# Scalable Coordinated Charging of Large EV Fleet via Imitation Learning
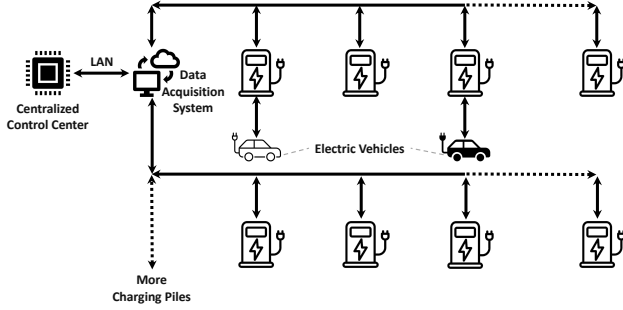


Fig. 1: Illustration of the proposed system

*Abstract*—

*Index Terms*—

## I. INTRODUCTION

COORDINATED charging of Large EV Fleet via Imitation Learning.

## II. SYSTEM MODELING AND PROBLEM FORMULATION

As depicted in Fig. 1, the proposed coordinated charging system is structured around a large parking lot, which is equipped with many electric vehicle (EV) charging piles. Once the connection between the EV and the charging pile is established, a wired data communication link is formed [1], [2]. The EV can transmit information such as rated charging power, current state of charge (SoC), and battery capacity to the charging pile, while the charging pile will transmit the current charging power to the corresponding EV. In addition to the charging piles, the parking lot incorporates a data acquisition system (DAS), responsible for the regular collection of pertinent charging data. This data is then transmitted to the centralized control centre (CCC) via a local area network (LAN) [3]–[5]. The CCC employs the proposed coordinated charging strategy based on imitation learning and reinforced with safety post-processing measures to swiftly determine optimal charging actions for each charging pile. These decisions are communicated back to the DAS, prompting adjustments in power output at each charging pile. This continuous real-time loop ensures the efficiency of the charging operations.

### A. System Modeling

Consider a parking lot with $N$ charging piles, each indexed by $k \in \{1, 2, ..., N\}$. Since EVs continuously arrive at and depart from the parking lot, a discrete-time model is usually utilized to describe the dynamic behavior of the system to facilitate receding horizon optimization (RHC) [6]–[9]. Let each time slot occupy time of length $\Delta t$ and the total number of time slots across the scheduling horizon be $T$. For convenience, the first time step of a scheduling horizon is always recorded as $t = 1$, which is particularly useful when integrating our algorithm into a receding horizon optimization framework [9], [10]. That is, we adopt a relative time system specific to each scheduling period in this study, and all other timestamps such as the arrival and departure time of an EV are shifted implicitly according to the relative time origin. After one time step passes, we start a new scheduling horizon, and the optimization for coordinated charging is redone in a RHC framework.

Consider an arbitrary scheduling horizon Let us consider hereafter an arbitrary index $k \in [1, N]$ and an arbitrary time step $t \in [1, T]$ without loss of generality. After an EV establishes a connection to charging pile $k$, it is denoted as $EV_k$. Note that $EV_k$ and charging pile $k$ share identical parameters upon the connection. Each $EV_k$ is characterized by 7 parameters $\langle t_k^a, t_k^d, E_k, p_k^{rated}, e_k, SoC_k^a, SoC_k^d \rangle$, corresponding to $EV_k$'s arrival time, departure time, battery capacity, rated charging power, charging efficiency, SoC at arrival, and expected SoC upon departure. In practice, $t_k^d$ and $SoC_k^d$ are provided by the EV driver, e.g., through a mobile phone application [7], [11], and the other parameters are typically fetched automatically through communication with the EV's battery management system once the EV is plugged in [4], [9], [12]. In our convention, $EV_k$ is qualified for charging at $t_k^a$ but cannot get charged at the departure time step $t_k^d$. Besides, it is assumed that $EV_k$ does not leave earlier before its provided departure time $t_k^d$.

Let $p_{t,k}$ denote the charging power of $EV_k$ at time $t$. The rated power of $EV_k$ is $p_k^{rated}$, commonly treated as the maximum charging power. We scale $p_{t,k}$ into the unit range $[0, 1]$ to facilitate subsequent algorithm design [6], which is termed the charging action $a_{t,k}$ and computed by

$$a_{t,k} = p_{t,k}/p_k^{rated} \tag{1}$$

Evidently, $EV_k$ can only be charged when it is connected to charging pile $k$. The range of $a_{t,k}$ is thus

$$\begin{cases} a_{t,k} \in [0,1], & t_k^a \leq t < t_k^d \\ a_{t,k} = 0, & \text{otherwise.} \end{cases} \tag{2}$$

