

# Digital Twin Empowered PV Power Prediction

Xiaoyu Zhang, Yushuai Li, *Member, IEEE*, Tianyi Li, *Member, IEEE*, Yonghao Gui, *Member, IEEE*, Qiuye Sun, *Senior Member, IEEE*, and David Wenzhong Gao, *Fellow, IEEE*

**Abstract**—The accurate prediction of photovoltaic (PV) power generation is significant to ensure the economic and safe operation of power systems. To this end, the paper establishes a new digital twin (DT) empowered PV power prediction framework that is capable of ensuring reliable data transmission and employing the DT to achieve high accuracy of power prediction. With this framework, considering potential data contamination in the collected PV data, a generative adversarial network is employed to restore the historical dataset, which offers a prerequisite to ensure accurate mapping from the physical space to the digital space. Further, a new DT-empowered PV power prediction method is proposed. Therein, we model a DT that encompasses a digital physical model for reflecting the physical operation mechanism and a neural network model (i.e., a parallel network of convolution and bidirectional long short-term memory model) for capturing the hidden spatiotemporal features. The proposed method enables the use of the DT to take advantages of the digital physical model and the neural network model, resulting in enhanced prediction accuracy. Finally, a real dataset is conducted to assess the effectiveness of the proposed method.

**Index Terms**—Photovoltaic power prediction, digital twin, hybrid prediction, data recovery.

## I. INTRODUCTION

WITH the increasing integration of PV power generation, its nonlinearity, periodicity, and volatility pose great challenges to the stable operation of power systems. The uncertainty of the PV power generation and the randomness of the power demand may lead to the imbalance between the power supply and demand. Accurate prediction models can mitigate the impacts of uncertainty of PV power

generation, improve power system stability, and reduce the maintenance costs of additional equipments [1]–[3].

Currently, several studies on PV power prediction have been proposed, which can be roughly divided into three categories: ① physical methods; ② statistical methods; and ③ artificial intelligence (AI)-based methods. The concept of physical methods is to use physical models to construct the relationship between PV power output and other factors such as numerical weather prediction (NWP) data [4], sky images [5], and satellite images [6]. The concept of statistical methods is to apply statistical principles such as Bayesian model averaging (BMA) [7], exponential smoothing [8], and autoregressive integrated moving average (ARIMA) [9] to extract correlations and variation patterns from historical data. Both physical and statistical methods have the advantage of obtaining stable PV power prediction. However, it is very difficult to establish a physical model that can obtain high-accuracy prediction results for every prediction scenario, since there exist several hidden features that are hard to capture via mechanism analysis. Meanwhile, statistical methods mainly focus on using historical data of power generation, which ignores weather conditions and results in limited prediction accuracy [10].

To cope with shortcomings of physical and statistical methods, the AI-based methods for PV power generation have been proposed and gained significant attentions. For instance, convolutional neural networks (CNNs) were used for extracting spatial features [11]–[13], while long short-term memory (LSTM) networks were used for extracting temporal features [14]–[16]. CNNs do not fully consider the temporal characteristics of the input data, and LSTM networks have limited ability to capture the causal relationships between input factors. To address this issue, hybrid models based on CNN and LSTM were proposed in [17] and [18]. Additionally, graph neural networks (GNNs) [19], [20], particularly graph convolutional networks (GCNs) [21], [22], are often combined with graph modeling methods to explore the causal relationships among input factors [23]. Nevertheless, GNNs and GCNs primarily focus on the neighboring information of nodes and have limited modeling capabilities for time-series data, which may pose challenges when dealing with the graphs with complex topological structures. Recently, generative adversarial networks (GANs) with capabilities in image restoration and data completion have also been used to address PV power prediction problems. In [24], a generator based on recurrent neural network (RNN) was employed to predict solar power, while a CNN discriminator

Manuscript received: May 24, 2023; revised: August 19, 2023; accepted: October 15, 2023. Date of CrossCheck: October 15, 2023. Date of online publication: November 17, 2023.

This work was supported by European Horizon 2020 Marie Skłodowska-Curie Actions (No. 101023244).

This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>).

X. Zhang and Q. Sun are with the School of Information Science and Engineering, Northeastern University, Shenyang 110004, China (e-mail: 2100689@stu.neu.edu.cn; sunqiuye@ise.neu.edu.cn).

Y. Li (corresponding author) is with the Department of Informatics, University of Oslo, Oslo 0316m, Norway (e-mail: yushuai.li@ieee.org).

T. Li is with the Department of Computer Science, Aalborg University, Aalborg 9220, Denmark (e-mail: tianyi@cs.aau.dk).

Y. Gui is with the Electrification and Energy Infrastructures Division, Oak Ridge National Laboratory, Oak Ridge, TN 37830, USA (e-mail: guiy@ornl.gov).

D. W. Gao is with the Department of Electrical and Computer Engineering, University of Denver, Denver, CO 80208, USA (e-mail: Wenzhong.Gao@du.edu).

DOI: 10.35833/MPCE.2023.000351



was utilized to enhance the prediction accuracy of the generator. However, when the training data are unbalanced or samples are scarce, it may lead to unreliable power prediction results generated by a GAN.

The aforementioned PV power prediction models are built up on the assumption that the dataset is complete [11]-[24]. In fact, varying degrees of pollution are usually observed in the collected measurement data, which may be caused by data logger failures, communication network failures, and inaccurate instruments, etc. Learning samples with these unexpected pollutants may lead to the bias in prediction results. To address this issue, the PV power generation is predicted based on a recursive long short-term memory network in [25], which considers the possible quality problems of the dataset. The missing data are estimated using a recursive process. However, the robustness of the method is reduced when the testing data loss rate significantly differs from the training data loss rate. Moreover, this method does not consider the continuous missing data patterns in the dataset. In [26], a super-resolution perception CNN was employed to recover missing data, and a stochastic configuration network (SCN) was utilized for PV power prediction. However, the quality of data recovery needs to be improved, which subsequently affects the accuracy of PV power prediction. In addition, even if the dataset is complete, there may exist data imbalance. In [27], data augmentation methods, e.g., noise injection, color space transformations, and mixing of images, were used to expand a small amount of sky image data under cloudy conditions. Meanwhile, the CNN was used to predict short-term PV output. In [28], the dataset was augmented with complementary exogenous features including the periodic properties of the production, altitude, azimuth, and irradiance of solar, and clear and overcast days, etc. Then, a hybrid neural network model was proposed to predict PV power generation.

The aforementioned AI-based methods have yielded remarkable outcomes. However, there exist two challenges. On the one hand, although the consideration of data recovery is presented in [25] and [26], the implementation of these methods is challenging; meanwhile, the quality of the recovered data is insufficient. To address this issue, a potential method is the utilization of the GAN. It is a framework for training parameter generation models, which is capable of learning arbitrarily complex probability distributions. The success of GANs in image restoration [29] and traffic data completion [30] serves as inspiration for applying GANs to learn the distribution of PV data, thus tackling the challenging task of recovering large-scale historical data. On the other hand, these AI-based methods [11]-[26] are predominantly developed using historical data such as power generation and meteorology data, without taking into account the specific physical characteristics of the PV system itself. It should be noted that the actual state of the PV power station, particularly the physical condition of the PV panels, significantly impacts the power generation process. To address this issue, the DT technology provides an alternative solution. The DT refers to the construction of a virtual system in a virtual space that utilizes physical models and operational historical data to ac-

curately represent and map the physical entity or process [31]. The advantages of DTs can be divided into the following three points: ① the DT is an accurate virtual simulation of a real-world entity, process, or system, which allows us to perform various tests, predictions, and optimizations in a virtual environment without actually manipulating real-world objects. It results in saved time and money [32]; ② the DT is capable of sharing information with the physical entities in real time, resulting in the information synchronization. This is helpful to make fast and accurate decision-making [33]; and ③ the DT can leverage the digital physical models to describe the behavior of real systems and combine with data-driven machine learning methods to achieve accurate modeling and prediction of real systems [34]. These advantages make DT become an innovative method and tool that can be applied to multiple fields. For instance, a two-level hierarchical learning process using the real-time model state stored on the DT server was proposed in [35], aiming to enhance the machine learning (ML)-based product design on a DT-aided Internet of Things (IoT) platform. An intelligent context-aware medical system was implemented in [36] by using a DT-based framework. Meanwhile, an electrocardiogram (ECG) heart rhythms classifier model was built by using ML to diagnose heart disease and detect heart problems. In addition, the DT was also used for product quality prediction [37], intelligent transportation [38], and smart home [39]. Although the DT has gained broad applications, it has not been applied to PV power prediction. Based on the advantages of the DT, our aim is to jointly create the digital physical model to reflect the inherent mechanism of PVs and use the neural networks to capture the hidden features that are hard to be modeled by physical model. In the sense, we can create a high-fidelity DT to reflect the reality well by taking advantages of physical knowledge and learned data knowledge, resulting in enhanced prediction accuracy. However, no attention has so far been paid to this aspect.

To tackle those challenges, this paper establishes the DT empowered PV power prediction framework and proposes a DT-empowered PV power generation prediction method. The main contributions are described as follows.

1) We propose a new DT-empowered PV power prediction framework, which is composed of a physical layer, a data transmission layer, a DT layer, and a service layer, while defining the detailed functionality of each layer. This is a universal reference framework that enables the integration of the DT to empower the PV power prediction.

