

A hybrid policy gradient and rule-based control framework for electric vehicle charging

Brida V. Mbuwir^{a,b}, Lennert Vanmunster^b, Klaas Thoelen^{a,b,*}, Geert Deconinck^{a,b}

^a*EnergyVille, Thor Park 8310, 3600 Genk, Belgium*

^b*ESAT-Electa, KU Leuven, Kasteelpark Arenberg 10 bus 2445, 3001 Leuven, Belgium*

Abstract

Recent years have seen a significant increase in the adoption of electric vehicles, and investments in electric vehicle charging infrastructure and rooftop photo-voltaic installations. The ability to delay electric vehicle charging provides inherent flexibility that can be used to compensate for the intermittency of photo-voltaic generation and optimize against fluctuating electricity prices. Exploiting this flexibility, however, requires smart control algorithms capable of handling uncertainties from photo-voltaic generation, electric vehicle energy demand and user's behaviour. This paper proposes a control framework combining the advantages of reinforcement learning and rule-based control to coordinate the charging of a fleet of electric vehicles in an office building. The control objective is to maximize self-consumption of locally generated electricity and consequently, minimize the electricity cost of electric vehicle charging. The performance of the proposed framework is evaluated on a real-world data set from EnergyVille, a Belgian research institute. Simulation results show that the proposed control framework achieves a 62.5% electricity cost reduction compared to a business-as-usual or passive charging strategy. In addition, only a 5% performance gap is achieved in comparison to a theoretical near-optimal strategy that assumes perfect knowledge on the required energy and user behaviour of each electric vehicle.

Keywords: electric vehicles, smart charging, proximal policy optimization, reinforcement learning

*Corresponding author

Email address: `klaas.thoelen@kuleuven.be` (Klaas Thoelen)

1 List of symbols

\mathcal{A}	Action space
a_t	Action at time step t
${}^i C_t^{anx}$	Anxiety cost of EV_i at time step t
$C_{t,m}^{cons}$	Consumption cost in minute m during time step t
$C_{t,m}^{elec}$	Electricity cost (consumption + injection) in minute m of t
C^{inj}	Electricity injection cost
\mathbb{E}	Expected value
${}^i E_t^{ch}$	Energy used to charge EV_i during time step t
${}^i E_t^{ch,ses}$	Total energy charged in the session of EV_i up to time step t
${}^i E_t^{rem}$	Remaining energy [kWh] for EV_i in the RTC algorithm, i.e. the energy that still needs to be charged before the end of time step t
${}^i E_{t,m}^{req}$	Required energy [kWh], i.e. energy that still needs to be charged before departure, of EV_i in minute m of time step t
G	Expected (discounted) return
K_1, K_2, K_3	Coefficients in anxiety cost function ${}^i C_t^{anx}$
L_c	Length of a control time step [minutes]
L^{CLIP}	Clipped loss function in PPO
L^{CLIP2}	Augmented, clipped loss function in PPO
L^{ORIG}	Unconstrained loss function in PPO
m	Minute index in a control slot
${}^i M_t^{end}$	Index of last minute of EV_i in time step t
${}^i M_t^{start}$	Index of first minute of EV_i in time step t
N_{ev}	Fleet size, i.e. number of EVs controlled by the agent
N_{fut}	Number of future time steps in state vectors for PV and electricity price forecasts
N_{par}	Number of power partitions in the heuristic dispatch algorithm in the aggregate MDP
N_{past}	Number of past time steps in state vectors of the PV forecast
P_t^a	Aggregate fleet charging power [kW] action (denormalised output of the actor network in the aggregate MDP) at time step t
${}^i P_t$	Charging power [kW] action (denormalised output of the actor network in the base and hidden MDPs) of EV_i at time step t

iP_t^b	Clipped charging power [kW] (output of the backup controller or heuristic dispatch) of EV _{<i>i</i>} at time step <i>t</i>
$iP_t^{b,max}$	Maximum charging power [kW] for backup controller or heuristic dispatch of EV _{<i>i</i>} at time step <i>t</i>
$iP_t^{b,min}$	Minimum charging power [kW] for backup controller or heuristic dispatch of EV _{<i>i</i>} at time step <i>t</i>
P^{cons}	Historical grid power consumption [kW]
P^{inj}	Historical grid power injection [kW]
iP^{lim}	Absolute maximum charging power [kW] of EV _{<i>i</i>}
\mathbf{P}_t^{fpv}	State vector at time <i>t</i> of PV power forecast [kW]
$\mathbf{P}_t^{fpv,fut}$	State vector at time <i>t</i> of future PV power forecasts [kW] (after <i>t</i>)
$\mathbf{P}_t^{fpv,past}$	State vector at time <i>t</i> of past PV power forecasts [kW] (before <i>t</i>)
$P_t^{fpv,rest}$	State parameter containing the average PV power forecast [kW] of the rest of the day after $\mathbf{P}_t^{fpv,fut}$
³ $P_m^{max,\Sigma}$	Total maximum charging power [kW] in the RTC algorithm for minute <i>m</i>
$P_m^{min,\Sigma}$	Total minimum charging power [kW] in the RTC algorithm for minute <i>m</i>
P_t^{pv}	PV power generation [kW] at time <i>t</i> , minute 0 ($P_{t,m=0}^{pv}$)
\mathbf{P}_t^{pv}	State vector of PV power generation [kW] at time step <i>t</i>
$P_{t,m}^{pv}$	PV power generation [kW] in minute <i>m</i> of time step <i>t</i>
$P^{pv,net}$	Historic net PV power generation [kW] available for charging the EV fleet
$\mathbf{P}_t^{pv,past}$	State vector of past PV power generation [kW] at time step <i>t</i>
$P^{pv,scaled}$	Scaled PV power generation [kW], $P^{pv,scaled} \equiv 0.2 \times P^{pv,tot}$
$P^{pv,tot}$	Historic PV power generation [kW]
$i\mathbf{P}_t^r$	Vector with the charging powers [kW] of EV _{<i>i</i>} in each minute <i>m</i> of time step <i>t</i>

$iP_{t,m}^r$	Charging power [kW] of EV_i in minute m of time step t
$iP_m^{r,max}$	Maximum charging power [kW] in the RTC algorithm for EV_i in minute m in time step t
$iP_m^{r,min}$	Minimum charging power [kW] in the RTC algorithm for EV_i in minute m in time step t
pr	Probability ratio in PPO algorithm
Q	Action-value function
ρ	Reward function
\mathcal{S}	State space
s_t	State at time t
t	Time step index in an episode
T	Final time step or number of time steps in an episodic environment
f	Transition function
iT^{arr}	Arrival time of EV_i (unit = time slot)
⁴ iT^{dep}	Departure time of EV_i (unit = time slot)
V	State-value function
$i\mathbf{X}_t$	State vector of EV_i at time step t
\mathbf{Z}_t	Vector containing the aggregated fleet state parameters at time step t
α	Learning rate
γ	Discount factor
ΔiT^{dep}	Time left until departure for EV_i at time t (unit = time steps)
θ	Function approximator parameters (e.g. neural network weights)
κ	Flexibility factor in the RTC algorithm
λ^{Belpex}	Belpex day-ahead electricity price [€/kWh]
$\bar{\lambda}^{cons}$	Average TOU grid consumption price in the current day (between 7:00 and 20:00)
λ_t^{cons}	State vector of TOU grid consumption price [€/kWh] at time step t