To assess the charging cost, $m_{t,k}$ is defined as the electricity cost incurred by $EV_k$ at time $t$:

$$m_{t,k} = \rho_t p_{t,k} \Delta_t = \rho_t a_{t,k} p_k^{rated} \Delta_t \tag{3}$$

where $\rho_t$ represents the real-time electricity price at time $t$. One goal of coordinated charging is to reduce the total electricity bill $C$ computed by

$$C = \sum_{k=1}^{N} \sum_{t=1}^{T} m_{t,k}. \tag{4}$$

Apart from reducing the charging cost, the equality of

charging among customers is also worth considering. Roughly speaking, we want the charging of each EV to progress at a similar pace instead of charging one EV first for some intervals and then another. The SoC of $EV_k$ at time $t$ is termed $SoC_{t,k} \in [0,1]$. If the EV keeps charging at the rated power, then the minium number of time steps still required to fulfill the expected $SoC_k^d$ at departure is given by

$$r_{t,k} = \frac{E_k(SoC_k^d - SoC_{t,k})}{e_k p_k^{\text{rated}} \Delta t}. \qquad (5)$$

Note here that we allow $r_{t,k}$ to be fractional, because its main purpose is to compute the charging urgency of $EV_k$, which is denoted by $g_{t,k}$ at time $t$ and calculated as

$$g_{t,k} = \begin{cases} \frac{r_{t,k}}{t_k^d - t} & t_k^a \le t < t_k^d \\ 0, & \text{otherwise.} \end{cases} \qquad (6)$$

In Eq. (6) above, $(t_k^d - t)$ counts the remaining number of time steps within the expected during of stay of $EV_k$. Intuitively, the charging task becomes increasingly urgent as $g_{t,k}$ increases. A reasonable objective is to maintain the balancing of urgency among the EVs being served. That is, the charging pace of each EV should be similar such that the charging fairness is assured [5], [11].

A straightforward rule to maintain charging fairness is that, if an EV's charging urgency is higher, then we allocate it more power to charge provided that all constraints are satisfied. Formally, the proposed method will try to proportionally allocate charging action $a_{t,k} \in [0,1]$ according to the corresponding charging urgency $g_{t,k}$ for each $EV_k$ at time $t$ [13]. Let $G_{t,k}$ represent the proportion of $EV_k$'s charging urgency $g_{t,k}$ over the total urgency of all EVs:

$$G_{t,k} = \frac{g_{t,k}}{\sum_{k=1}^N g_{t,k}} \qquad (7)$$

Clearly, the range of the normalized charging urgency above is $G_{t,k} \in [0,1]$. Given a candidate set of charging actions $a_{t,k}, k \in [1, N]$ at time $t$, we compute the *fair* charging action of each EV according to proportion $G_{t,k}$ as follows:

$$a_{t,k}^{\text{fair}} = G_{t,k} \sum_{k=1}^N a_{t,k}. \qquad (8)$$

Ideally, a fair charging scheme should attain $a_{t,k} = a_{t,k}^{\text{fair}}$ for every $t$ and $k$, but the strict equality requirement is practically impossible due to a variety of constraints for coordinated charging. In addition, another objective of our charging strategy is to minimize the total electricity cost, which may conflict with the charging fairness. Thus, we try to minimize the difference between $a_{t,k}$ and $a_{t,k}^{\text{fair}}$ instead, which is quantified for the EV fleet by

$$D = \sum_{k=1}^N \sum_{t=1}^T (a_{t,k} - a_{t,k}^{\text{fair}})^2. \qquad (9)$$

### B. Optimization Problem Formulation

Coordinated charging of EVs is typically formulated as a constrained optimization problem. In this study, the objective is to minimize the charging cost An optimal charging scheduling is proposed to minimize the electricity cost $C$ while maximizing the charging fairness for all connected EVs (or equivalently, minimizing the fairness discrimination $D$). Since the two objectives are conflicting in general, we combine them into a single one via weighted sum as follows

$$F = (1 - \beta)C + \beta D, \qquad (10)$$

where the weight coefficient $\beta \in [0,1]$ specifies the importance of each objective, whose value depends on the practical preference of electricity cost or charging fairness [11], [13]. Obviously, as $\beta$ increases, the model prioritizes fairness over electricity cost.

There are primarily two constraints associated with the optimization problem. The first is the satisfaction of the EV's SoC demand, that is, to guarantee that the SoC of $EV_k$ at departure time matches the expected SoC $SoC_k^d$ as specified previously by the driver when plugging in the EV. To formulate this constraint, we first describe the SoC dynamics during the stay period of an EV. Since the SoC keeps increasing even within a single time slot if it gets charged, it should be clarified further that $SoC_{t,k}$ refers to the SoC of $EV_k$ at the very beginning of the time step $t$. We then have

$$SoC_{t,k} = \begin{cases} SoC_k^a, & t = t_k^a \\ SoC_{t-1,k} + \frac{p_{t-1,k} e_k \Delta_t}{C_k}, & t_k^a < t \le t_k^d. \end{cases} \qquad (11)$$