2) To ensure accurate mapping from the physical to the digital space, a GAN is employed to restore the historical dataset, considering potential data contamination in the collected PV data. This restoration process serves as a prerequisite for reliable data analysis and prediction within the DT framework.

3) A DT-empowered PV power prediction method is proposed, where the DT is constructed with a digital physical model and a parallel CNN and bidirectional long short-term memory (CNN-BiLSTM) model. The proposed method captures both the physical operation mechanism and hidden spatiotemporal features, leveraging the strengths of both models

to increase the prediction accuracy.

The remainder of this paper is summarized as follows. Section II presents the DT-empowered PV power prediction framework. Section III provides the DT-empowered prediction method within the proposed framework. Section IV presents the simulations to evaluate the performance of the proposed method. Finally, Section V concludes this paper.

## II. DT-EMPOWERED PV POWER PREDICTION FRAMEWORK

Figure 1 shows the proposed DT-empowered PV power prediction framework, which is composed of a physical layer, a data transmission layer, a DT layer, and a service layer. The variables in Fig. 1 will be explained in Section III.

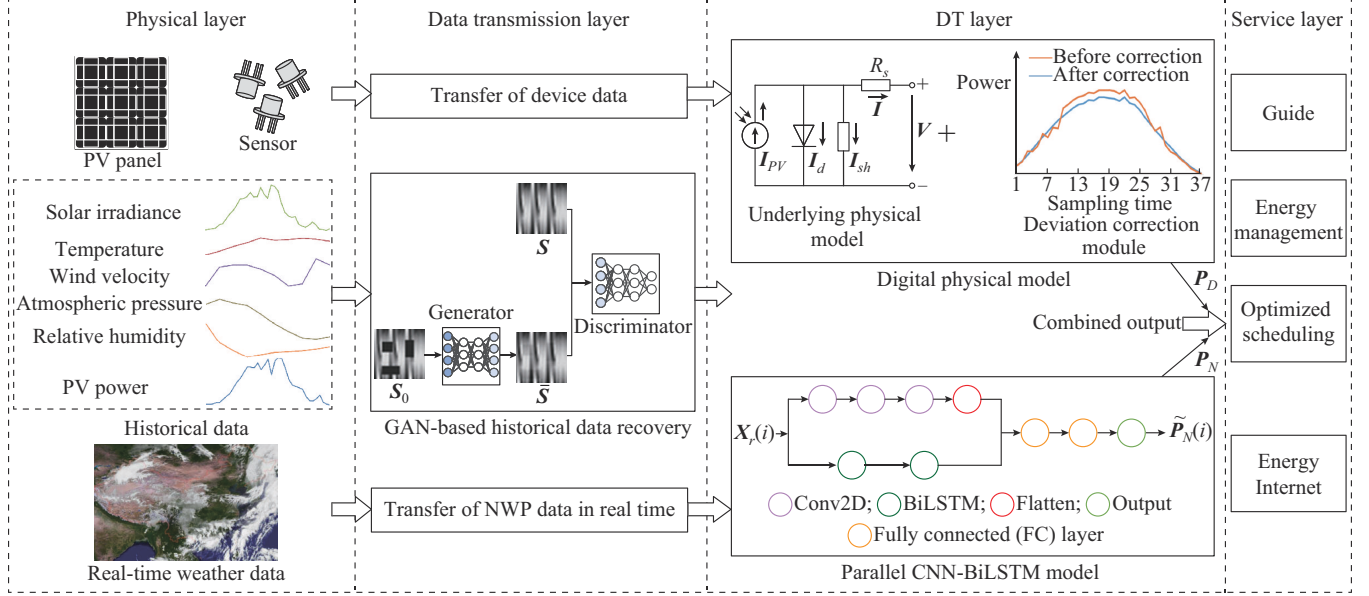


Fig. 1. Proposed DT-empowered PV power prediction framework.

### A. Physical Layer

This layer refers to physical objects in the real world such as PV panels and sensors. The layer will collect and store device parameters, the PV power generation data, and the meteorological data. Device parameters include short-circuit current  $I_{SC}$ , open-circuit voltage  $U_{OC}$ , data at the maximum power point (current  $I_m$ , voltage  $U_m$ , and the maximum power  $P_{om}$ ), and volt-ampere characteristic curve of the PV panel. According to different sampling time points, historical datasets can be expressed as:

$$\mathbf{D}_{xp} = [\mathbf{D}_{xp}(1), \mathbf{D}_{xp}(2), \dots, \mathbf{D}_{xp}(n_T)]^T \quad (1)$$

$$\mathbf{X} = [\mathbf{X}(1), \mathbf{X}(2), \dots, \mathbf{X}(n_T)]^T \quad (2)$$

$$\mathbf{P} = [\mathbf{P}(1), \mathbf{P}(2), \dots, \mathbf{P}(n_T)]^T \quad (3)$$

where  $\mathbf{D}_{xp}(j)$  is the historical data including temperature, wind speed, solar radiation, relative humidity, and PV power generation, etc., collected at the  $j^{\text{th}}$  sampling;  $\mathbf{D}_{xp}(j) = \{\mathbf{X}(j), \mathbf{P}(j)\}$ ;  $n_T$  is the time dimension of the data; and  $\mathbf{X}(j)$  and  $\mathbf{P}(j)$  are the historical meteorological data and historical PV power generation data collected at the  $j^{\text{th}}$  sampling time point, respectively.

### B. Data Transmission Layer

This layer serves as the connection channel between the physical and virtual spaces, enabling the collection and transmission of relevant data information from the PV power station. During data collection, the loss of data packets is possible, leading to incomplete time series data in the analysis of

historical PV power generation data. To address this issue, we propose the utilization of a GAN for data recovery, which will be discussed in Section III-A. The historical data restored by GAN and the parameter data of PV panels, sensors, and other devices are transmitted from the physical space to the virtual space at one time, participating in the construction of the DT model of the PV power station. The real-time weather data are transmitted in real time from the physical space to the virtual space, which enables participate in the power prediction of the DT layer.

### C. DT Layer

As the main part of this paper, this layer focuses on creating the DT model and using it to achieve PV power generation. In order to accurately reflect the real world and create a high-fidelity DT model, it is necessary to consider the physical characteristics of the PV system and extract the inherent relationships within the historical data simultaneously. In the virtual space, we set up a digital physical model that can reflect the physical operation mechanism and a parallel CNN-BiLSTM model to capture hidden spatiotemporal features. These components are combined using a fusion formula to accomplish the prediction of PV power. The detailed DT modeling process and the prediction procedure will be discussed in Section III-B and Section III-C, respectively.

### D. Service Layer

This layer receives the prediction results from the DT layer to meet diverse services such as: ① providing reference for energy dispatch and optimization; ② optimizing the

charging/discharging control for battery; and ③ facilitating demand response programs.

### III. DT-EMPOWERED PV POWER PREDICTION METHOD

Within the proposed framework, we propose the DT-empowered PV power prediction method that contains three phases: ① data preparation phase; ② DT modeling phase; and ③ power prediction phase. Figure 2 illustrates the overall flowchart. Next, we proceed to elaborate the design of each phase.

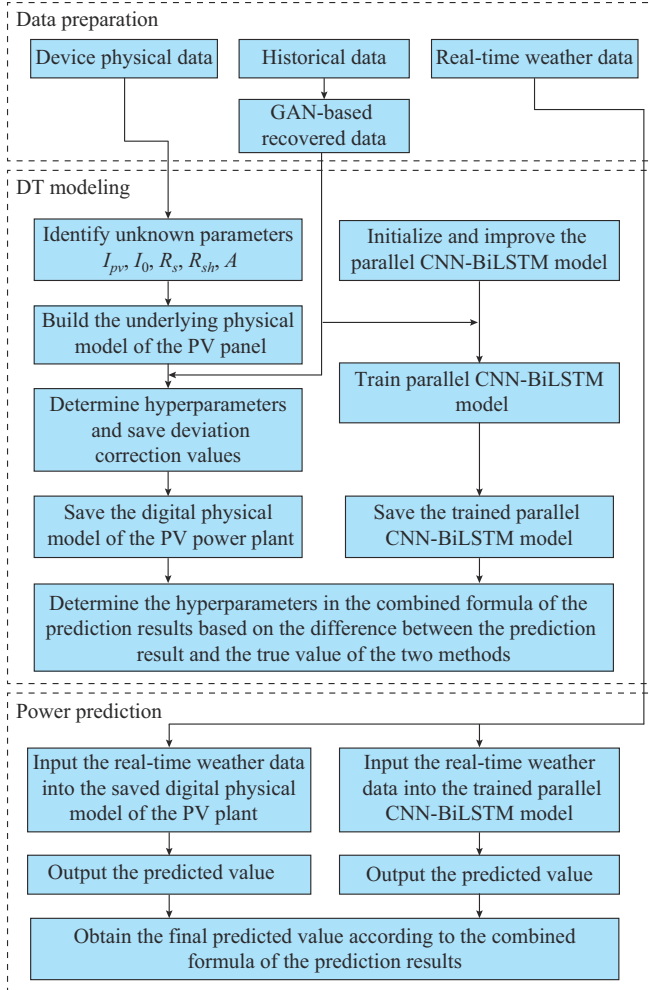


Fig. 2. Flowchart of DT empowered PV prediction method.