$\lambda_{t,m}^{cons}$	TOU grid consumption price [€/kWh] in minute m of time step t
$\boldsymbol{\lambda}_t^{cons,fut}$	State vector at time step t of future TOU grid consumption prices [€/kWh] (after t)
$\lambda_t^{cons,rest}$	State parameter containing the average TOU grid consumption prices [€/kWh] of the rest of the day after $\boldsymbol{\lambda}_t^{cons,fut}$
λ^{inj}	Electricity injection price [€/kWh]
Λ	Advantage function
π	Policy function
$i_{\mathcal{T}}$	Flexibility of EV_i in the heuristic dispatch algorithm of MDP_{agg}
∇	Nabla-operator

6 1. Introduction

7 The increasing concern on the effects of greenhouse gas emissions has
8 led to an increase in the use of renewable energy sources (RESs) and the
9 electrification of transport. While the decreasing cost of photo-voltaic (PV)
10 installations has led to an increase in the number of buildings with rooftop
11 PV installations, the electrification of transport has led to several incentive
12 programs to encourage the use of electric vehicles (EVs). For example, the
13 EV30@30 campaign has set a target of at least 30% market share of EVs in
14 the Electric Vehicle Initiative member states by 2030 [1]. The increase in the
15 number of EVs, however, significantly alters the electricity demand curve [2].
16 A typical example of an alteration of the electricity demand curve can be
17 seen in office buildings with EV charging infrastructure where EVs tend to
18 arrive at the same time, see Fig. 1.

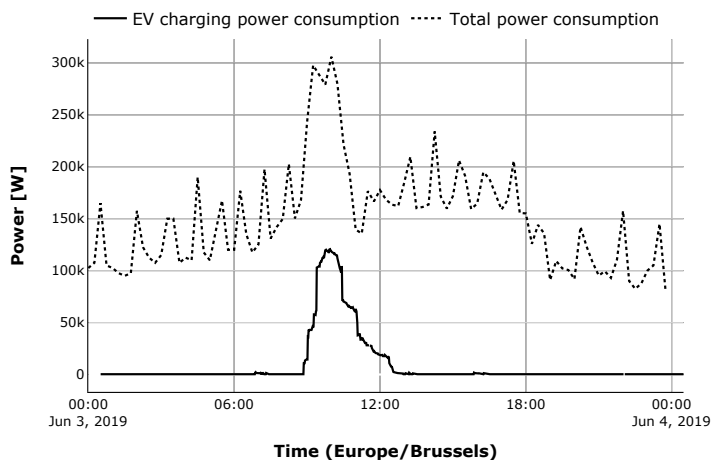


Figure 1: The typical power demand caused by EV charging in the morning can have a significant impact on the total power consumption of an office building. For example, on the morning of June 3rd, 2019, at the EnergyVille building.

19 In the past years, EV charging has mainly been passive or uncontrolled
20 i.e. charging is activated immediately and at maximum charging power when
21 an EV is plugged into a charging station. In this paper, this type of charging
22 is referred to as business-as-usual (BAU) charging. Such uncontrolled EV

23 charging does not exploit the inherent flexibility¹ of EV charging; typically
24 an EV is plugged in longer than the time needed to fully charge its battery.
25 If EV fleet charging is controlled, the flexibility harnessed can be used for a
26 range of objectives.

27 Recently, several control methods have emerged for controlling the charg-
28 ing of EVs in order to harness flexibility for objectives such as: avoiding grid
29 congestion problems [3], maximising self-consumption of local electricity gen-
30 eration [4] and load flattening [5]. These methods range from rule-based to
31 model-based and model-free (data-driven). Rule-based methods rely on pre-
32 defined rules and conditions expressed in the form “if *condition*, do *action*”
33 statements to determine a control policy for the control agent. These rules
34 are typically handcrafted and to guarantee an adequate performance of a
35 rule-based controller, considerable expert knowledge is required to correctly
36 set the threshold values, and tune the system parameters. Rule-based control
37 has been widely used in EV charging due to its simplicity and computational
38 efficiency for uninterrupted EV charging [6] and prevention of grid overload-
39 ing [7]. However, since the rules are tailored towards a specific system and
40 objective, the method cannot be easily generalised.

41 Model-based control methods, on the other hand, require an explicit def-
42 inition of the system dynamics in order to establish a control policy. Model
43 predictive control, for instance, has been extensively applied in literature
44 for EV charging to minimize energy costs [8] and for voltage control [9, 10].
45 While successful, their performance relies on the accuracy of the model, and
46 a mismatch between the model and the real system will result in sub-optimal
47 operation. In the EV charging context, identifying a (sufficiently) accurate
48 model is challenged by the heterogeneity of EV models and unpredictability
49 of EV-user behavior.

50 In contrast to model-based methods, model-free methods do not rely on
51 explicit knowledge of a system model. Instead, they learn a control pol-
52 icy from system observations collected a-priori (batch or offline learning) or
53 through online interactions with the system. These methods are therefore
54 data-driven, which renders them flexible and more generalisable compared
55 to model-based and rule-based methods. The most popular model-free tech-

¹The available EV flexibility can be described based on the number of hours that the EV charging can be delayed while meeting the user’s departure deadline and respecting the constraints on battery capacity, maximum charging rate, and additional constraints of the charging infrastructure.

56 nique in EV charging literature is reinforcement learning (RL) [11, 12]. Fitted
57 Q-iteration, a batch RL technique, has been used to control EV charging for
58 load flattening purposes [5], electricity cost savings based on day-ahead mar-
59 ket prices [13] and long-term cost optimization [14]. Wang *et al.* [15] used
60 an online RL algorithm, SARSA, to schedule EV charging for minimizing
61 electricity cost. Deep Q-learning, another RL algorithm, has been used to
62 minimize electricity cost based on real-time electricity pricing [16], minimize
63 long-term operating cost [16], as well as for load flattening purposes [17].
64 The above mentioned RL algorithms are based on the standard Q-learning
65 algorithm [18], which relies on Q-values (of state-action pairs) for the evalu-
66 ation and selection of control actions, and consequently, learning of a control
67 policy. To efficiently compute these Q-values, the action space for the learn-
68 ing agent is required to be finite and discrete to avoid a heavy computational
69 burden and the curse of dimensionality. Even though the above mentioned
70 methods employ regression algorithms such as neural networks to approxi-
71 mate the Q-values through a Q-function, fine-grained discrete actions lead to
72 large (continuous) action spaces, which in turn lead to an intractable com-
73 putation of Q-values. In the context of EV charging, discretizing the action
74 space limits full exploitation of EV flexibility since the charging powers are
75 continuous values. To allow the learning of control policies in systems with
76 large or continuous action spaces, policy gradient methods were introduced
77 [11].

78 Policy gradient techniques directly optimize a control policy without a
79 need for Q-values to select control actions. Several policy gradient methods
80 have been employed in EV charging literature. Yu *et al.* [19] used deep
81 deterministic policy gradient to minimize electricity costs in a smart home
82 with electricity generation from RES and multiple loads including EVs and
83 HVAC systems. In [20], the authors proposed prioritised deep deterministic
84 policy gradient for the coordination of EV fleet charging by an aggregator
85 with the aim of maximizing profit (through vehicle-to-grid capabilities) and
86 minimizing EV charging electricity costs. The authors showed that priori-
87 tised deep deterministic policy gradient outperforms standard Q-learning and
88 deep Q-learning. Trust region policy optimization was used for home energy
89 management in which the charging of an EV was controlled to minimize elec-
90 tricity cost [21]. Moonens and Nowé [22] used policy proximal optimization
91 (PPO) to coordinate charging of an EV fleet for load balancing purposes. The
92 authors showed that this method outperformed the BAU scheme by reducing
93 the number of electricity consumption peaks caused by EV charging.