The SoC demand constraint is thus enforced by

$$SoC_{t_k^d, k} = SoC_k^d. \qquad (12)$$

Additionally, the sum of all $EV_k$'s charging power $p_{t,k}$ cannot exceed the parking station's power capacity $P_{\text{max}}$:

$$\sum_{k=1}^N p_{t,k} \le P_{\text{max}} \qquad (13)$$

Now we can present the complete optimization problem called **P1** for coordinated charging towards both cost reduction and fairness improvement as follows:

$$(\textbf{P1}) \quad \min_{a} F \qquad (14a)$$

$$\text{subject to } (1) - (13). \qquad (14b)$$

The decision variables of **P1** are collected into a vector as

$$a = [a_{t,k}], \ \forall t \in [1, T], k \in [1, N]. \qquad (15)$$

Since the fairness related term in the objective function (14a) involves fractions and squares, **P1** forms a general nonlinear program (and more precisely, not a quadratic program) that is notoriously difficult to solve because of its excessive computation burden [10], [11], [13], [14]. Though some techniques such as piecewise linear approximation [14], [15], greedy search [9], and metaheuristic algorithms like particle swarm optimization (PSO) [11] have been proposed to handle such nonlinear programming (NLP) problems, they are not applicable here because of either the dependence on special problem properties [9], [10] or the poor scalability to a large number of EVs [11]. For instance, we have conducted a small-scale case study with 100 EVs, and two well-established NLP solvers, namely Bonmin and Couenne[1], both failed to solve **P1** to optimality even after 24 hours. However, the each time step for practical charging scheduling is usually of length $\Delta_t = 10$ min or 15 min [5], [9], [10], and the solving time is expected to be only a fraction of $\Delta_t$ to yield charging actions

---

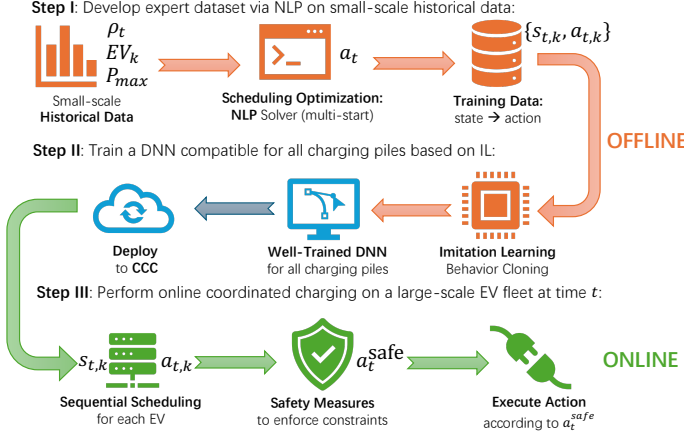[1] Please check the COIN-OR Foundation: https://github.com/coin-or.

Fig. 2: Overview of the proposed framework

fast enough for receding horizon optimization. Obviously, this high-demanding requirement is far beyond the capacity of traditional numerical optimization algorithms like branch and bound [14] or metaheuristic search algorithms like PSO [11].

To sum up, the fundamental challenge of coordinated charging considering both cost and fairness is how to solve **P1** accurately and efficiently such that the optimization strategy is practical and performant even in the case of a large EV fleet, which is the main motivation of the current study. In the next section, we will present a novel solution based on imitation learning to overcome this challenge effectively.

## III. PROPOSED METHOD

Solving large-scale **CO** requires significant computational resources, making it impractical for real-time decision-making scenarios. Therefore, we propose a novel approach based on imitation learning that effectively reduces the computational costs, enabling the real-time provision of near-optimal solutions even for ultra-large-scale EV fleet ($N \geq 500$).

### A. Framework Overview

The framework of the proposed method is illustrated in Fig. 2 for a more concise exposition. It can be delineated into three principal phases. The first stage is expert data generation. An NLP solver is utilized to obtain an optimal solution of the historical data by solving **CO**. Upon getting the optimal solution, a re-simulation of the charging process is conducted. At each time step $t$, the corresponding state and action of $EV_k$ are recorded according to the sequence of charging piles, thus forming an expert training dataset. For the second stage, the dataset is used for training an imitation learning (IL) agent, which is appproximated by a Deep Neural Network (DNN). This single IL agent will be responsible for making decisions for all charging piles. Phase one and two of the proposed method run completely offline. Finally, the well-trained DNN is deployed online within a large-scale parking lot equipped with charging facilities. Like the simulation in step one, the IL agent sequentially estimates the optimal charging action according to the real-time state of each $EV_k$. Note that for each time step $t$, after the IL agent completes the sequential decision-making for all EVs, the CCC will employ safety

measures to adjust the decisions, ensuring all hard constraints about charging are met. Through these steps, the proposed method efficiently handles large-scale coordinated charging problems, generating near-optimal solutions while ensuring computational feasibility and practical applicability.

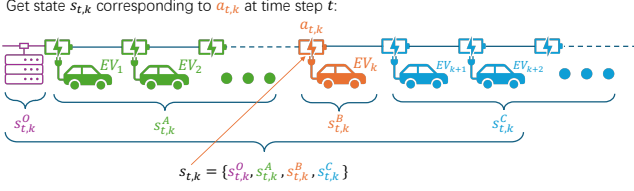### B. Expert Dataset Generation

One significant characteristic of the proposed approach is the excellent scalability, allowing for the flexible expansion of both the scale of EV fleets and sampling interval $\Delta t$. This feature endows the proposed method with the full advantage of optimization-based methods. To construct an expert training dataset, several small-scale EV charging scenarios will be generated. By learning from these small-scale scenarios, the proposed method can accommodate EV charging scheduling problems of any scale.