#### A. Data Preparation Phase

The data preparation phase is performed at the data transmission layer. In this phase, the data transmission layer retrieves pertinent data from the PV power station. The historical meteorological and power data are fed into the GAN. Subsequently, the recovered historical data and device parameters are transferred from the physical layer to the DT layer at one time to participate in the construction of DT model. Real-time weather data are transmitted from the physical layer to the DT layer, which are used for subsequent PV power prediction.

##### 1) Tensor Modeling of Historical Data

Historical weather and power data collected from PV sites

are combined and modeled as a tensor.

Firstly,  $c$  adjacent vectors  $\mathbf{D}_{xp1}, \mathbf{D}_{xp2}, \dots, \mathbf{D}_{xpc}$  are established for  $\mathbf{D}_{xp}$ , with a time interval of one sampling interval, i. e., 15 min. The adjacent vectors  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_c$  and  $\mathbf{P}_1, \mathbf{P}_2, \dots, \mathbf{P}_c$  adjacent to  $\mathbf{X}$  and  $\mathbf{P}$  are also established. For example,  $\mathbf{D}_{xp2} = \{\mathbf{X}_2, \mathbf{P}_2\} = [\mathbf{D}_{xp}(2), \mathbf{D}_{xp}(3), \dots, \mathbf{D}_{xp}(l+1)]^T$  has two adjacent vectors  $\mathbf{D}_{xp}(1) = \{\mathbf{X}(1), \mathbf{P}(1)\} = [\mathbf{D}_{xp}(1), \mathbf{D}_{xp}(2), \dots, \mathbf{D}_{xp}(l)]^T$  and  $\mathbf{D}_{xp}(3) = \{\mathbf{X}_3, \mathbf{P}_3\} = [\mathbf{D}_{xp}(3), \mathbf{D}_{xp}(4), \dots, \mathbf{D}_{xp}(l+2)]^T$ , where  $l$  is the time step of the vector. It means that each vector contains data of  $l$  sampling time points. The corresponding adjacent vectors  $\mathbf{M}_1, \mathbf{M}_2, \dots, \mathbf{M}_c$  of the mask matrix can also be obtained by using its procedure.

We represent  $\mathbf{S}_i$  as the  $i^{\text{th}}$  training sample input into the GAN. Then, we have:

$$\mathbf{S}_i = [\mathbf{D}_{xpi}, \mathbf{D}_{xp(i+1)}, \dots, \mathbf{D}_{xpi}, \mathbf{P}_j, \mathbf{P}_j, \mathbf{P}_j, \dots] \quad (4)$$

where  $\mathbf{S}_i \in \mathbb{R}^{l \times l}$ ,  $\mathbf{S}_i \in \mathbf{S}$ , and  $\mathbf{S}$  is the set of all training samples input into GAN;  $i < j$ , and  $j - i + 1 = \lfloor l/(n_x + 1) \rfloor$ .  $f = l\%(n_x + 1)$  is the number of padding vectors  $\mathbf{P}_j$ , where  $n_x$  is the number of meteorological factors, and  $\%$  is to obtain the value of remainder. An binary mask matrix  $\mathbf{M}$  with the same shape as  $\mathbf{S}$  is created to mark the positions of missing elements. For the missing elements in  $\mathbf{S}$ , the corresponding elements in  $\mathbf{M}$  are set to be 0; meanwhile, the rest of elements are set to be 1.

After modeling the historical data into a tensor, the problem of historical data recovery becomes the recovery of missing elements in the tensor.

##### 2) Data Recovery

To achieve effective data recovery, we employ the GAN consisting of a generator and a discriminator, which is capable of learning the temporal features of the data and capturing the intrinsic relationship between meteorological data and power data. The generator uses a CNN-based encoder-decoder structure. The encoder takes the missing dataset  $\mathbf{S}_0 = \mathbf{S} \odot \mathbf{M}$  as input and generates the latent feature representation of  $\mathbf{S}_0$ , where  $\odot$  is the dot product operator. Then, the decoder obtains the latent feature representation and outputs  $\bar{\mathbf{S}}$ , which includes the recovered part of the missing data. Furthermore, to maximize the utilization of the reliable data that are already presented in set  $\mathbf{S}_0$  during the data generation process, a U-net is adopted in the generator to enhance feature extraction. The discriminator takes the restored matrix  $\bar{\mathbf{S}}$  and the original complete matrix  $\mathbf{S}$  as inputs. The generator is trained to generate the restored matrix  $\bar{\mathbf{S}}$ , while discriminator is trained to judge whether the quality of the missing data recovery is realistic enough. The employed generator and discriminator network structures are shown in Fig. 3 and Fig. 4, respectively. Through the game between the generator and the discriminator, effective data recovery can be achieved.

##### 3) Loss Function of GAN

Based on the description of the model structure, the loss function of the generator and discriminator is proposed.

The loss function of generator includes the adversarial loss and the recovery loss. The adversarial loss is defined based on the output of discriminator, which represents the quality of recovery of missing data, i.e.,



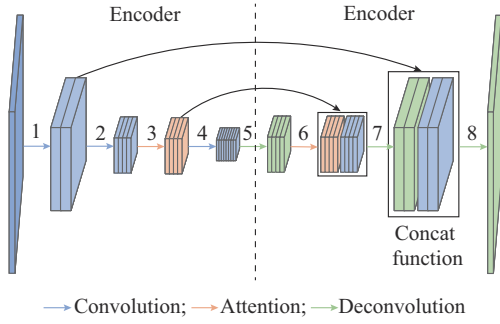


Fig. 3. Generator network structure.

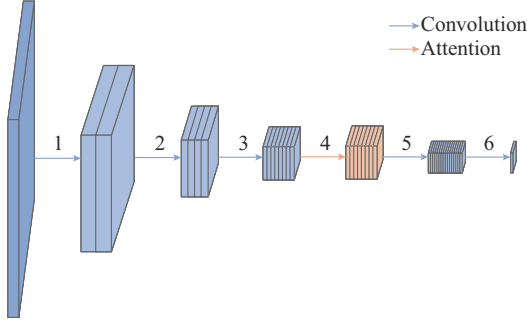


Fig. 4. Discriminator network structure.

$$L_a = -\mathbb{E}_{S,M}[\mathbf{D}(\bar{\mathbf{S}})] \quad (5)$$

where  $\mathbf{D}(\cdot)$  is the discriminant value of the output of discriminator. The recovery loss is defined as the masked root-mean-squared error (RMSE) between  $\bar{\mathbf{S}}$  and  $\mathbf{S}$ . Since  $L_a$  has already dealt with the missing data part, the recovery loss mainly focuses on the part of intact data. The mathematical expression of the recovery loss is given by:

$$L_r = \mathbb{E}_{S,M}[\|\mathbf{S} \odot \mathbf{M} - \bar{\mathbf{S}} \odot \mathbf{M}\|] \quad (6)$$

Next, the loss function of the generator is defined as:

$$L_G = L_a + L_r \quad (7)$$

The objective of the discriminator is to maximize the discriminative value of real historical data and minimize the discriminative value of the output of the generator. Therefore, the loss function of the discriminator is defined as:

$$L_D = -\mathbb{E}_S[\mathbf{D}(\mathbf{S})] + \mathbb{E}_{S,M}[\mathbf{D}(\bar{\mathbf{S}})] \quad (8)$$

### B. DT Modeling Phase

In order to build a virtual model at the DT layer that can accurately reflect the process of PV power generation in the real world, we receive the device parameters and the historical dataset after the data recovery from the transmission layer. First, we construct a digital physical model to simulate the internal mechanism of PV panel power generation. Then, a parallel CNN-BiLSTM model is built and trained to extract the inherent characteristics of meteorological factors and PV power generation. Eventually, a combination formula is applied to connect the two models to form the DT model.

#### 1) Digital Physical Model

This part is composed of the underlying physical model and the power deviation correction module. Specifically, the

PV power plant is a device designed to convert solar radiation into direct current electricity. It primarily consists of solar cells, which are semiconductor thin films that directly generate electricity when exposed to sunlight of a specific irradiance. These solar cells can produce voltage and current when connected in a circuit. The power output of solar cells varies due to fluctuations under weather conditions. Solar radiation plays a crucial role in determining the power output. Higher temperatures can reduce the efficiency of power generation components, while strong winds can help to reduce the temperature of solar cells, thereby increasing power generation. This behavior can be effectively modeled using an equivalent circuit.

The formula for describing the output current of a single diode equivalent circuit is given by:

$$I = I_{pv} - I_{sh} - I_d \quad (9)$$

where  $I_{pv}$  is the photocurrent generated by the battery due to incident solar radiation;  $I_{sh}$  is the short-circuit current caused by leakage at the edge of the battery and the formation of metal bridges; and  $I_d$  is the diode current that comes from the Shockley equation. The mathematical expressions of  $I_{sh}$  and  $I_d$  are given by:

$$I_{sh} = \frac{IR_s + V}{R_{sh}} \quad (10)$$

$$I_d = I_0 \left[ \exp \left( \frac{q(IR_s + V)}{AbT_m} \right) - 1 \right] \quad (11)$$

where  $V$  is the voltage drop across the battery due to incident solar radiation;  $R_s$  is the series resistance;  $R_{sh}$  is the shunt resistance;  $I_0$  is the reverse saturation current;  $q$  is the electron charge;  $A$  is the ideality factor of the diode;  $b$  is the Boltzmann constant; and  $T_m$  is the actual temperature of the PV module defined as:

$$T_m = T + \frac{G}{\mu_0 + \mu_1 v} \quad (12)$$

where  $T$  is the ambient temperature;  $G$  is the real-time irradiance;  $\mu_0$  is the irradiance-induced shading effect;  $\mu_1$  is the effect of wind speed; and  $v$  is the real-time wind speed.