94 Motivated by the success of RL, and particularly policy gradient tech-
95 niques for EV charging, this work builds on existing literature and proposes
96 a novel control framework for EV charging in a work environment with the
97 objective of maximising self-consumption of local electricity generation and
98 minimizing electricity cost. The proposed control framework combines PPO
99 and rule-based control allowing a quick response of the control agent to the
100 stochastic PV generation. The main contributions of this work are sum-
101 marised as follows:

- 102 • A novel control framework is proposed combining the strengths of PPO
103 (model-free, data-driven, ability to deal with continuous actions) with
104 those of rule-based control (low computational complexity). In the
105 proposed framework, a RL agent learns a control policy in a low time
106 resolution (60 minutes, 15 minutes or 5 minutes), which is refined by
107 a rule-based controller during real-time operation (one minute time
108 step) to ensure a more optimal real-time control. This contrasts with
109 existing literature in which control actions are predominantly taken at
110 an hourly resolution [5, 21].
- 111 • A demonstration of the scalability of the proposed framework by using
112 the three-step approach introduced by Vandael *et al* [23] in which a
113 RL agent learns the optimal aggregate charging power for an entire EV
114 fleet. In contrast to the original work, this paper uses PPO to learn
115 the aggregate charging power.
- 116 • Case study demonstrating that the proposed control method achieves
117 an increase in self-consumption of local electricity generation and reduc-
118 tions in the net electricity costs compared to the BAU scheme. Com-
119 pared to a “perfect information optimum” (PIO) strategy, the proposed
120 control method is slightly less performant. The PIO strategy is based
121 on sequential quadratic programming [24] and assumes full knowledge
122 of the EV user behaviour and PV electricity generation.

123 The performance of the proposed framework is evaluated on a real-world
124 data set from EnergyVille², a research institute in Belgium. The EV charg-
125 ing infrastructure at EnergyVille contains charging stations from 7 different

²www.energyville.be

126 brands, totalling in 27 connection points. While most charging stations are
127 conventional 22kW AC charging stations, the setup includes an AC/DC fast-
128 charger and a vehicle-to-grid (V2G) charging station. In total the charging
129 stations amount to an installed EV charging power of 530kVA, while the ca-
130 pacity at the electrical cabinet is only 436kVA. At the moment, 12 charging
131 stations are fully monitored and controllable via OCPP version 1.6³, and
132 more are expected in the short term. A custom IT infrastructure has been
133 deployed to monitor and control the charging sessions at EnergyVille. It pro-
134 vides the following data: the arrival time detected when an EV is plugged
135 in to a charging station, the maximum charging power measured three min-
136 utes into the charging session, and estimates of the departure time and total
137 energy needed to fully charge the EV are collected as user input via a web
138 app. A data set of all OCPP-controlled charging sessions since August 2018
139 is available. In addition, a 368kWp PV system is installed on the rooftop
140 of the EnergyVille building. The energy yield of this installation can either
141 be used for consumption within the building or injected into the grid. For
142 the latter, a fixed fee is paid per kWh; on the other hand, consumption
143 from the grid has a dynamic component with hourly periodicity based on
144 the day-ahead market price. A data set is available for the PV production
145 at EnergyVille (since April 2016). Finally, the work described in this paper
146 also relies on the availability of day-ahead market prices and a regional PV
147 production forecast prior to a charging session.

³Open Charge Point Protocol, www.openchargealliance.org/protocols/ocpp-16/

148 2. Problem description and Markov decision process

149 This work considers the problem of coordinating the charging of a fleet of
150 EVs in an office building with a rooftop PV installation. The following prob-
151 lem description aims to be generic for this type of buildings, but employs
152 some specific details of the EnergyVille building where needed. Charging
153 transactions are characterised by: the arrival time T^{arr} , the departure time
154 T^{dep} , the energy E^{req} [kWh] required to fully charge the EV, and the maxi-
155 mum charging power P^{lim} [kW] of the EV. The objective is to charge a fleet
156 of N_{ev} EVs while maximising self-consumption of the locally generated elec-
157 tricity - from the PV installation - and minimising electricity cost; the EVs
158 can be charged using locally generated electricity or directly from the grid.
159 To achieve this objective, the charging power of each connected EV has to be
160 decided at every time step based on the PV electricity generation, the elec-
161 tricity price, and estimates on the departure time of the EV and the energy
162 required to fully charge the EV by that time. This decision making problem
163 encountered at every time step can be expressed as a Markov decision process
164 (MDP), which is the basis for formulating RL problems. However, the deci-
165 sion making problem is challenged by the uncertainty in the PV generation,
166 and the arrival and departure times of the EVs.

167 2.1. Markov decision process formulation

168 A Markov decision process is characterized by: (i) a state space \mathcal{S} de-
169 scribing the finite set of states that the system can be in, (ii) an action space
170 \mathcal{A} consisting of a finite set of possible actions that can change the state
171 of the system, (iii) transition function f representing the system dynamics
172 or probabilities for a stochastic state evolution, and (iv) a reward function,
173 ρ , evaluating each state transition. Three MDP formulations - MDP_{base} ,
174 MDP_{hid} and MDP_{agg} - are presented in the following subsections. The base
175 MDP, MDP_{base} , formulates the EV fleet charging problem described in the
176 previous paragraph. The hidden MDP, MDP_{hid} , is similar to MDP_{base} but
177 does not include information on the estimates of the departure time and en-
178 ergy required to fully charge the EV. The aggregate MDP, MDP_{agg} , builds
179 on MDP_{base} to improve scalability to larger fleet sizes using the three-step
180 approach introduced by Vandael *et al.* [23].

181 2.2. Base Markov decision process

182 The base Markov decision process, MDP_{base} , aims at providing optimal
183 charging schedules for individual EVs and assumes the widest range of infor-

184 mation to be available.

185 *State space*

186 The state space has three components:

- *PV component*: consists of the current PV generation \mathbf{P}^{pv} , and a forecast of the PV generation \mathbf{P}^{fpv} . \mathbf{P}_t^{pv} contains: (i) $\mathbf{P}_t^{pv,past}$ a vector with the average⁴ over each hour of the measured PV generation of the previous N_{past} hours, as shown in (1); and (ii) P_t^{pv} , the PV measurement at time stamp t .

$$\begin{aligned} \mathbf{P}^{pv} &= (\mathbf{P}_t^{pv,past}, P_t^{pv}), \\ \mathbf{P}_t^{pv,past} &= \left[\text{av}^5 \left(P^{pv}, t - j \times \frac{60}{L_c}, t - (j - 1) \times \frac{60}{L_c} \right), j = N_{past}, \dots, 1 \right], \end{aligned} \quad (1)$$

187 with L_c representing the length of a control time step in minutes.