Given a charging scenario $\chi$, the optimal scheduling $A$ is obtained by solving **CO** on scenario $\chi$. We need to extract the corresponding state $s_{t,k}$ for each action $a_{t,k}$ to form the expert dataset, which will be used for behaviour cloning later in imitation learning. The proposed method's principal attribute lies in its scalability. When formulating the state $s_{t,k}$ for $EV_k$, it is imperative to ensure its adaptability to EV fleets of varying magnitudes. In the context of an EV fleet comprising $N$ vehicles, partitioning it into three distinct groups as follows:

1) $Group_A$ encompasses EVs connected to $pile_1$ to $pile_{k-1}$, representing vehicles whose charging actions have already been determined.
2) $Group_B$ constitutes a specialized subset containing solely EV connected to $pile_k$, for which the charging action $a_{t,k}$ is presently under determination.
3) $Group_C$ encompasses EVs connected to $pile_1$ to $pile_{k-1}$, representing vehicles whose charging actions are pending decisions.

Consequently, irrespective of the EV fleet's scale, obtaining the state $s_{t,k}$ for $EV_k$ entails abstracting it into three distinct groups. As illustrated in Figure 3, $s_{t,k}$ is composed of $s_{t,k}^O$, $s_{t,k}^A$, $s_{t,k}^B$ and $s_{t,k}^C$.

$s_{t,k}^A$, $s_{t,k}^B$ and $s_{t,k}^C$ represent states for the $Group_A$, $Group_B$ and $Group_C$ mentioned above, respectively. These groups can be regarded as three 'large' EVs, and they share some same attributes and have some unique attributes. Consider the same attributes first, given a specific $Group_X$ where consists of $N_X$ piles, the state $s_{t,k}^X$ always contains average battery capacity $C_X$, average rated charging power $p_X^{rated}$, weighted average expected state of charge $SoC_X^d$ when leaving, weighted average state of charge $SoC_{t,X}$, and average remain parking time $r_{t,X}$.

Fig. 3: Illustration of the state $s_{t,k}$

There parameters are calculated by the following equations:

$$s_{t,k}^X = \{C_X, p_X^{rated}, SoC_X^d, SoC_{t,X}, r_{t,X}\}, \quad X \in \{A, B, C\}$$

$$C_X = \frac{\sum_{k \in X} C_k}{N_X}$$

$$p_X^{rated} = \frac{\sum_{k \in X} p_k^{rated}}{N_X}$$

$$SoC_X^d = \frac{\sum_{k \in X} SoC_k^d \cdot C_k}{\sum_{k \in X} C_k}$$

$$SoC_{t,X} = \frac{\sum_{k \in X} SoC_{t,k} \cdot C_k}{\sum_{k \in X} C_k}$$

$$r_{t,X} = \sum_{k \in X} (t_k^d - t)$$

(16)

Please note that, aside from the aforementioned shared attributes, Group$_A$, Group$_B$ and Group$_C$ also have unique attributes. For Group$_A$, where the charging actions have already been determined, it is pivotal to get the proportion of total power consumption $p_A^{total}$ in relation to the maximum power capacity of the charging station:

$$p_A^{total} = \sum_{k \in A} p_{t,k}$$

(17)

Similarly, for Group$_A$ and Group$_C$, which have multiple piles, it is imperative to obtain the occupancy rate of connected EVs Occ$_A$ and Occ$_C$ relative to the total number of charging piles:

$$\begin{aligned} Occ_A &= N_A/N \\ Occ_C &= N_C/N \end{aligned}$$

(18)

As for EV group B, which only contains one EV$_k$, it is essential to know the relative location of the connected charging pile within the overall charging station. Given the sequential decision-making model, location of Group$_B$ $l_B$ can be easily calculated through the following equation:

$$l_B = \frac{k}{N}$$

(19)

In addition of the three groups' states, $s_{t,k}^O$ denotes the overall state, including the time $t$ and current electricty price $\rho_t$:

$$s_{t,k}^O = \{t, \rho_t\}$$

(20)

Due to the complexity of the state $s_{t,k}$ involved, all components are meticulously listed in Table I below for easy reference.

Due to the special design of the state $s_{t,k}$, the proposed method follows a sequential decision-making process. At each time step $t$, when deciding the charging action $a_{t,k}$ for pile$_k$, it's essential to ensure that the charging actions for pile$_1$ to

| Type | Name | Description |
|------|------|-------------|
| $s_{t,k}^O$ | $t$ | time |
| | $\rho^t$ | current electricity price |
| $s_{t,k}^A$ | $C_A$ | average battery capacity |
| | $p_A^{rated}$ | average rated charging power |
| | $SoC_{t,A}$ | weighted average current SoC |
| | $SoC_A^d$ | weighted average target SoC |
| | $r_{t,A}$ | average remain parking time |
| | $Occ_A$ | occupancy rate of piles |
| | $p_A^{total}$ | total power consumption |
| $s_{t,k}^B$ | $C_B$ | battery capacity |
| | $p_B^{rated}$ | rated charging power |
| | $SoC_{t,B}$ | current SoC |
| | $SoC_B^d$ | target SoC |
| | $r_{t,B}$ | remain parking time |
| | $l_B$ | relative location of connected pile |
| $s_{t,k}^C$ | $C_C^{norm}$ | average battery capacity |
| | $p_C^{rated}$ | average rated charging power |
| | $SoC_{t,C}$ | weighted average current SoC |
| | $SoC_C^d$ | weighted average target SoC |
| | $r_{t,C}$ | average remain parking time |
| | $Occ_C$ | occupancy rate of piles |

