE Trans Smart Grid* 2023;15(2):2271–81.
- [35] Dempe S, Franke S. Solution of Bilevel optimization problems using the KKT approach. *Optimization* 2019;68(8):1471–89.
- [36] Niu T, Hu B, Xie K, Pan C, Jin H, Li C. Spacial coordination between data centers and power system considering uncertainties of both source and load sides. *Int J Electr Power Energy Syst* 2021;124:106358.
- [37] Zhang W, Wei W, Chen L, Zheng B, Mei S. Service pricing and load dispatch of residential shared energy storage unit. *Energy* 2020;202:117543.

- [38] Han O, Ding T, Mu C, Jia W, Ma Z. Coordinative optimization between multiple data center operators and a system operator based on two-level distributed scheduling algorithm. *IEEE Internet Things J* 2023;10(9):7517–27.
- [39] Zhong W, Su W, Huang X, Kang J, Yuen C, Deng R, Zhang Y, Xie S. Joint energy-computation management for electric vehicles under coordination of power distribution networks and computing power networks. *IEEE Trans Smart Grid* 2025;16(2):1549–61. <https://doi.org/10.1109/TSG.2024.3498945>
- [40] Hall S, Micheli F, Belgioioso G, Radovanović A, Dörfler F. Carbon-aware computing for data centers with probabilistic performance guarantees. *arXiv:2410.21510*. 2025.
- [41] Gu C, Li Z, Huang H, Jia X. Energy efficient scheduling of servers with multi-sleep modes for cloud data center. *IEEE Trans Cloud Comput* 2018;8(3): 833–46.
- [42] Hogan WW. Strengths and weaknesses of the PJM market model. In: *Handbook on electricity markets*. Edward Elgar Publishing; 2021. pp. 182–204.
- [43] Wang Y, Lin X, Pedram M. A Stackelberg game-based optimization framework of the smart grid with distributed PV power generations and data centers. *IEEE Trans Energy Convers* 2014;29(4):978–87.
- [44] Wang K, Ye L, Yang S, Deng Z, Song J, Li Z, Zhao Y. A hierarchical dispatch strategy of hybrid energy storage system in internet data center with model predictive control. *Appl Energy* 2023;331:120414.
- [45] Boyd S, Vandenberghe L. *Convex optimization*. Cambridge University Press; 2004.
- [46] Smart DR. *Fixed point theorems*, vol. 66. Cup Archive; 1980.
- [47] Yang L, Chen S, Dong Z. Multi-period locally-facet-based MIP formulations for unit commitment problems. *IEEE Trans Power Syst* 2022;38(4):3733–47.
- [48] Ausgrid. *Electricity usage research; 2024*. <https://www.ausgrid.com.au/Industry/Our-Research/Data-to-share/Electricity-Research>.
- [49] Ding Y, Sun X, Ruan J, Shi W, Wu H, Xu Z. Customized decentralized autonomous organization based optimal energy management for smart buildings. *Appl Energy* 2024;376:124223.