188 As shown in (2), the vector \mathbf{P}_t^{fpv} contains a forecast on the PV genera-
 189 tion in terms of $\mathbf{P}_t^{fpv,fut}$ (the average over each hour⁶ of the forecasted
 190 PV generation P^{fpv} for the next N_{fut} hours), and $P_t^{fpv,rest}$ the average
 191 of P^{fpv} for the rest of the day, and information on the forecast
 192 of the PV generation for the past N_{past} hours, $\mathbf{P}_t^{fpv,past}$. In this work,
 193 N_{fut} and N_{past} are set to a value of 2 through manual tuning and ex-
 194 perimentation by considering the trade-off between training time of the
 195 control agent and gains in electricity cost and self-consumption. Setting
 196 N_{fut} to larger values increases the uncertainty on the variables and also
 197 increases the state space dimension leading to curse of dimensionality
 198 issues and an increase in the training time.

$$\mathbf{P}_t^{fpv} = (\mathbf{P}_t^{fpv,past}, \mathbf{P}_t^{fpv,fut}, P_t^{fpv,rest}), \quad (2)$$

$$\begin{aligned} \text{with } \mathbf{P}_t^{fpv,past} &= \left[\text{av} \left(P^{fpv}, t - j \frac{60}{L_c}, t - (j - 1) \frac{60}{L_c} \right), j = N_{past}, \dots, 1 \right], \\ \mathbf{P}_t^{fpv,fut} &= \left[\text{av} \left(P^{fpv}, t + 1 + (j - 1) \frac{60}{L_c}, t + 1 + j \frac{60}{L_c} \right), j = 1, \dots, N_{fut} \right], \end{aligned}$$

⁴The average is used instead of the actual values for dimensionality reduction purposes.

⁵ $\text{av}(\mathbf{x}, \mathbf{t1}, \mathbf{t2})$ returns the mean value of $\mathbf{x}(\mathbf{t})$ between $\mathbf{t}=\mathbf{t1}$ and $\mathbf{t}=\mathbf{t2}$ and returns $\mathbf{x}(\mathbf{t2})$ if $\mathbf{t1} \geq \mathbf{t2}$

⁶The PV forecast values have a periodicity of 15 minutes.

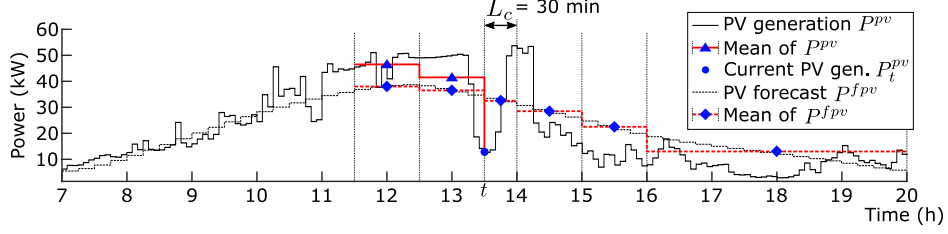


Figure 2: Example of state parameter vectors \mathbf{P}_t^{pv} (triangles + circle) and \mathbf{P}_t^{fpv} (diamonds) with $N_{past} = N_{fut} = 2$ and $L_c = 30$ minutes

201 and $P_t^{fpv,rest} = \text{av} \left(P^{fpv}, t + 1 + N_{fut} \frac{60}{L_c}, T \right)$, with T the final time
 202 step in the optimization horizon.

203 An example of \mathbf{P}^{pv} and \mathbf{P}^{fpv} is shown in Fig. 2.

- 204 • *Price component λ^{cons}* : represents the price of importing a kilowatt-
 205 hour of energy from the grid at time t . λ^{cons} as shown in (3) contains
 206 the price at time step t and information on the forecasted price ($\lambda_t^{cons,fut}$
 207 and $\lambda_t^{cons,rest}$, defined in a similar manner as $\mathbf{P}_t^{fpv,fut}$ and $P_t^{fpv,rest}$).

$$\lambda_t^{cons} = (\lambda_t^{cons}, \lambda_t^{cons,fut}, \lambda_t^{cons,rest}), \quad (3)$$

$$\lambda_t^{cons,fut} = \left[\text{av} \left(\lambda^{cons}, t + 1 + (j - 1) \frac{60}{L_c}, t + 1 + j \frac{60}{L_c} \right), j = 1, \dots, N_{fut} \right]$$

208
 209 and $\lambda_t^{cons,rest} = \text{av} \left(\lambda^{cons}, t + 1 + N_{fut} \times \frac{60}{L_c}, T \right)$.

210 It is important to note that in the event where a price, λ^{inj} , is set for
 211 injecting power to the grid, this price would become part of the price
 212 component of the state space. However, in this work a constant price
 213 for injecting power to the grid is considered.

- 214 • *EV transaction component ${}^i\mathbf{X}_t$* contains three parameters directly re-
 215 lated to an EV charging transaction as shown in (4): (i) ${}^iE_t^{req}$ the
 216 energy required left to fully charge EV i (EV_i) at time t before its
 217 departure; (ii) ${}^iP^{lim}$ the maximum charging power for EV_i and (iii)
 218 $\Delta {}^iT_t^{dep} = {}^iT_t^{dep} - t$ the time left until departure (in number of control
 219 time steps) for EV_i at time t .

$${}^i\mathbf{X}_t = \begin{cases} (0, 0, 0), & \text{if station } i \text{ is unused at time } t, \\ \left({}^iE_t^{req}, {}^iP^{lim}, \Delta {}^iT_t^{dep} \right), & \text{otherwise.} \end{cases} \quad (4)$$

220 In summary, at any time t , the state of the system is defined as follows:

$$\forall s_t \in \mathcal{S}, s_t = \left({}^1\mathbf{X}_t, \dots, {}^{N_{ev}}\mathbf{X}_t, t, \mathbf{P}_t^{pv}, \mathbf{P}_t^{fpv}, \boldsymbol{\lambda}_t^{cons} \right) \quad (5)$$

221 **Action space**

222 At any time step, t , an action $a_t \in \mathcal{A}$ is a vector containing the charging
223 powers iP_t [kW] for all the connected EVs as shown below:

$$\forall a_t \in \mathcal{A}, a_t = \left({}^1P_t, \dots, {}^{N_{ev}}P_t \right). \quad (6)$$

224

225 **Reward function**

226 The reward function consists of two components C^{cons} and C^{inj} represent-
227 ing the cost incurred for charging the EVs from the grid and for injecting
228 power to the grid respectively. This reward function is based on the electric-
229 ity price λ^{cons} and a grid injection price λ^{inj} at time t . The total electricity
230 cost during minute m in control time step t is given as follows:

$$C_{t,m}^{elec} = \begin{cases} C_{t,m}^{cons} = \frac{\sum_{i=1}^{N_{EV}} {}^iP_{t,m} - P_{t,m}^{pv}}{60} \times \lambda_{t,m}^{cons}, & \text{if } \sum_{i=1}^{N_{EV}} {}^iP_{t,m} - P_{t,m}^{pv} \geq 0, \\ C_{t,m}^{inj} = \frac{\sum_{i=1}^{N_{EV}} {}^iP_{t,m} - P_{t,m}^{pv}}{60} \times \lambda^{inj}, & \text{otherwise,} \end{cases} \quad (7)$$

231 where ${}^iP_{t,m}$ is the charging power for EV_i and $P_{t,m}^{pv}$ is the PV generation at
232 minute m of time step t .