TABLE I: COMPONENTS OF STATE $s_{t,k}$

pile$_{k-1}$ have already been determined, while pile$_{k+1}$ to pile$_N$ are still awaiting decisions. The state $s_{t,k}$ of pile$_k$ depends on the actions already decided before it. Conseqently, acquiring the corresponding states sequentially, which is outlined in Algorithm 1, is necessary.

---

**Algorithm 1** Extract $s_{t,k}$ of each pile$_k$ sequentially

---

1: **Input:** EV charging scenario $\chi$
2: **Output:** state-action matrix $S$
3: **for** $i \leftarrow 1$ to 64 **do**
4:     Randomly initialize charging scheduling $a_i$ on $\chi$
5:     Find the local optimal $a_i^*$ on $\chi$ by solving **CO**
6:     Calculate the objective value $f_i$ of $a_i^*$
7: **end for**
8: Select $a_i^*$ with the lowest $f_i$ as $a^*$
9: **for** $t = \tau : \tau + T - 1$ **do**
10:     **for** $k = 1 : N$ **do**
11:         **if** EV$_k$ is connected to pile$_k$ **then**
12:             Extract the state $s_{t,k}$ of pile$_k$ from $\chi$ and $a^*$
13:             Get charging action $a_{t,k}^*$ of pile$_k$ from $a^*$
14:             Add $\langle s_{t,k}, a_{t,k}^* \rangle$ to $S$
15:         **end if**
16:     **end for**
17: **end for**

---

With the definition of $s_{t,k}$ and state extraction algorithm, we can easily extract the pair of state vector and scalar optimal action $\langle s, a^* \rangle$. Combining all pairs to the state-action matrix $S$, we then get the expert dataset that can be used for training
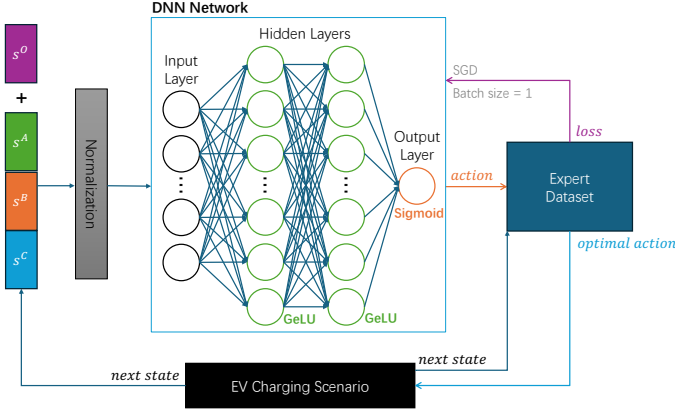
Fig. 4: Architecture of the DNN Network.

in IL later.

### C. Behaviour Cloning

Due to the significant computational requirements of optimization methods like **P1**, it is impractical to make real-time decisions for charging actions of each pile in real-world applications. Therefore, our goal is to identify an alternative method that significantly reduces the computational load while maintaining optimization effectiveness as much as possible. Deep neural networks (DNN), with their powerful high-dimensional representation capabilities and rapid inference speed, are the preferred solution for addressing this issue. The EV charging scheduling scenario, which is defined in **P1**, can be regarded as a finite Markov Decision Process (MDP) over discrete time steps. At each time step $\tau$, we can observe the real-time state $s_{\tau,0}$ of the first charging pile$_0$. Based on this information, the most suitable charging action $a_{\tau,0}$ is determined. Since the proposed model is sequential, once the charging action for pile$_0$ is decided, the new state $s_{\tau,1}$ of pile$_1$ can be obtained. After all piles have made their decisions at time step $t$, the state $s_{\tau+1,0}$ of pile$_0$ at time step $\tau+1$ can be calculated. In summary, after executing an action, a new system state can always be obtained. For a typical MDP, there are two main approaches to training the DNN: reinforcement learning (RL) and imitation learning (IL). Given the availability of a well-defined optimal solution (the expert dataset in section III-B), IL is evidently the more stable and cost-effective choice.

Utilizing behaviour cloning (BC), IL also tries to learn a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ compatible with all charging piles. The policy $\pi$ maps the state $s_{\tau,k}$ of each pile$_k$ to the optimal action $a_{\tau,k}$ at time step $\tau$. In this particular application, we use the state $s_{t,k}$ as the input features and the optimal action $a_{t,k}^*$ obtained from the expert dataset as the expected labels to train the DNN as the approximator to the optimal policy $\pi^*$. Theoretically, DNN is very powerful and has the ability to approximate any bounded continous function. However, there's no theoretical ganrantee that the optimal policy $\pi^*$ can be written as a bounded continous function.