The generated power of the solar cell, denoted as  $P_0$ , is calculated as:

$$P_0 = VI \quad (13)$$

There exist five unknown parameters, i.e.,  $I_{pv}$ ,  $I_0$ ,  $R_s$ ,  $R_{sh}$ , and  $A$ . By establishing five equations based on the short-circuit current  $I_{SC}$ , open-circuit voltage  $U_{OC}$ , the maximum power  $P_{0m} = U_m I_m$ ,  $\frac{dP}{dV} = 0$  at the maximum power point, and  $\frac{dI}{dV} = -\frac{1}{R_{sh}}$  at the short-circuit point, the unknown parameters can be obtained. With those components, a physical model of the PV power station can be constructed. The input data are  $T$ ,  $G$ , and  $v$ , while the output data are  $I$ ,  $V$ , and  $P_0$ .

Based on the predicted PV power data obtained from the aforementioned underlying physical model, the model considers only environmental temperature, real-time irradiance, and real-time wind speed as inputs. However, this method fails to account for the complex practical conditions of the PV power station and other weather factors, leading to certain

deviations in the prediction results. To address this issue, a deviation correction process is introduced. In this process, the similarity in PV power output under the influence of external climate conditions is taken into account, considering different seasons and sampling times within a day. By calculating and storing the difference between the output power of the underlying physical model and the actual historical power, it is possible to determine a correction value. This correction value is then used to adjust the predicted power from the underlying physical model, resulting in more accurate prediction results within the digital physical model. To implement the deviation correction, the historical weather data that have been restored through the use of GAN are employed as input to the underlying physical model. Let  $\mathbf{P}_0 = [P_0(1), P_0(2), \dots, P_0(n_T)]^T$  represent the output power, and  $\mathbf{P}$  is the actual power. The difference between the predicted power and the actual power of the underlying physical model can be calculated as:

$$\boldsymbol{\theta} = \mathbf{P}_0 - \mathbf{P} = [\theta(1), \theta(2), \dots, \theta(n_T)]^T \quad \boldsymbol{\theta} \in \mathbb{R}^{n_T} \quad (14)$$

According to (14), the revised value, denoted as  $\mathbf{E} = [E(1), E(2), \dots, E(n_T)]^T$ , can be calculated as:

$$E(0) = 0 \quad (15)$$

$$E(i) = \frac{\beta E(i-1) + (1-\beta)\theta(i)}{1-\beta^i} \quad (16)$$

where  $\beta$  is an adjustable hyperparameter between 0 and 1.

As the historical dataset used for constructing the digital physical model typically contains a large amount of data, spanning more than one year, it is essential to fully utilize this dataset while ensuring the stability of the revised value calculation. To achieve this, the calculation result of the revised value is averaged on a yearly basis, resulting in  $\bar{\mathbf{E}} = [\bar{E}(1), \bar{E}(2), \dots, \bar{E}(365t)]^T$ . The formula for calculating the elements in the  $\bar{\mathbf{E}}$  array is expressed as:

$$\bar{E}(j) = \begin{cases} \frac{1}{m} \sum_{i=1}^m E(j + (i-1) \times 365t) & 1 \leq j \leq j_0 \\ \frac{1}{m-1} \sum_{i=1}^{m-1} E(j + (i-1) \times 365t) & j_0 < j \leq 365t \end{cases} \quad (17)$$

where  $j_0 = n_T \% (365t)$ ;  $m = \lceil n_T / (365t) \rceil$  and  $\lceil \cdot \rceil$  is used to round up to an integer; and  $t$  is the number of data sampling times per day.

After  $\bar{\mathbf{E}}$  is obtained, the power value of the corrected output power at the  $j^{\text{th}}$  sampling time point, representing the predicted power of the digital physical model, can be calculated as:

$$\tilde{P}_D(j) = P_0(j) - \bar{E}(j) \quad (18)$$

## 2) Parallel CNN-BiLSTM Model

To capture the underlying relationships among diverse meteorological data and the temporal dependencies within the data, we propose a parallel CNN-BiLSTM model, as depicted in Fig. 4. The parallel CNN-BiLSTM network can simultaneously process different parts of the input data and fully leverage the capabilities of parallel computing. This significantly enhances computational efficiency and speeds up both model training and inference processes. Furthermore, the parallel CNN and BiLSTM layers facilitate the extraction and

integration of data features concurrently. This enables us to capture the information pertaining to various aspects of the data and to provide a potent model representation, thereby improving prediction accuracy. To be specific, the CNN component is employed to extract intrinsic features between different data types within a defined time step. Meanwhile, the BiLSTM is utilized to capture deeper temporal features by considering the information from both the “forward” and “backward” directions. The parallel architecture of the CNN and BiLSTM allows independent extraction of intrinsic features from various data types and deeper temporal features from the input data. These features are then concatenated into a final feature vector, which is used for predicting PV power generation.

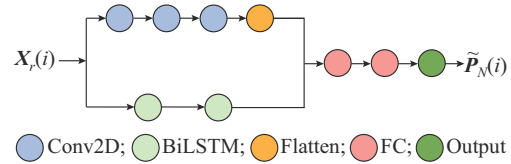


Fig. 5. Structure of parallel CNN-BiLSTM model.

Tensor modeling is conducted on the recovered historical meteorological data and historical power data, denoted as  $\mathbf{X}_r$  and  $\mathbf{P}_r$ , respectively. Meanwhile, the predicted power of the neural network model is defined as  $\tilde{\mathbf{P}}_N$ :

$$\mathbf{X}_r = [\mathbf{X}_r(1), \mathbf{X}_r(2), \dots, \mathbf{X}_r(n_T)] \quad (19)$$

$$\mathbf{P}_r = [\mathbf{P}_r(1), \mathbf{P}_r(2), \dots, \mathbf{P}_r(n_T)] \quad (20)$$

$$\tilde{\mathbf{P}}_N = [\tilde{\mathbf{P}}_N(1), \tilde{\mathbf{P}}_N(2), \dots, \tilde{\mathbf{P}}_N(n_T)] \quad (21)$$

$$\mathbf{X}_r(i) = [\mathbf{x}_r(i-L), \mathbf{x}_r(i-L+1), \dots, \mathbf{x}_r(i-2), \mathbf{x}_r(i-1)]^T \quad (22)$$

where  $\mathbf{X}_r(i) \in \mathbb{R}^{L \times n_s}$  is the input data required to predict the power at the  $i^{\text{th}}$  sampling time point;  $\mathbf{x}_r(i)$  is the meteorological data collected at the  $i^{\text{th}}$  sampling time point;  $\mathbf{P}_r(i)$  is the power at the  $i^{\text{th}}$  sampling time point;  $\tilde{\mathbf{P}}_N(i)$  is the predicted power by the parallel CNN-BiLSTM model; and  $L$  is the time step of the input data of the parallel CNN-BiLSTM model. Note that  $\mathbf{X}_r(i)$  contains the time-series data. After normalization, it is viewed as the grayscale image and served as the input of Conv2D. In addition,  $\mathbf{X}_r(i)$  is flattened and served as the input of BiLSTM. The loss function is set to be:

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (\tilde{\mathbf{P}}_N(i) - \mathbf{P}_r(i))^2} \quad (23)$$

where  $N$  is the batch of samples for each training.

## 3) Combination Formula of PV Power Prediction Results

In order to leverage the advantages of the digital physical model and the parallel CNN-BiLSTM model, we design the combination formula that is a linear combination of the prediction results from the two models. The combined result is used as the final predicted PV power. We define  $\mathbf{P}_D$  and  $\mathbf{P}_N$  to represent the predicted values of the digital physical model and the parallel CNN-BiLSTM model. The difference  $\boldsymbol{\theta}_1$  between the real power and the predicted power from the digital physical model as well as the difference  $\boldsymbol{\theta}_2$  between the real power and the predicted power from the parallel

CNN-BiLSTM model can be calculated as (24) and (25), respectively:

$$\theta_1 = P_D - P = [\theta_1(1), \theta_1(2), \dots, \theta_1(n_T)]^T \quad (24)$$

$$\theta_2 = P_N - P = [\theta_2(1), \theta_2(2), \dots, \theta_2(n_T)]^T \quad (25)$$

In order to reduce the amount of data, maximize the use of recovered historical data, and avoid the contingency of calculation results, the above two difference values are averaged annually. The calculation formula is given by:

$$\bar{\theta}_\kappa(j) = \begin{cases} \frac{1}{m} \sum_{i=1}^m \theta(j + (i-1) \times 365t) & 1 \leq j \leq j_0 \\ \frac{1}{m-1} \sum_{i=1}^{m-1} \theta(j + (i-1) \times 365t) & j_0 < j \leq 365t \end{cases} \quad (26)$$

where  $\kappa = 1$  or  $2$ .