233 Based on the cost incurred per minute, the negative reward function for
234 each state transition is defined as shown in (8).

$$\rho(s_t, a_t, s_{t+1}) = - \sum_{m=0}^{L_c-1} C_{t,m}^{elec} \quad (8)$$

235 In this work, we focus on a data-driven RL algorithm, as such, the defini-
236 tion of a transition function for the system dynamics is not required. How-
237 ever, the state of each connected EV is updated as follows:

$${}^i E_t^{ch} = \sum_{m={}^i M_t^{start}}^{{}^i M_t^{end}} {}^i P_{t,m}/60 \quad (9)$$

$${}^i \mathbf{X}_{t+1} = \begin{cases} (0, 0, 0), & \text{if } {}^i E_t^{req} - {}^i E_t^{ch} \leq 0 \text{ or } \Delta {}^i T_t^{dep} - 1 \leq 0, \\ \left({}^i E_t^{req} - {}^i E_t^{ch}, {}^i P^{lim}, \Delta {}^i T_t^{dep} - 1 \right), & \text{otherwise,} \end{cases} \quad (10)$$

238 where ${}^i E_t^{ch}$ is the energy charged by EV_{*i*} during time step *t*. New EVs that
 239 arrived between *t* and *t* + 1 are added to ${}^i \mathbf{X}_{t+1}$. The remaining state param-
 240 eters $\mathbf{P}_t^{pv}, \mathbf{P}_t^{fpv}, \boldsymbol{\lambda}_t^{cons}$ are updated by reading the data from the database.

241 2.3. MDP with unknown SoC and departure time

242 As mentioned earlier, the availability of estimates of T^{dep} and E^{req} is based
 243 on the willingness of the EV users to provide these inputs in practice. They
 244 can be termed as hidden parameters as they cannot be measured directly
 245 without user interaction, and as such, are hidden from the learning agent.
 246 Therefore, to allow learning adequate control policies in scenarios without
 247 user interaction, MDP_{hid} is proposed in this section. The definition of MDP_{hid}
 248 is similar to that of MDP_{base}. However, the state space is different since ${}^i E_t^{req}$
 249 and ${}^i T_t^{dep}$ are not included.

250 **State space**

251 Compared to MDP_{base}, the hidden state parameters ${}^i E_t^{req}$ and ${}^i T_t^{dep}$ in
 252 the state s_t are replaced by two known state parameters: (i) ${}^i E_t^{ch,ses}$ the total
 253 energy already charged in the current session of EV_{*i*} at time *t*, and (ii) the
 254 arrival time ${}^i T^{arr}$. The state at time *t*, s_t , is now defined as:

$$s_t = ({}^1 \mathbf{X}_t, \dots, {}^{N_{ev}} \mathbf{X}_t, t, \mathbf{P}_t^{pv}, \mathbf{P}_t^{fpv}, \boldsymbol{\lambda}_t^{cons}) \quad (11)$$

$$\text{with } {}^i \mathbf{X}_t = \begin{cases} (0, 0, 0), & \text{if station } i \text{ empty at time } t \\ \left({}^i E_t^{ch,ses}, {}^i P^{lim}, {}^i T^{arr} \right), & \text{otherwise.} \end{cases} \quad (12)$$

255 The remaining state parameters ${}^i P^{lim}, t, \mathbf{P}_t^{pv}, \mathbf{P}_t^{fpv}, \boldsymbol{\lambda}_t^{cons}$ are the same as in
 256 MDP_{base}.

257 **Action space**

258 The action space is the same as MDP_{base} .

259 **Reward function**

260 The reward function consists of two terms as shown in (13).

$$\rho(s_t, a_t, s_{t+1}) = - \sum_{m=0}^{L_c-1} C_{t,m}^{elec} - \sum_{i=1}^{N_{ev}} {}^i C_t^{anx} \quad (13)$$

$$\text{with } {}^i C_t^{anx} = \begin{cases} 0, & \text{if } \Delta {}^i T_t^{dep} > 0, \\ K_1 \times \frac{({}^i E_t^{req})^{K_2}}{({}^i E_t^{req})^{K_3}} \times \bar{\lambda}^{cons}, & \text{if } \Delta {}^i T_t^{dep} = 0, \end{cases} \quad (14)$$

261 where $\sum_{m=0}^{L_c-1} C_{t,m}^{elec}$ is the electricity cost in (7) and ${}^i C_t^{anx}$ is the “range anx-
 262 iety” cost, a penalty for not charging EV_i with its ${}^i E^{req}$ [25]. $\frac{{}^i E_t^{req}}{{}^i E^{req}}$ is the
 263 fraction of uncharged energy at time t . The coefficients K_1, K_2 and K_3 are
 264 hyperparameters and $\bar{\lambda}^{cons}$ is the average of the (dynamic) electricity con-
 265 sumption price profile of the current day. In this case, the objective of the
 266 RL algorithm is to minimize the charging cost and the uncharged energy
 267 fraction in each session.

268 The coefficient K_1 weighs the trade-off between the charging cost and
 269 the average amount of uncharged energy. The coefficients K_2 and K_3 are
 270 exponents that allow testing the performance with different types of the
 271 anxiety function. For example, the anxiety cost with $K_2 = 2$ and $K_3 = 2$
 272 is proportional to the square of the fraction of the uncharged energy, while
 273 with $K_2 = 1$ and $K_3 = 1$ the relationship is linear. Notice that, to compute
 274 the anxiety at time t , knowledge of ${}^i E_t^{req}$ and ${}^i T_t^{dep}$ is required. However,
 275 that knowledge is only required in the training phase of the RL algorithm
 276 when ${}^i E_t^{req}$ and ${}^i T_t^{dep}$ are readily available from historical data. During policy
 277 execution, computation of the reward is not required - since we will focus on
 278 a policy gradient RL algorithm. Thus, knowledge of ${}^i E_t^{req}$ and ${}^i T_t^{dep}$ is not
 279 required during policy execution.

280 For each connected EV, its state is updated as follows:

$${}^i \mathbf{X}_{t+1} = \begin{cases} (0, 0, 0), & \text{if battery full or } EV_i \text{ has departed,} \\ \left({}^i E_t^{ch,ses} + {}^i E_t^{ch}, {}^i P_{lim}, \Delta {}^i T_{arr} \right), & \text{otherwise,} \end{cases} \quad (15)$$

281 where ${}^iE_t^{ch}$ is the energy charged by EV $_i$ during slot t as defined in (9).
 282 Similar to the base MDP, new EVs that arrive between t and $t + 1$ are
 283 added to ${}^i\mathbf{X}_{t+1}$, and the remaining state parameters $t, \mathbf{P}_t^{pv}, \mathbf{P}_t^{fpv}, \boldsymbol{\lambda}_t^{cons}$ are
 284 also updated.

285 2.4. MDP with aggregated state-action space

286 To improve the scalability of the base MDP for larger fleet sizes, MDP $_{agg}$ is
 287 proposed. This MDP builds on the three-step method for smart EV charging
 288 proposed in [23]. This approach consists of three steps: (i) an aggregation
 289 step in which the individual EV charging constraints are established in the
 290 form of priorities and aggregated; (ii) an optimization step that uses the
 291 aggregated constraints to compute a collective charging plan for all the EVs,
 292 with the aim of maximizing self-consumption and minimising electricity costs;
 293 and (iii) a real-time control step dividing and dispatching the charging plan
 294 to all the EVs. In this aggregate MDP, the reward function is the same as
 295 that in Section 2.2, therefore only the state and action spaces are described
 296 below.