The architecture of the proposed DNN is shown in Figure 4, which is a fully-connected feed-forward network with two hidden layers. The input is a 21-dimention normalized vector $s$ defined in Table I, and the output is a scalar action $a_{t,k}^* \in [0, 1]$. The activation function used in the hidden layers is the Gaussian Error Linear Unit (GeLU), which demonstrates significantly superior performance compared to the Rectified Linear Unit (ReLU) in this specific problem. For the output layers, as the range of output is $[0, 1]$, Sigmoid activation function, which can transform any value to a scalar between $0$ and $1$, is chosen.

As a regression task, the mean squared error (MSE) function is used to evaluate the loss. To minimize the loss function, the Adam optimizer with mini-batch mode is commonly used. However, this training method is not suitable for the problem at hand. The need to consider fairness involves collaboration among different charging piles. This collaboration is implicit since the network is provided with states that do not contain complete information. Maintaining fairness in charging through implicit collaboration represents one solution space of the model. Simultaneously, the model also need to minimize the charging costs. Unfortunately, the solution spaces for cost reduction and fairness maintenance are in conflict. In practice, cost reduction typically receives more attention, thus it is given greater weight. The significant contribution of the highly weighted cost reduction objective to the total loss function causes its solution space's attraction basins to become deeper or broader. Consequently, the Adam optimizer with an adaptive learning rate easily falls into the local optima.

---

**Algorithm 2** Training the DNN network via IL

---

1: **Input:** Training dataset $S$ and validation dataset $V$
2: **Output:** DNN's parameters $\theta^*$ of best performance
3: Randomly initialize $\theta$
4: **for** epoch $\leftarrow 1$ to $100$ **do**
5:      **for** $<s, a^*> \in S$ **do**
6:          Calculate the estimated action $\pi(s; \theta)$
7:          Calculate the loss $L(\theta, S) = \frac{1}{n} \sum_s^S (\pi(s; \theta) - a_s^*)^2$
8:          Update $\theta$ by SGD optimization algorithm
9:      **end for**
10:      Calculate the estimated action vector $\pi(V; \theta)$
11:      Calculate the loss $L(\theta, V) = \frac{1}{n} \sum_v^V (\pi(v; \theta) - a_v^*)^2$
12:      **if** $L(\theta, V)$ doesn't decrease for $5$ epochs **then**
13:          **Break**
14:      **end if**
15: **end for**

---

As an alternative, we employ the Stochastic Gradient Descent (SGD) optimization algorithm. By utilizing different samples for each update, the gradient update direction can vary with each iteration. This stochastic nature enhances SGD's ability to explore the global parameter space, thereby reducing the risk of becoming trapped in local optima. Furthermore, SGD does not incorporate an adaptive learning rate mechanism. Consequently, all objectives contribute equally to the overall gradient during each update, mitigating the bias that may arise from disproportionate weights assigned to specific objectives.

The early termination technology is inculded to avoid possible overfitting. To verify the scalability of the DNN agent, the validation dataset is bulid from a larger scale EV fleet
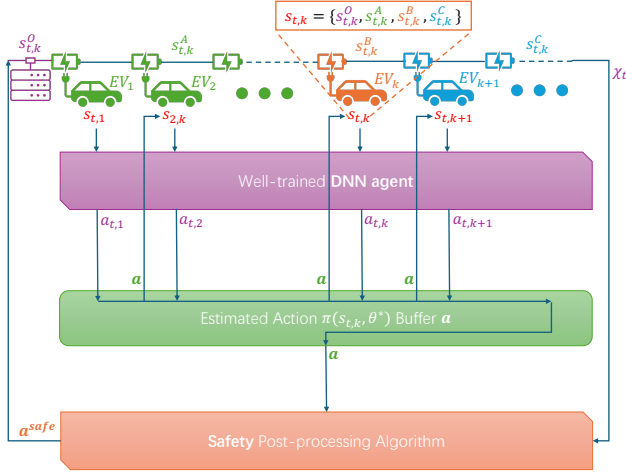
Fig. 5: Illustration of the proposed sequential scheduling.

than the training dataset's. If MSE loss doesn't decrease for a certain number of epochs, the early termination is activated and return the $\boldsymbol{\theta}^*$ which has the best performance on the validation dataset.

### D. Real-Time Online EV Charging Scheduling

Please note that the construction of the expert dataset in Section III-B, as well as the training of the DNN in Section III-C, are entirely conducted offline. Once the well-trained DNN agent is obtained, it is deployed to the CCC to commence online real-time scheduling for charging of EVs. As shown in Figure 5, the order of the decision-making and the extraction of states from the optimal solution $\boldsymbol{A}$ remain consistent. At each time step $\tau$, each pile$_k$'s state $s_{\tau,k}$ is computed sequentially as shown in Table I. The state $s_{\rho,k}$ is used as the input of the DNN agent, which returns the estimated charging action $\pi(s_{\tau,k}; \boldsymbol{\theta}^*)$.