According to (26), we can obtain the averaged differences  $\bar{\theta}_1$  and  $\bar{\theta}_2$  as:

$$\bar{\theta}_1 = [\bar{\theta}_1(1), \bar{\theta}_1(2), \dots, \bar{\theta}_1(365t)]^T \quad (27)$$

$$\bar{\theta}_2 = [\bar{\theta}_2(1), \bar{\theta}_2(2), \dots, \bar{\theta}_2(365t)]^T \quad (28)$$

where  $\bar{\theta}_1(j) \in \bar{\theta}_1$ ; and  $\bar{\theta}_2(j) \in \bar{\theta}_2$ .

The combined formula of the power prediction results is defined as:

$$\hat{P}(j) = w_1(j)P_D(j) + w_2(j)P_N(j) \quad (29)$$

where  $w_1$  and  $w_2$  are the weight coefficients of the predicted power from the digital physical model and the parallel CNN-BiLSTM model, respectively. The mathematical definitions of the weight coefficients are designed as:

$$w_1(j) = \frac{k\bar{\theta}_2^2(j)}{\bar{\theta}_1^2(j) + k\bar{\theta}_2^2(j)} \quad (30)$$

$$w_2(j) = \frac{\bar{\theta}_1^2(j)}{\bar{\theta}_1^2(j) + k\bar{\theta}_2^2(j)} \quad (31)$$

where  $k > 0$  is a hyperparameter.

### C. Power Prediction Phase

After finishing the phases of data preparation and DT modelling, we proceed the final power prediction phase. Taking the real-time weather data as input, we use the digital physical model and the parallel CNN-BiLSTM model to calculate the prediction results  $\hat{P}_D$  and  $\hat{P}_N$ . Then, the final predicted power  $\hat{P} = w_1\hat{P}_D + w_2\hat{P}_N$  is obtained through the calculation of the combined formula of the power prediction results.

Remark 1: the data augmentation method is a kind of data preprocessing technique for expanding training data through a series of transformations and extensions of the original dataset to generate new training samples. Distinguished from the data augmentation methods, the DT focuses on creating the data counterpart of the physical systems to provide simulation and analysis. In the aspect of solving the prediction problem, the data augmentation method enables the extension of training data to deal with the data imbalance and improve the prediction accuracy. In this paper, the DT is used to create digital physical models that reflect the intrinsic mechanisms of physical systems, and use machine learning models to capture hidden features that are difficult to ana-

lyze based on physical models. This enables the integration of physical knowledge and data-driven methods to achieve accurate modeling and prediction of real systems. In this paper, we have complete real dataset without the requirement of generating new dataset. Thus, we intend to introduce DT to increase the prediction accuracy.

## IV. SIMULATIONS

### A. Preparation

#### 1) Dataset

The real dataset comes from the global intelligent evolution simulation experiment platform and engineering demonstration application project of distributed information energy system at Northeastern University in China. This dataset contains historical records of relevant information on power generation and weather conditions. Specifically, it covers the period between 2016 and 2018 and includes the data recorded from 08:00 a.m. to 17:00 p.m. daily. The sampling interval is 15 min. The data types include temperature, wind speed, solar irradiance, relative humidity, and PV output power. The first 24 months and the last 12 months of the historical dataset are taken as training and testing samples, respectively. The time dimension of the data, the number of meteorological factors, and the data sampling frequency per day are  $n_T = 40515$ ,  $n_x = 12$ , and  $t = 37$ , respectively. To handle the missing and abnormal data, invalid data are identified and set to be zero in the mask matrix  $M$ . In order to eliminate data dimensions and enhance data features, the historical dataset is normalized and then inputted into the GAN for data recovery to improve the quality of the dataset.

#### 2) Network Parameters

The parameters of the generator and discriminator networks are listed in Tables I and II, respectively.

TABLE I  
PARAMETERS OF GENERATOR NETWORK

Layer	Part	Kernel size	Number
1	Convolution	$4 \times 4$	16
2	Convolution	$3 \times 3$	32
3	Attention		
4	Convolution	$4 \times 4$	64
5	Deconvolution	$4 \times 4$	64
6	Attention		
7	Deconvolution	$3 \times 3$	64
8	Deconvolution	$4 \times 4$	32

TABLE II  
PARAMETERS OF DISCRIMINATOR NETWORK

Layer	Part	Kernel size	Number
1	Convolution	$3 \times 3$	8
2	Convolution	$3 \times 3$	16
3	Convolution	$5 \times 5$	32
4	Attention		
5	Convolution	$3 \times 3$	64
6	Convolution	$4 \times 4$	1

The generator takes input data with a time step of  $l=92$  and  $n_x+1=13$ , resulting in  $j-i+1=7$  and  $f=1$ . The convolutional layers in the generator network employ SAME padding with a stride of  $s=2$ . Similarly, the convolutional layers in the discriminator network also adopt SAME padding, with a stride of  $s=2$ , except for the last convolutional layer, which has a stride of  $s=1$ . The Adam optimizer is used for the GAN with the activation function Leaky ReLU and keep-probability of 0.8. The input data of the parallel CNN-BiLSTM model in the DT layer have a time step of  $l=12$  and  $n_x=12$ . We chose the batch size as 64, the epochs as 50, and the learning rate as 0.0002. The parameters of the parallel CNN-BiLSTM model are shown in Table III, where the convolutional layer has no padding (i.e.,  $p=0$ ) with a stride of  $s=1$ .

TABLE III  
PARAMETERS OF PARALLEL CNN-BiLSTM MODEL

Part	Kernel size or hidden size	Number of convolutional kernels
Conv2D 1	$4 \times 4$	6
Conv2D 2	$4 \times 4$	6
Conv2D 3	$3 \times 3$	8
BiLSTM 1	64	
BiLSTM 2	64	
FC 1	128	
FC 2	64	

### 3) Performance Evaluation Metrics for Prediction

We evaluate the accuracy of PV power prediction models by using the RMSE and the mean absolute error (MAE), which are defined as:

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (P_i - \hat{P}_i)^2} \quad (32)$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |P_i - \hat{P}_i| \quad (33)$$

where  $P_i$  is the measured PV power at the  $i^{\text{th}}$  sampling time;  $\hat{P}_i$  is the corresponding predicted value; and  $n$  is the total number of samples.

### 4) Determination of Hyperparameters

The hyperparameter  $\beta$  of the deviation correction module in the digital physical model and the hyperparameter  $k$  in the combination formula of power prediction results are determined by using the grid searching method. The decision principle of  $\beta$  and  $k$  is that the higher the accuracy of the predicted power, the better the determination of hyperparameters. It means that the hyperparameter should be determined to minimize the RMSE. The searching results for RMSE of  $\beta$  and  $k$  are shown in Fig. 6 and Fig. 7, respectively. Specifically, it can be observed from Fig. 6 that the optimal value of hyperparameter  $\beta$  is 0.8061. When  $\beta$  is 0, RMSE is 12.360. As  $\beta$  increases to 0.8061, RMSE decreases to 9.669. As  $\beta$  further increases to 1, RMSE increases to 17.291. According to Fig. 7, the optimal value of hyperparameter  $k$  is 0.9809. When  $k$  is 0, RMSE is 5.367. As  $k$  increases to 0.9809, RMSE decreases to 4.293. As  $k$  further increases to 2, RMSE increases to 4.596.

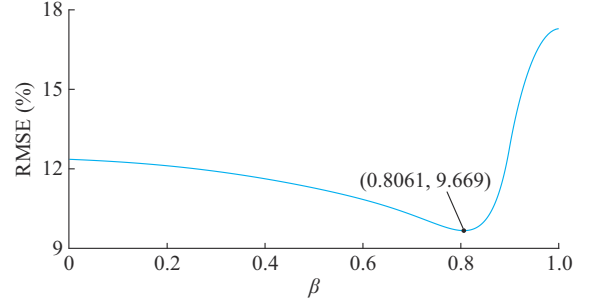


Fig. 6. Searching result for RMSE of  $\beta$ .

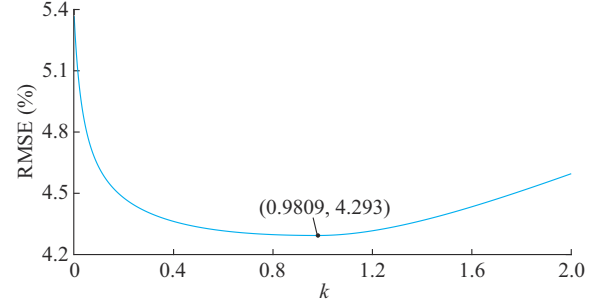


Fig. 7. Searching result for RMSE of  $k$ .

## B. Performance Evaluation and Comparison Analysis

In this case study, we focus on evaluating the performance of the proposed DT-empowered PV power prediction method by comparing with several baselines. The baselines include CNN [11], LSTM [14], CNN-LSTM [18], and GCN [22]. We compare the prediction accuracy of those methods for different weather types (i.e., sunny, rainy, and extreme weather) and different seasons (i.e., spring, summer, autumn, and winter). Meanwhile, typical days are selected as shown in Table IV.