297 **State space**

The state s_t at time t is defined as:

$$s_t = (\mathbf{Z}_t, t, \mathbf{P}_t^{pv}, \mathbf{P}_t^{fpv}, \boldsymbol{\lambda}_t^{cons}), \quad (16)$$

298 where \mathbf{Z}_t , as shown in (17), is the aggregated fleet state, which is obtained
 299 through manual feature extraction consisting of the total fleet required energy
 300 and the total fleet maximum charging power with ${}^iP^{lim} = 0$ if station i is
 301 unused. The state parameters $t, \mathbf{P}_t^{pv}, \mathbf{P}_t^{fpv}, \boldsymbol{\lambda}_t^{cons}$ are the same as those in the
 302 base MDP.

$$\mathbf{Z}_t = \left(\sum_{i=1}^{N_{ev}} {}^iE_t^{req}, \sum_{i=1}^{N_{ev}} {}^iP^{lim} \right) \quad (17)$$

303 It is important to note that even though we do not consider the dynamics
 304 of the system, the aggregate fleet state can be updated as follows:

$$\mathbf{Z}_{t+1} = \left(\sum_{i=1}^{N_{ev}} {}^iE_t^{req} - {}^iE_t^{ch}, \sum_{i=1}^{N_{ev}} {}^iP^{lim} \right), \quad (18)$$

305 where ${}^iE_t^{ch}$ is the energy charged by EV $_i$ during slot t as defined in (9). New
 306 EVs that arrive between t and $t + 1$ are added to \mathbf{Z}_{t+1} , and the remaining

307 state parameters $t, \mathbf{P}_t^{pv}, \mathbf{P}_t^{fpv}, \boldsymbol{\lambda}_t^{cons}$ are also updated by reading the values
 308 from the database.

309 **Action space**

310 An action $a_t \in \mathcal{A}$ at time step t represents an aggregate charging power
 311 P^a for the entire fleet as shown in (19).

$$\forall a_t \in \mathcal{A}, a_t = P_t^a = \sum_{i=1}^{N_{ev}} {}^i P_t \quad (19)$$

312 This aggregate charging power is divided and dispatched to all the EVs
 313 using a heuristic dispatch described below.

314 *Heuristic dispatch*

315 The aggregate charging power is divided to all the EVs in the fleet using
 316 a heuristic dispatch based on (20). This dispatch also ensures that each
 317 EV leaves with its required energy charged and is not charged with a power
 318 greater than its maximum charging power.

$$\mathbf{P}_t^b = ({}^1 P_t^b, \dots, {}^{N_{ev}} P_t^b) \quad (20)$$

$$\text{s.t. } {}^i P_t^{b,min} \leq {}^i P_t^b \leq {}^i P_t^{b,max}, \quad \forall i \in [1, N_{ev}] \quad (21)$$

$${}^i P_t^{b,min} = \max \left(0, ({}^i E_t^{req} - (\Delta {}^i T_t^{dep} - 1) \times {}^i P^{lim}) \times \frac{60}{{}^i M_t^{end} - {}^i M_t^{start}} \right), \quad (22)$$

$$\text{with } {}^i P_t^{b,max} = \min \left({}^i P^{lim}, {}^i E_t^{req} \times \frac{60}{{}^i M_t^{end} - {}^i M_t^{start}} \right) \quad (23)$$

319 ${}^i P_t^{b,min}$ is the minimum charging power required to guarantee that EV_i
 320 leaves with its battery fully charged, i.e. to guarantee ${}^i E_t^{req} = 0$ when
 321 $\Delta {}^i T_t^{dep} = 0$. Overcharging - charging above maximum charging power -
 322 is prevented by the maximum charging power ${}^i P_t^{b,max}$, which is limited by
 323 ${}^i P^{lim}$ and ${}^i E_t^{req}$. ${}^i M_t^{start}$ and ${}^i M_t^{end}$ are the first minute and last minute re-
 324 spectively of EV_i in time step t where $0 \leq {}^i M_t^{start} < L_c$ and $0 < {}^i M_t^{end} \leq L_c$
 325 and are computed as follows:

$${}^iM_t^{start} = \begin{cases} \lfloor ({}^iT^{arr} - t) \times L_c \rfloor, & \text{if } t \leq {}^iT^{arr} < t + 1, \\ 0, & \text{otherwise.} \end{cases} \quad (24)$$

$${}^iM_t^{end} = \begin{cases} \lceil ({}^iT^{dep} - t) \times L_c \rceil, & \text{if } t \leq {}^iT^{dep} < t + 1, \\ L_c, & \text{otherwise.} \end{cases} \quad (25)$$

326 The operation of this heuristic dispatch can be described as “least-flexible
 327 first” scheduling. EVs are assigned partitions of the aggregate charging power
 328 P^a in order of priority. EVs with a lower flexibility are given a higher priority.
 329 The flexibility ${}^i\tau$ of each EV_i is calculated according to (26).

$${}^i\tau := \Delta {}^iT^{dep} - 1 - \frac{{}^iE^{req} - {}^iE^{ch}}{{}^iP^{lim}} \times \frac{60}{L_c}. \quad (26)$$

330 This priority represents the number of time steps until the next time
 331 step $t+1+{}^i\tau$ at which ${}^iP_{t+1+{}^i\tau}^{b,min} > 0$ (to ensure EV_i leaves with ${}^iE_t^{req} = 0$),
 332 assuming EV_i charges ${}^iE^{ch}$ at time step t . To illustrate this priority computa-
 333 tion for $L_c = 60$ minutes, consider an EV with $\Delta {}^iT^{dep} = 8$ time steps, ${}^iE^{req} =$
 334 10 kWh, ${}^iE^{ch} = 0$ kWh, ${}^iP^{lim} = 5$ kW that needs exactly 2 time steps to fully
 335 charge at a power of ${}^iP^{lim}$. The earliest time step when its ${}^iP_{t+1+{}^i\tau}^{b,min} > 0$ is
 336 thus ${}^i\tau = 8 - 1 - 2 = 5$ slots from $t+1$. The dispatch algorithm assigns par-
 337 titions of P_t^a heuristically with the aim to maximize the minimum value of
 338 ${}^i\tau$:

$$\text{maximize } \min_i {}^i\tau. \quad (27)$$

339 **Reward function**

340 The reward function is the same as that of MDP_{base} .

341
 342 It is important to note that even though the MDPs described above do
 343 not explicitly take into account the system constraints, which come in the
 344 form of ensuring that the EV is fully charged before its departure and that
 345 the EV is not charged at a power greater than its maximum charging power,
 346 this work proposes using a backup controller to ensure these constraints are
 347 respected. This backup controller is described in the next section.

348 **3. Algorithms**

349 To solve the EV fleet charging problem described in Section 2, a policy
 350 gradient RL algorithm, proximal policy optimization (PPO) [26], is used.
 351 The goal of the learning agent is to find a (parameterised) control policy π
 352 (i.e. a mapping from a given or perceived state to the action that has to
 353 be taken in that state, $\pi : \mathcal{S} \rightarrow \mathcal{A}$), which maximises the return over the
 354 optimization horizon from some initial state s_0 as shown in (28).

$$G^\pi(s_0) = \sum_{t=1}^{T-1} \gamma^t \rho(s, \pi(s)), \quad (28)$$

355 where $\gamma \in [0; 1]$ is a discount factor that takes into account the uncertainty
 356 in the future reward and T is the length of the finite optimisation horizon.
 357 This return is the discounted cumulative reward along a trajectory generated
 358 by the policy.

359 *3.1. Proximal policy optimization*

360 Proximal policy optimization is a policy gradient RL algorithm based
 361 on the actor-critic algorithm [11] that directly optimises a parameterised
 362 and differentiable policy. The policy must be differentiable with respect to
 363 its parameters to allow computation of the gradient required for the policy
 364 parameter updates. Typically, the policy is represented by a neural network
 365 and is expressed as follows:

$$\pi(a|s, \theta) = Pr(a|s, \theta), \quad (29)$$

366 where θ represents the weights of the neural network. The goal is therefore
 367 to find the values of θ that maximise the return G .