---

**Algorithm 3** Real-Time Online EV Charging Scheduling

1: **Input:** DNN's parameters $\boldsymbol{\theta}^*$, real-time EV scenario $\chi_t$
2: **Output:** Real-time EV charging scheduling
3: Load the parameters $\boldsymbol{\theta}^*$ of the well-trained DNN.
4: Initiallize the action buffer $\boldsymbol{a}$
5: **for** $t = \tau : \tau + T - 1$ **do**
6:     Clear $\boldsymbol{a}$
7:     **for** $k = 1 : N$ **do**
8:         **if** EV$_k$ is connected to charging pile$_k$ at time step $t$ **then**
9:             Calculate pile$_k$'s state $s_{t,k}$ from $\chi_t$ and $\boldsymbol{a}$
10:            Get estimated action $\pi(s_{t,k}; \boldsymbol{\theta}^*)$ from DNN
11:            Store $\pi(s_{t,k}; \boldsymbol{\theta}^*)$ in $\boldsymbol{a}$
12:         **end if**
13:     **end for**
14:     Apply the safety post-processing by Algorithm 4 to $\boldsymbol{a}$ and get the safe charging action vector $\boldsymbol{a_{\text{safe}}}$
15:     Excute the final action $\boldsymbol{a_{\text{safe}}}$ on each charging pile
16: **end for**

---

Please note that the estimated charging action $\pi(s_{t,k}; \boldsymbol{\theta}^*)$ provided by the DNN agent is not the final decision. All estimated charging actions $\pi(s_{t,k}; \boldsymbol{\theta}^*)$ have to be re-uploaded to

the CCC for safety post-processing, which will be introduced in Section III-E. After post-processing, the CCC will issue the ultimate safe charging decision $a_{t,k}^{\text{safe}}$ to all charging piles for execution. This process above will keep iterating until the last EV departs.

### E. Safety Measures

---

**Algorithm 4** Safety post-processing algorithm

1: **Input:** Estimated action vector $\boldsymbol{a}^*$ by DNN, real-time EV scenario $\chi_t$
2: **Output:** Ultimate safe action vector $\boldsymbol{a}^{\text{safe}}$
3: **for** $t = 0 : t_{\text{end}}$ **do**
4:     **for** $k = 1 : N$ **do**
5:         **if** EV$_k$ is connected to charging pile $k$ at time step $t$ **then**
6:             Update SoC$_{t+1,k}$ of EV$_k$ by $\boldsymbol{a}_k^*$ and $\chi_t$
7:             **if** SoC$_{t+1,k} > $ SoC$_k^{\text{d}}$ **then**
8:                 Adjust $\boldsymbol{a}_k^*$ to make SoC$_{t+1,k} = $ SoC$_k^{\text{d}}$
9:             **end if**
10:         **end if**
11:     **end for**
12:     Save the adjusted $\boldsymbol{a}^*$ to the intermediate vector $\boldsymbol{a}'$
13:     **for** $k = 1 : N$ **do**
14:         **if** EV$_k$ is connected to charging pile $k$ at time step $t$ **then**
15:             Calculate the estimated charging power $p_{t,k}$ by $\boldsymbol{a}'[k]$ and $\chi_t$
16:         **else**
17:             Set $p_{t,k}$ to 0
18:         **end if**
19:     **end for**
20:     **if** $\sum_k^N p_{t,k} > P_{\text{rated}}$ **then**
21:         Update $\boldsymbol{a}'$ by scaling it down equally to make $\sum_k^N p_{t,k} = P_{\text{rated}}$
22:     **end if**
23:     **for** $k = 1 : N$ **do**
24:         **if** EV$_k$ is connected to charging pile $k$ at time step $t$ **then**
25:             Update SoC$_{t+1,k}$ of EV$_k$ by $\boldsymbol{a}_k'$ and $\chi_t$
26:             **if** SoC$_{t+1,k} + \frac{p_{t,k}^{\text{rated}} \cdot e_k \cdot \Delta t \cdot (t_k^{\text{d}} - (t+1))}{C_k} < $ SoC$_{t+1,k}^{\text{d}}$ **then**
27:                 Allocate the remain power $P_{\text{rated}} - \sum_k^N p_{t,k}$ to EV$_k$, until the charging requirement SoC$_{t+1,k}^{\text{d}}$ is fulfilled.
28:                 Update $\boldsymbol{a}_k'$ according to the allocated power and make sure $\boldsymbol{a}_k' \leq 1$
29:             **end if**
30:         **end if**
31:     **end for**
32:     Save the adjusted $\boldsymbol{a}'$ to $\boldsymbol{a}^{\text{safe}}$
33: **end for**

---

The estimated charging actions $\pi(s; \boldsymbol{\theta}^*)$ obtained from the DNN agent can't be excuted directly due to the hard constraints. In the particular charging scheduling problem, violating the constraints could lead to serious consequences

like grid instability, battery overcharge and failing to meet user demands. There will inevitably be some error in any neural network's predictions for real world applications. Therefore, a safety post-porcessing algorithm is proposed to ensure the constraints are fully met for any policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$. After the safety post-porcessing, even a stochastic policy can be adjusted to satisfy the constraints. The safety measures can be mainly divided into three steps:

1) Overcharge protection: to check if any $EV_k$'s $SoC_{t+1,k}$ will be larger than the user's expected $SoC_k^d$ after executing the estimated charging action $\pi(s; \boldsymbol{\theta}^*)$. If such vehicles do exist, make appropriate adjustments and save the changes to intermediate action vector $\boldsymbol{a}'$.