TABLE IV  
TYPICAL DAYS IN 2018

Weather type	Winter	Spring	Summer	Autumn
Sunny	January 9	April 7	July 13	October 27
	January 10	April 8	July 14	October 28
	January 11	April 9	July 15	October 29
	January 12	April 10	July 16	October 30
Rainy	January 2	April 16	July 2	October 14
	January 5	April 29	July 3	October 15
	January 19	May 17	July 4	October 16
	January 21	May 27	July 24	October 21
Extreme	January 3	March 7	June 28	November 5
	January 4	April 5	July 6	November 7
	January 27	April 23	July 22	November 8
	February 19	May 5	August 17	November 26

The results of PV power prediction on typical days using the proposed method and baselines are presented in Fig. 8. The comparison of performance evaluation metrics using the proposed method and baselines, including all prediction results of the testing set, is shown in Table V. In order to com-



prehensively compare the predictive performance of the proposed method, Tables VI and VII list the values of RMSE and MAE after performing the proposed method and baselines in different weather types and seasons, respectively.

Meanwhile, Table VIII presents the values of RMSE and MAE after performing the proposed method and baselines on all testing samples. The following conclusions can be drawn.

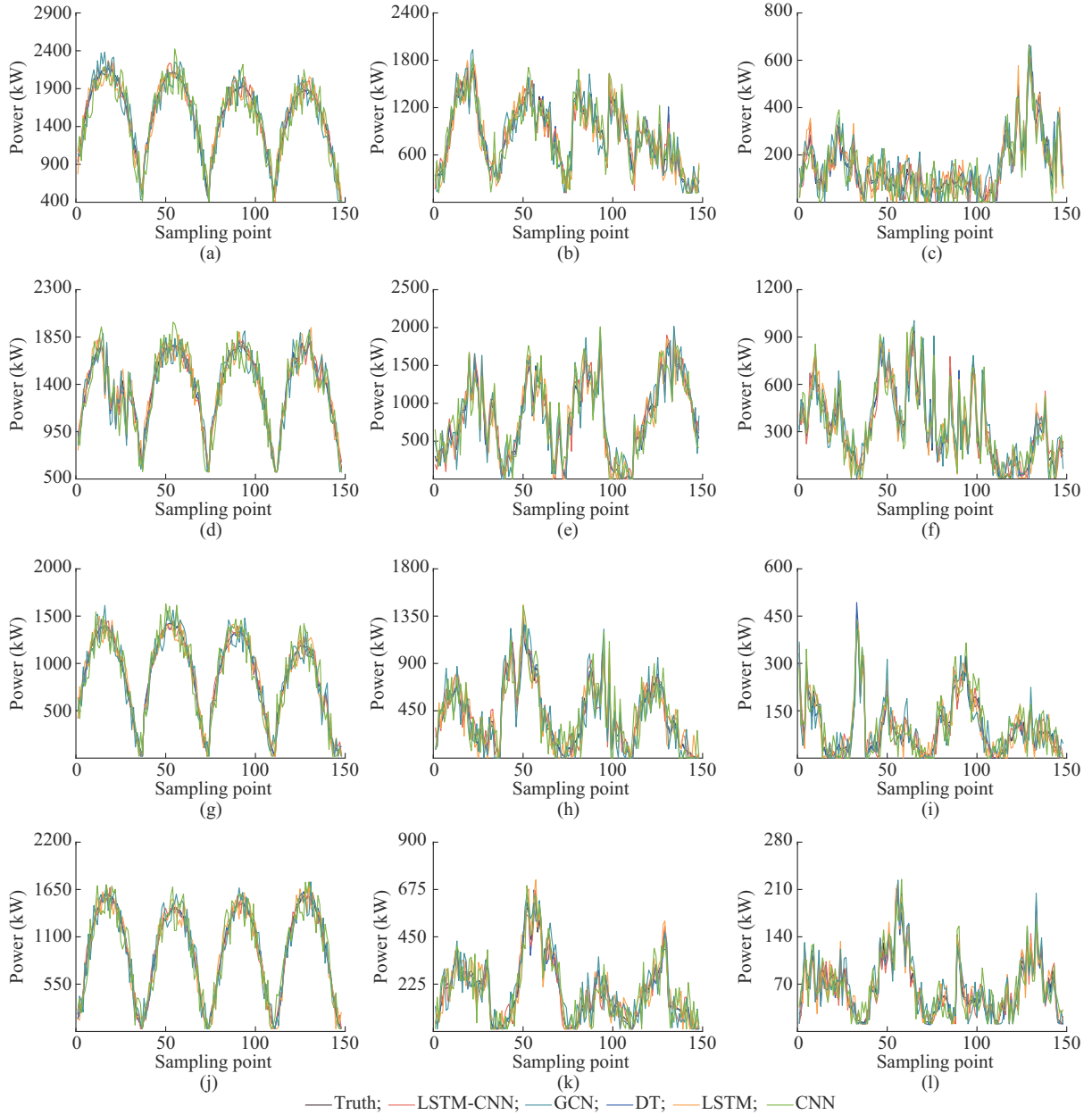


Fig. 8. Results of PV power prediction on typical days using proposed method and baselines. (a) Sunny in spring. (b) Rainy in spring. (c) Extreme weather in spring. (d) Sunny in summer. (e) Rainy in summer. (f) Extreme weather in summer. (g) Sunny in autumn. (h) Rainy in autumn. (i) Extreme weather in autumn. (j) Sunny in winter. (k) Rainy in winter. (l) Extreme weather in winter.

1) The proposed method obtains the lowest RMSE and MAE values compared with the baselines, regardless of the season and weather conditions. The lowest RMSE value means that the prediction performance of the proposed method is the most stable and the error fluctuation range is small. The lowest MAE value denotes that the difference between the predicted results of the proposed method and the actual observed values are the smallest.

2) In comparison to the LSTM and CNN models, the proposed method is capable of extracting spatio-temporal fea-

tures from the dataset more effectively and has stronger abilities in mining data features. Compared with the CNN-LSTM model, the proposed method considers not only the inherent hidden features of weather and power data, but also takes into account the practical conditions of PV panels and other devices. For the GCN, it relies primarily on the adjacency relationships of nodes, which limits information propagation and leads to lower prediction accuracy. Consequently, the prediction accuracy of the proposed method is significantly superior to that of baselines.

TABLE V  
COMPARISONS OF PERFORMANCE EVALUATION METRICS USING PROPOSED METHOD AND BASELINES

Season	Weather type	Evaluation indicator	DT	CNN	LSTM	CNN-LSTM	GCN
Spring	Sunny	RMSE	5.4841	13.8077	9.1976	7.6780	11.3677
		MAE	3.8838	7.3492	5.3634	5.3930	5.2880
	Rainy	RMSE	6.9357	12.8645	10.2409	8.0464	11.7896
		MAE	5.0822	7.9312	5.9123	5.2222	6.2912
	Extreme	RMSE	3.6891	9.3243	7.4335	5.6493	8.7926
		MAE	2.3016	5.3195	4.3060	3.2325	5.0059
Summer	Sunny	RMSE	4.2208	7.0925	6.6126	5.5649	7.7595
		MAE	3.2375	4.3243	4.2353	3.4159	5.1515
	Rainy	RMSE	5.4557	12.8795	10.3968	8.6488	13.4175
		MAE	3.5505	8.0703	5.9831	4.8015	7.7678
	Extreme	RMSE	4.2045	13.2859	9.5298	7.7692	12.3455
		MAE	2.7993	7.4703	5.4042	4.4423	7.0225
Autumn	Sunny	RMSE	3.9254	10.3079	6.8957	6.2815	8.5103
		MAE	2.8185	5.6227	4.4287	4.0632	4.7555
	Rainy	RMSE	3.3121	11.1165	8.2280	6.9399	9.8849
		MAE	2.2533	6.3080	4.7879	4.4101	5.4118
	Extreme	RMSE	2.7799	6.1071	4.8034	3.6499	5.7009
		MAE	1.5666	3.4169	2.6687	2.0296	3.1732
Winter	Sunny	RMSE	4.9257	12.4006	7.9110	6.6768	9.4003
		MAE	3.6179	7.7115	5.3117	4.3161	6.2213
	Rainy	RMSE	3.0111	6.4340	5.3785	4.7058	6.0889
		MAE	1.7395	4.1226	3.3040	2.5032	3.8803
	Extreme	RMSE	3.0069	5.6353	3.7990	3.3466	5.1164
		MAE	1.9092	3.3758	2.6538	2.1242	3.0018

TABLE VI  
VALUES OF RMSE AND MAE AFTER PERFORMING PROPOSED METHOD AND BASELINES IN DIFFERENT WEATHER TYPES

Weather type	Evaluation indicator	DT	CNN	LSTM	CNN-LSTM	GCN
Sunny	RMSE	4.6787	11.7701	7.7518	6.6092	9.3590
	MAE	3.3787	6.2900	4.5375	4.2887	5.2065
Rainy	RMSE	4.7853	11.4384	8.6185	7.0423	10.4900
	MAE	2.7853	6.4790	4.9233	4.1279	5.8891
Extreme	RMSE	3.3973	8.7974	6.5928	5.2062	8.2042
	MAE	2.1973	5.3658	3.8859	3.1043	4.7576

TABLE VII  
VALUES OF RMSE AND MAE AFTER PERFORMING PROPOSED METHOD AND BASELINES IN DIFFERENT SEASONS

Season	Evaluation indicator	DT	CNN	LSTM	CNN-LSTM	GCN
Spring	RMSE	5.2201	11.2352	8.2694	6.1115	9.8017
	MAE	3.4201	7.1627	4.7939	4.1282	5.6451
Summer	RMSE	4.5776	10.9193	8.6484	7.4751	9.6631
	MAE	3.1776	6.6673	5.3002	4.1911	6.7146
Autumn	RMSE	3.1688	8.9542	6.3748	4.9967	7.7242
	MAE	2.1688	5.3829	3.5284	3.2643	4.4611
Winter	RMSE	3.6216	7.7566	5.1221	4.9871	6.6926
	MAE	2.4216	4.9165	3.7249	2.8665	4.1703

TABLE VIII  
COMPARISONS OF VALUES OF RMSE AND MAE AFTER PERFORMING PROPOSED METHOD AND BASELINES ON ALL TESTING SAMPLES

Method	RMSE	MAE
DT	4.2934	2.7841
CNN	9.8195	6.2591
LSTM	6.9598	4.4019
CNN-LSTM	5.8476	3.7675
GCN	9.1588	5.3338

### C. Ablation Analysis

In order to further demonstrate the effectiveness of the proposed method, ablation analysis is conducted in this case study. Figure 9 shows the prediction results for three different weather types under ablation analysis, where one typical day is selected for each weather type in four seasons.