368 The algorithm uses the clipped surrogate objective - the probability ratio
 369 pr of the old policy and new policy - with the aim of providing a more stable
 370 update of the policy parameters [26]. The unconstrained objective function
 371 that PPO aims to maximize is as shown in (30).

$$L^{ORIG}(\theta) = \mathbb{E}_t \left[\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} \hat{\Lambda}_t \right] = \mathbb{E} \left[pr_t(\theta) \hat{A}_t \right], \quad (30)$$

372 where $pr_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$ and $\hat{\Lambda}_t$ is an estimate of the advantage ($\Lambda(s_t, a_t) =$
 373 $Q(s, a) - V(s)$), $Q(s, a)$ is the state-action value function and $V(s)$ is the state

374 value function). $L^{ORIG}(\theta)$ can still lead to large gradient updates and con-
 375 sequently, instability during learning. This is remedied by using the clipped
 376 surrogate objective shown in (31).

$$L_t^{CLIP}(\theta) = \mathbb{E}_t \left[\min \left(\rho_t(\theta) \hat{\Lambda}_t, \text{clip}(\rho_t(\theta), 1 - \delta, 1 + \delta) \hat{\Lambda}_t \right) \right], \quad (31)$$

377 where $\delta \in [0.1, 0.3]$ according to [26].

378 The clipped surrogate objective is further augmented for applying it to
 379 a neural network architecture with shared parameters for representing the
 380 policy and value functions⁷. Typically, the policy and value network share
 381 the first few hidden layers, which perform feature extraction of the state
 382 space. Additionally, an entropy term is included to the objective to increase
 383 exploration and as such more coverage of the state space. The new objective
 384 is as shown in (32).

$$L_t^{CLIP2}(\theta) = \mathbb{E}_t \left[L_t^{CLIP}(\theta) - c_1 (V_\theta(s_t) - V_t^{\text{targ}})^2 + c_2 H[\pi_\theta](s_t) \right], \quad (32)$$

385 where c_1 and c_2 are hyper-parameters and H is an entropy measure. The
 386 resulting PPO algorithm is described in Algorithm 1. The algorithm par-
 387 allelises the sampling of the agent-environment interactions - by using N
 388 parallel actors - and uses multiple epochs of stochastic gradient ascent per
 389 policy update. Parallel sampling significantly speeds up training times by us-
 390 ing parallel processors while the ability to use multiple epochs when updating
 391 the neural network increases sample efficiency.

392 3.2. Backup controller

393 The backup controller is an overrule mechanism that ensures that the
 394 system constraints are respected. Recall that in the context of this paper,
 395 these constraints are the charging power limits of the EV, and the need to
 396 fully charge the EV before its departure. The backup controller therefore
 397 clips ${}^i P_t$ - the charging power or action suggested by the RL agent - at each
 398 time step t for each EV_i according to (33).

⁷The value function is represented by a neural network for function approximation to make it more generalisable over unseen state(-action pairs).

Algorithm 1: Proximal policy optimization [26]

Input : policy parameters θ_0 , clipping threshold δ , initial network parameters θ_0

1 **for** $k := 0, 1, \dots$ **do**

2 **for** $actor := 1, \dots, N$ **do**

3 Run policy π_k in environment for T time steps;

4 Compute rewards \hat{r}_t ;

5 Compute advantage estimates $\hat{\Lambda}_t$;

6 Compute policy update $\theta_{k+1} := \arg \max_{\theta} L^{CLIP2}(\theta)$

7 with K epochs of mini-batch stochastic gradient ascent, where

8 $L^{CLIP2}(\theta) := \mathbb{E}_{\tau \sim \pi_k} [\sum_{t=0}^T L_t^{CLIP2}(\theta)]$;

$${}^iP_t^b = \begin{cases} {}^iP_t^{b,min}, & \text{if } {}^iP_t \leq {}^iP_t^{b,min}, \\ {}^iP_t, & \text{if } {}^iP_t^{b,min} < {}^iP_t \leq {}^iP_t^{b,max}, \\ {}^iP_t^{b,max}, & \text{if } {}^iP_t > {}^iP_t^{b,max}. \end{cases} \quad (33)$$

399 Recall that ${}^iP_t^{b,min}$ is the minimum charging power required to guaran-
400 tee that EV_i leaves with the required energy charged, i.e. ${}^iE_t^{req} = 0$ when
401 $\Delta {}^iT_t^{dep} = 0$, as shown in (22), and ${}^iP_t^{b,max}$ is limited by ${}^iP^{lim}$ and ${}^iE_t^{req}$, as
402 shown in (23).

403 It is important to note that when EV_i , arrives in control slot t ($t \leq$
404 ${}^iT^{arr} < t + 1$), it is entirely controlled by the backup controller starting from
405 ${}^iT^{arr}$ up to $t + 1$. No action is taken by the agent (${}^iP = 0$), and ${}^iM_t^{start}$
406 is set to $\lfloor ({}^iT^{arr} - t) \times L_c \rfloor$ (24). The backup controller determines ${}^iP^b$ by
407 clipping the ${}^iP = 0$ using (33) to ensure EV_i can still charge its required
408 energy during the remainder of the session starting from $t + 1$. Therefore,
409 the backup controller ensures that the charging power for EV does not exceed
410 ${}^iP^{lim}$ and that each EV_i is charged with exactly its ${}^iE^{req}$ at departure.

411 3.3. Real-time controller

412 The output of the backup controller, ${}^iP_t^b$ for $i = 1 \dots N_{ev}$, is a charging
413 power for each EV for control time step t . On cloudy days the PV electricity
414 generation can fluctuate rapidly during the entire time step. As a result, even
415 when the total charging power in a control time step, $\sum_{i=1}^{N_{ev}} {}^iP_t^b$, is equal to
416 the mean PV generation in the same control time step, $\text{av}(P^{pv}, t, t + 1)$, the

417 total charging cost for that slot C_t^{elec} is not zero because the PV generation
 418 fluctuations incur extra electricity consumption and injection costs. Hence,
 419 an algorithm that can learn on a fine timescale and can adapt the charging
 420 power to the rapidly fluctuating PV generation has a lower bound on the
 421 minimum charging cost it can achieve than an algorithm that learns on a
 422 lower time resolution. One solution is to learn on a fine timescale by reducing
 423 L_c , which results in charging decisions being taken more frequently. However,
 424 the time to simulate each episode is inversely proportional to L_c . A decrease
 425 in L_c thus leads to an increase in training time.

426 Instead of decreasing L_c , we propose a hybrid solution in which the con-
 427 trol period L_c is set relatively large (e.g. 60 minutes) and the RL algorithm
 428 is combined with a rule-based real-time controller (RTC) that dynamically
 429 adapts the charging powers from the low resolution RL algorithm to the cur-
 430 rent real-time PV generation on a higher resolution of one minute. The RTC
 431 takes the output of the backup controller, ${}^iP_t^b$ for $i = 1 \dots N_{ev}$, and for each
 432 minute m of time step t determines a real-time charging power ${}^iP_{t,m}^r$ for each
 433 EV_i by solving the following optimization problem:

$$\begin{aligned}
 \min \quad & \sum_{m=0}^{L_c-1} \left| P_{t,m}^{pv} - \sum_{i=1}^{N_{ev}} {}^iP_{t,m}^r \right|, \\
 \text{s.t.} \quad & \sum_{m=0}^{L_c-1} {}^iP_{t,m}^r / L_c = {}^iP_t^b \quad \forall i \in [1, N_{ev}], \\
 & {}^iP_{t,m}^r \leq {}^iP^{lim} \quad \forall i \in [1, N_{ev}], \quad \forall m \in [0, L_c - 1].
 \end{aligned} \tag{34}$$

434 By solving the above optimization problem through a set of manually
 435 defined rules, the RTC computes the power schedule for each EV_i for each
 436 minute m throughout the duration of time step t according to (35).

$${}^iP_t^r = ({}^iP_{t,m}^r = {}^iP_t^b \quad \text{for } {}^iM_t^{start} \leq m < {}^iM_t^{end}). \tag{35}$$

437 This ensures that the total charging power stays as close as possible to the
 438 PV generation in each minute of the time step while ensuring that the total
 439 energy for charging the EV during that time step is equal to the energy
 440 suggested by the backup controller. The RTC also ensures that the charging
 441 power for each EV in each minute m does not exceed the absolute maximum
 442 charging power ${}^iP^{lim}$ for that EV.