2) Station capacity protection: to check if the total charging power of all EVs exceed the charging station's power capacity after executing $\boldsymbol{a}'$. If such a situation occurs, update $\boldsymbol{a}'$ by scaling down it equally.

3) Undercharging protection: to check if any $EV_k$ can't fulfill the SoC target when leaving after executing $\boldsymbol{a}'_k$, even charging with maximum power $p_k^{\text{rated}}$ for the remaining time steps (except the current one). If such an occurrence arise, the forced charging module will be activated, ensuring the vehicle is charged as much as possible without exceeding the charging station's power capacity. After saving the changes to $\boldsymbol{a}'$, the ultimate safe action vector $\boldsymbol{a}^{\text{safe}} = \boldsymbol{a}'$ is upload to the charging station and all charging piles will excute it in the following time step.

Please refer to Algorithm 4 for the details.

## REFERENCES

[1] S. L. Xuezhong Hu and Y. J. Haibin Wang, "Ac charging pile of electric vehicle and intelligent charging control strategy research," *REVIEWS OF ADHESION AND ADHESIVES*, vol. 11, no. 3, 2023.

[2] A. Zhou, J. Yu, Z. Li, and J. Pu, "A security authentication method between the charging pile and battery management system," in *2021 IEEE 4th International Conference on Electronics Technology (ICET)*. IEEE, 2021, pp. 400–405.

[3] E. ElGhanam, M. Hassan, A. Osman, and I. Ahmed, "Review of communication technologies for electric vehicle charging management and coordination," *World Electric Vehicle Journal*, vol. 12, no. 3, p. 92, 2021.

[4] H. Yu, C. Xu, W. Wang, G. Geng, and Q. Jiang, "Communication-Free Distributed Charging Control for Electric Vehicle Group," *IEEE Transactions on Smart Grid*, vol. 15, no. 3, pp. 3028–3039, May 2024.

[5] L. Yao, W. H. Lim, and T. S. Tsai, "A Real-Time Charging Scheme for Demand Response in Electric Vehicle Parking Station," *IEEE Transactions on Smart Grid*, vol. 8, no. 1, pp. 52–62, Jan. 2017.

[6] H. M. Abdullah, A. Gastli, and L. Ben-Brahim, "Reinforcement learning based ev charging management systems–a review," *IEEE Access*, vol. 9, pp. 41 506–41 531, 2021.

[7] M. Dorokhova, Y. Martinson, C. Ballif, and N. Wyrsch, "Deep reinforcement learning control of electric vehicle charging in the presence of photovoltaic generation," *Applied Energy*, vol. 301, p. 117504, Nov. 2021.

[8] Z. Wan, H. Li, H. He, and D. Prokhorov, "Model-Free Real-Time EV Charging Scheduling Based on Deep Reinforcement Learning," *IEEE Transactions on Smart Grid*, vol. 10, no. 5, pp. 5246–5257, Sep. 2019.

[9] W. Wu, Y. Lin, R. Liu, Y. Li, Y. Zhang, and C. Ma, "Online EV Charge Scheduling Based on Time-of-Use Pricing and Peak Load Minimization: Properties and Efficient Algorithms," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 1, pp. 572–586, Jan. 2022.

[10] Z. Yi, D. Scoffield, J. Smart, A. Meintz, M. Jun, M. Mohanpurkar, and A. Medam, "A highly efficient control framework for centralized residential charging coordination of large electric vehicle populations," *International Journal of Electrical Power & Energy Systems*, vol. 117, p. 105661, May 2020.

[11] S. Hajforoosh, M. A. S. Masoum, and S. M. Islam, "Online optimal variable charge-rate coordination of plug-in electric vehicles to maximize customer satisfaction and improve grid performance," *Electric Power Systems Research*, vol. 141, pp. 407–420, Dec. 2016.

[12] S. Thangavel, D. Mohanraj, T. Girijaprasanna, S. Raju, C. Dhanamjayulu, and S. M. Muyeen, "A Comprehensive Review on Electric Vehicle: Battery Management System, Charging Station, Traction Motors," *IEEE Access*, vol. 11, pp. 20 994–21 019, 2023.

[13] H. Li, G. Li, T. T. Lie, X. Li, K. Wang, B. Han, and J. Xu, "Constrained large-scale real-time EV scheduling based on recurrent deep reinforcement learning," *International Journal of Electrical Power & Energy Systems*, vol. 144, p. 108603, Jan. 2023.

[14] C. A. Floudas, I. G. Akrotirianakis, S. Caratzoulas, C. A. Meyer, and J. Kallrath, "Global optimization in the 21st century: Advances and challenges," *Computers & Chemical Engineering*, vol. 29, no. 6, pp. 1185–1202, May 2005.

[15] S. Gao, R. b. Z. Lee, Z. Huang, C. Xiang, M. Yu, K. T. Tan, and T. H. Lee, "A Hybrid Approach for Home Energy Management With Imitation Learning and Online Optimization," *IEEE Transactions on Industrial Informatics*, pp. 1–13, 2023.