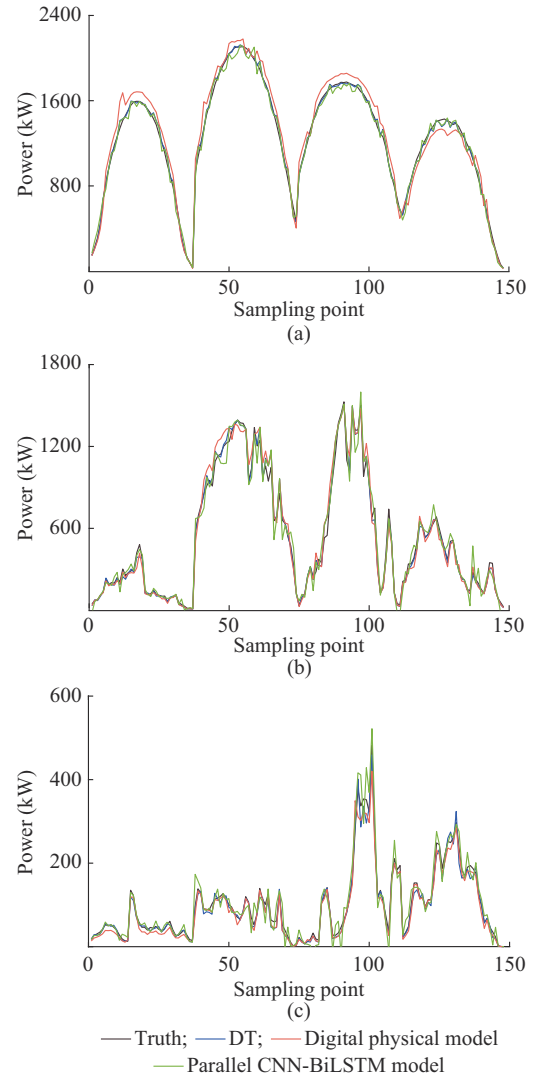


Fig. 9. Prediction results for three different weather types under ablation analysis. (a) Sunny. (b) Rainy. (c) Extreme.

Tables IX and X show the ablation analysis results in different weather types and seasons, respectively, with the digital physical model and the parallel CNN-BiLSTM model.

And Table XI shows the ablation analysis results in all testing samples.

TABLE IX  
ABLATION ANALYSIS RESULTS IN DIFFERENT WEATHER TYPES

Weather type	Evaluation indicator	DT	Digital physical model	Parallel CNN-BiLSTM model
Sunny	RMSE	4.6787	14.1536	6.1925
	MAE	3.3787	8.7514	3.2086
Rainy	RMSE	4.7853	7.8433	6.6547
	MAE	2.7853	4.8217	3.8305
Extreme	RMSE	3.3973	3.7055	3.2356
	MAE	2.1973	2.4293	2.3866

TABLE X  
ABLATION ANALYSIS RESULTS IN DIFFERENT SEASONS

Season	Evaluation indicator	DT	Digital physical model	Parallel CNN-BiLSTM model
Spring	RMSE	5.2201	11.3062	5.3032
	MAE	3.4201	7.2975	3.2701
Summer	RMSE	4.5776	9.7968	6.1744
	MAE	3.1776	5.8107	3.8261
Autumn	RMSE	3.1688	6.2832	3.9064
	MAE	2.1688	4.2208	2.2982
Winter	RMSE	3.6216	9.5294	4.9965
	MAE	2.4216	5.0249	2.7553

TABLE XI  
ABLATION ANALYSIS RESULTS IN ALL TESTING SAMPLES

Evaluation indicator	DT	Digital physical model	Parallel CNN-BiLSTM model
RMSE	4.2934	9.6687	5.3674
MAE	2.7841	5.5808	3.2338

The results indicate that the combined version achieves the highest prediction accuracy compared with the digital physical model and the parallel CNN-BiLSTM model. This is because the proposed method takes advantages of both the physical characteristics of PV power station and the inherent data features between meteorological and power data. This method enables better simulation of real-world PV power generation processes and achieves accurate PV power prediction.

## V. CONCLUSION

In the paper, we have established a DT-empowered PV power prediction framework to achieve reliable data transmission and power prediction with high accuracy. We have designed the use of GAN for data recovery from historical data, which is capable of significantly improving the quality of constructing a DT virtual power station. This enhances the reliability of mapping from the physical space to the digital space. We have proposed a new DT-empowered PV power prediction method. By integrating the digital physical model and the parallel CNN-BiLSTM model, the proposed

method effectively enhances the prediction accuracy for PV power generation. Finally, the testing results on the real dataset from Northeastern University show that the proposed method can achieve higher prediction accuracy than the baselines in different scenarios. In the future work, we would like to investigate the integration of federated learning to enhance the privacy of the proposed method.

## REFERENCES

- [1] A. Shafi, H. Sharadga, and S. Hajimirza, "Design of optimal power point tracking controller using forecasted photovoltaic power and demand," *IEEE Transactions on Sustainable Energy*, vol. 11, no. 3, pp. 1820-1828, Jul. 2020.
- [2] H. Zhang, Y. Li, D. W. Gao *et al.*, "Distributed optimal energy management for energy internet," *IEEE Transactions on Industrial Informatics*, vol. 13, no. 6, pp. 3081-3097, Dec. 2017.
- [3] Y. Li, H. Zhang, X. Liang *et al.*, "Event-triggered-based distributed cooperative energy management for multienergy systems," *IEEE Transactions on Industrial Informatics*, vol. 15, no. 4, pp. 2008-2022, Apr. 2019.
- [4] X. Zhang, Y. Li, S. Lu *et al.*, "A solar time based analog ensemble method for regional solar power forecasting," *IEEE Transactions on Sustainable Energy*, vol. 10, no. 1, pp. 268-279, Jan. 2019.
- [5] K. Hu, S. Cao, L. Wang *et al.*, "A new ultra-short-term photovoltaic power prediction model based on ground-based cloud images," *Journal of Cleaner Production*, vol. 200, pp. 731-745, Nov. 2018.
- [6] B. Kim, D. Suh, M.-O. Otto *et al.*, "A novel hybrid spatio-temporal forecasting of multisite solar photovoltaic generation," *Remote Sensing*, vol. 13, no. 13, pp. 2605, Jul. 2021.
- [7] K. Doubleday, S. Jascourt, W. Kleiber *et al.*, "Probabilistic solar power forecasting using bayesian model averaging," *IEEE Transactions on Sustainable Energy*, vol. 12, no. 1, pp. 325-337, Jan. 2021.
- [8] L. F. Tratar and E. Strmcnik, "The comparison of holt-winters method and multiple regression method: a case study," *Energy*, vol. 109, pp. 266-276, Aug. 2016.
- [9] H. T. C. Pedro and C. F. M. Coimbra, "Assessment of forecasting techniques for solar power production with no exogenous inputs," *Solar Energy*, vol. 86, no. 7, pp. 2017-2028, Jul. 2012.
- [10] Y. Tang, K. Yang, S. Zhang *et al.*, "Photovoltaic power forecasting: a hybrid deep learning model incorporating transfer learning strategy," *Renewable and Sustainable Energy Reviews*, vol. 162, no. 11, p. 112473, Jul. 2022.
- [11] J. Yan, L. Hu, Z. Zhen *et al.*, "Frequency-domain decomposition and deep learning based solar PV power ultra-short-term forecasting model," *IEEE Transactions on Industry Applications*, vol. 57, no. 4, pp. 3282-3295, Jul.-Aug. 2021.
- [12] H. Li, Z. Ren, Y. Xu *et al.*, "A multi-data driven hybrid learning method for weekly photovoltaic power scenario forecast," *IEEE Transactions on Sustainable Energy*, vol. 13, no. 1, pp. 91-100, Jan. 2022.
- [13] F. Wang, J. Li, Z. Zhen *et al.*, "Cloud feature extraction and fluctuation pattern recognition based ultrashort-term regional PV power forecasting," *IEEE Transactions on Industry Applications*, vol. 58, no. 5, pp. 6752-6767, Sept.-Oct. 2022.
- [14] Y. Zhang, C. Qin, A. K. Srivastava *et al.*, "Data-driven day-ahead PV estimation using autoencoder-LSTM and persistence model," *IEEE Transactions on Industry Applications*, vol. 56, no. 6, pp. 7185-7192, Nov.-Dec. 2020.
- [15] L. Cheng, H. Zang, Z. Wei *et al.*, "Short-term solar power prediction learning directly from satellite images with regions of interest," *IEEE Transactions on Sustainable Energy*, vol. 13, no. 1, pp. 629-639, Jan. 2022.
- [16] R. Zhang, H. Ma, T. K. Saha *et al.*, "Photovoltaic nowcasting with bi-level spatio-temporal analysis incorporating sky images," *IEEE Transactions on Sustainable Energy*, vol. 12, no. 3, pp. 1766-1776, Jul. 2021.
- [17] J. Li, C. Zhang, and B. Sun, "Two-stage hybrid deep learning with strong adaptability for detailed day-ahead photovoltaic power forecasting," *IEEE Transactions on Sustainable Energy*, vol. 14, no. 1, pp. 193-205, Jan. 2023.
- [18] S. Chai, Z. Xu, Y. Jia *et al.*, "A robust spatiotemporal forecasting framework for photovoltaic generation," *IEEE Transactions on Smart Grid*, vol. 11, no. 6, pp. 5370-5382, Nov. 2020.
- [19] J. Simeunovic, B. Schubnel, P.-J. Alet *et al.*, "Spatio-temporal graph

- neural networks for multi-site PV power forecasting,” *IEEE Transactions on Sustainable Energy*, vol. 13, no. 2, pp. 1210-1220, Apr. 2022.
- [20] L. Cheng, H. Zang, Z. Wei *et al.*, “Solar power prediction based on satellite measurements: a graphical learning method for tracking cloud motion,” *IEEE Transactions on Power Systems*, vol. 37, no. 3, pp. 2335-2345, May 2022.
- [21] M. Zhang, Z. Zhen, N. Liu *et al.*, “Optimal graph structure based short-term solar PV power forecasting method considering surrounding spatio-temporal correlations,” *IEEE Transactions on Industry Applications*, vol. 59, no. 1, pp. 345-357, Jan.-Feb. 2023.
- [22] T. Yao, J. Wang, Y. Wang *et al.*, “Very short-term forecasting of distributed PV power using GSTANN,” *CSEE Journal of Power and Energy Systems*, doi: 10.17755/CSEEJPES.2022.00110
- [23] L. Cheng, H. Zang, T. Ding *et al.*, “Multi-meteorological-factor-based graph modeling for photovoltaic power forecasting,” *IEEE Transactions on Sustainable Energy*, vol. 12, no. 3, pp. 1593-1603, Jul. 2021.
- [24] J. Choi, J.-I. Lee, I.-W. Lee *et al.*, “Robust PV-BESS scheduling for a grid with incentive for forecast accuracy,” *IEEE Transactions on Sustainable Energy*, vol. 13, no. 1, pp. 567-578, Jan. 2022.
- [25] Q. Li, Y. Xu, B. S. H. Chew *et al.*, “An integrated missing-data tolerant model for probabilistic PV power generation forecasting,” *IEEE Transactions on Power Systems*, vol. 37, no. 6, pp. 4447-4459, Nov. 2022.
- [26] W. Liu, C. Ren, and Y. Xu, “PV generation forecasting with missing input data: a super-resolution perception approach,” *IEEE Transactions on Sustainable Energy*, vol. 12, no. 2, pp. 1493-1496, Apr. 2021.
- [27] Y. Nie, A. S. Zamzam, and A. Brandt, “Resampling and data augmentation for short-term PV output prediction based on an imbalanced sky images dataset using convolutional neural networks,” *Solar Energy*, vol. 224, pp. 341-354, Aug. 2021.
- [28] T. Polasek and M. Cadik, “Predicting photovoltaic power production using high-uncertainty weather forecasts,” *Applied Energy*, vol. 339, p. 120989, Jun. 2023.
- [29] S. Goudarzi, A. Asif, and H. Rivaz, “Fast multi-focus ultrasound image recovery using generative adversarial networks,” *IEEE Transactions on Computational Imaging*, vol. 6, pp. 1272-1284, Aug. 2020.
- [30] L. Han, K. Zheng, L. Zhao *et al.*, “Content-aware traffic data completion in ITS based on generative adversarial nets,” *IEEE Transactions on Vehicular Technology*, vol. 69, no. 10, pp. 11950-11962, Oct. 2020.
- [31] Y. Li and Y. Zhang, “Digital twin for industrial internet,” *Fundamental Research*, vol. 4, no. 1, pp. 21-24, Jan. 2024.
- [32] A. Marot, A. Kelly, M. Naglic *et al.*, “Perspectives on future power system control centers for energy transition,” *Journal of Modern Power Systems and Clean Energy*, vol. 10, no. 2, pp. 328-344, Mar. 2022.
- [33] S. Mihai, M. Yaqoob, D. V. Hung *et al.*, “Digital twins: a survey on enabling technologies, challenges, trends and future prospects,” *IEEE Communications Surveys and Tutorials*, vol. 24, no. 4, pp. 2255-2291, Sept. 2022.
- [34] Y. Wu, K. Zhang, and Y. Zhang, “Digital twin networks: a survey,” *IEEE Internet of Things Journal*, vol. 8, no. 18, pp. 13789-13804, Sept. 2021.
- [35] C. Wang and Y. Li, “Digital-twin-aided product design framework for IoT platforms,” *IEEE Internet of Things Journal*, vol. 9, no. 12, pp. 9290-9300, Jun. 2022.
- [36] H. Elayan, M. Aloqaily, M. Guizani *et al.*, “Digital twin for intelligent context-aware IoT healthcare systems,” *IEEE Internet of Things Journal*, vol. 8, no. 23, pp. 16749-16757, Dec. 2021.
- [37] D. Liu, Y. Du, W. Chai *et al.*, “Digital twin and data-driven quality prediction of complex die-casting manufacturing,” *IEEE Transactions on Industrial Informatics*, vol. 18, no. 11, pp. 8119-8128, Nov. 2022.
- [38] H. Xu, A. Berres, S. B. Yoganath *et al.*, “Smart mobility in the cloud: enabling real-time situational awareness and cyber-physical control through a digital twin for traffic,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 24, no. 3, pp. 3145-3156, Mar. 2023.
- [39] L. Cascone, M. Nappi, F. Narducci *et al.*, “DTPAAL: digital twinning pepper and ambient assisted living,” *IEEE Transactions on Industrial Informatics*, vol. 18, no. 2, pp. 1397-1404, Feb. 2022.

**Xiaoyu Zhang** received the B.S. degree in automation from North China Electric Power University, Beijing, China, in 2018. She is currently a Graduate Student in electrical engineering at Northeastern University, Shenyang, China. Her main research interests include machine learning, digital twin, and data prediction.

**Yushuai Li** received the B.S. degree in electrical engineering and automation, and the Ph.D. degree in control theory and control engineering from Northeastern University, Shenyang, China, in 2014 and 2019, respectively. He is currently a Marie Curie Researcher at the Department of Informatics, University of Oslo, Norway. He received the Best Paper Awards from Journal of Modern Power Systems and Clean Energy and 2023 International Conference on Cyber-energy Systems and Intelligent Energies (ICCSIE). His main research interests include distributed optimization and control, machine learning, digital twin, and their applications in integrated energy and transportation systems.

**Tianyi Li** received the Ph.D. degree from Aalborg University, Aalborg, Denmark, in 2022. She is currently an Assistant Professor with the Department of Computer Science, Aalborg University. Her research interests include spatio-temporal data management and analytics, knowledge integration, graph neural network and digital energy.

**Yonghao Gui** received the B.S. degree in automation from Northeastern University, Shenyang, China, in 2009, and the M.S. and Ph.D. degrees in electrical engineering from Hanyang University, Seoul, South Korea, in 2012 and 2017, respectively. From 2017 to 2022, he was with Aalborg University, Aalborg, Denmark, as a Postdoctoral Researcher and an Assistant Professor. In 2022, he was with Vestas, Aarhus, Denmark. He has been working with Oak Ridge National Laboratory, Oak Ridge, USA. His research interests include control of power electronics in power systems, renewable energy integration, and smart grids.

**Qiuye Sun** received the Ph.D. degree in 2007. He is currently a Full Professor with Northeastern University, Shenyang, China, and obtained Special Government Allowances from the State Council in China. His current research interests include optimization analysis technology of power distribution network, network control of Energy Internet, integrated energy systems and microgrids.

**David Wenzhong Gao** received the M.S. and Ph.D. degrees in electrical and computer engineering, specializing in electric power engineering, from Georgia Institute of Technology, Atlanta, USA, in 1999 and 2002, respectively. He is now with the Department of Electrical and Computer Engineering, University of Denver, Denver, USA. He is an Associate Editor for IEEE Journal of Emerging and Selected Topics in Power Electronics, and Journal of Modern Power Systems and Clean Energy. He was an Editor of IEEE Transactions on Sustainable Energy. His research interests include renewable energy and distributed generation, microgrid, smart grid, power system protection, power electronics applications in power systems, power system modeling and simulation, and hybrid electric propulsion systems.