443 The RTC algorithm is described in Algorithm 2. For readability, the
 444 subscript t is dropped for most variables in the algorithm. For each EV to

445 be charged, ${}^iP^{r,min}$, ${}^iP^{r,max}$ and ${}^iE^{rem}$ are calculated. The variables ${}^iP^{r,min}$
 446 and ${}^iP^{r,max}$ are minimum and maximum charging powers that apply to each
 447 minute m during time step t , while ${}^iE^{rem}$ is the remaining energy to be
 448 charged in the current control time step as specified by ${}^iP^b$ (the backup
 449 controller output). ${}^iP^{r,min}$ and ${}^iP^{r,max}$ are limited by the charging flexibility
 450 factor κ , a hyper-parameter that dictates by how much the charging power
 451 in each minute can deviate from ${}^iP^b$. A value of $\kappa = 1$ is equivalent to not
 452 using the RTC.

453 Figure 3 illustrates the interactions between the different algorithms in
 454 the proposed control scheme.

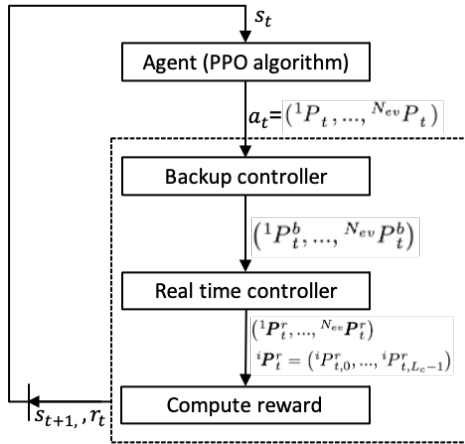


Figure 3: An illustration of the interaction between the different algorithms. The agent uses the PPO algorithm for action selection. Everything outside the agent is considered as its environment.

Algorithm 2: Real-Time Controller

Input : $(iP^b, iM^{start}, iM^{end}$ for $i = 1 \dots N_{ev}$),
 P_m^{pv} for $m = 0 \dots L_c - 1$
Parameter: Charging flexibility κ
Output : $iP_t^r = (iP_{t,0}^r, \dots, iP_{t,L_c-1}^r)$

- 1 EVsToCharge := set of all i where $iP^b > 0$;
- 2 **for** i *in* EVsToCharge **do**
- 3 $iPr,max := \min(iP^{lim}, iP^b \times \kappa)$;
- 4 $iPr,min := iP^b / \kappa$;
- 5 $iErem := iP^b \times (iM^{end} - iM^{start}) / 60$;
- 6 **for** $m := 0 \dots L_c - 1$ **do**
- 7 EVsThisMinute := set of all i where $iM^{start} \leq m < iM^{end}$;
- 8 **for** i *in* EVsToCharge \cap EVsThisMinute **do**
- 9 $iPr,max_m := \min(iErem \times 60, iPr,max)$;
- 10 $iPr,min_m :=$
 clip $(iErem \times 60 - (iM^{end} - m - 1) \times iPr,max, iPr,min, iPr,max_m)$;
- 11 $P_m^{min,\Sigma} := \sum_{i=1}^{N_{ev}} iPr,min_m$;
- 12 $P_m^{max,\Sigma} := \sum_{i=1}^{N_{ev}} iPr,max_m$;
- 13 **for** i *in* EVsToCharge \cap EVsThisMinute **do**
- 14 **if** $P_m^{max,\Sigma} - P_m^{min,\Sigma} == 0$ **then**
- 15 $iPr_{t,m} := P_m^{min}$;
- 16 **else**
- 17 $iPr_{t,m} :=$
 clip $(iPr_{t,m}^{min} + \frac{P_m^{pv} - P_m^{min,\Sigma}}{P_m^{max,\Sigma} - P_m^{min,\Sigma}} (iPr_{t,m}^{max} - iPr_{t,m}^{min}), iPr_{t,m}^{min}, iPr_{t,m}^{max})$;
- 18 $iErem := iErem - iPr_{t,m} / 60$;
- 19 **return** iP_t^r ;

455 4. Case study and simulation results

456 The proposed framework is evaluated using the EnergyVille office building
457 as a case study. The various MDPs are evaluated on (i) self-consumption of
458 local PV-generated electricity and net electricity costs, and (ii) scalability in
459 terms of fleet size.

460 4.1. Case study

461 This work considers real-world data sets on energy consumption and gener-
462 ation at EnergyVille, complemented with data from Belgian electricity mar-
463 kets and regional PV forecasts. More specifically:

- 464 • PV generation measurements $P^{pv,tot}$ of the EnergyVille rooftop PV
465 installation measured at 5 minute intervals
- 466 • PV forecast with a 15 minute time step for the province of Limburg,
467 Belgium (location of EnergyVille)⁸
- 468 • EV charging transactions at EnergyVille: 586 valid historical charging
469 sessions across 171 days collected between 08/08/2018 and 20/09/2019.
470 For each EV transaction, T_{arr} is known when the EV is plugged in, P^{lim}
471 is detected during the first few minutes of charging, T^{dep} and E^{req} are
472 extracted from the historical data set.
- 473 • EV power consumption ${}^iP^{hist}$ of each historical charging session at
474 EnergyVille, measured every 20 seconds
- 475 • Historical grid consumption P^{cons} and injection P^{inj} ⁹ of the EnergyVille
476 building, collected every 15 minutes
- 477 • A grid injection tariff of $\lambda^{inj} = 1.46\text{€}/\text{MWh}$ as the rooftop PV instal-
478 lation at EnergyVille is larger than 10kVA.
- 479 • A TOU grid consumption price¹⁰ $\lambda^{Belpeax}$ with a one hour periodicity is
480 used to compute the grid consumption prices λ^{cons} as shown in (36).

⁸Solar-PV power forecasting for Belgium: <https://www.elia.be/en/grid-data/power-generation/solar-pv-power-generation-data>

⁹ P^{inj} is the surplus PV generation injected to the grid.

¹⁰Belgian day-ahead market prices: <https://transparency.entsoe.eu/transmission-domain/r2/dayAheadPrices/show>

$$\lambda^{cons} = (\lambda^{Belpex} + 0.045 \text{ €/kWh}) \times 1.21, \quad (36)$$

481 where 0.045 €/kWh are the estimated grid tariffs, and a 21% VAT
 482 is charged. By including taxes in λ^{cons} an estimate of the actual cost
 483 savings for EnergyVille is obtained in the experiments.

484 It is worth noting that in our simulations a scaled version of the PV
 485 generation, $P^{pv} = P^{pv,scaled} \equiv 0.2 \times P^{pv,tot}$ is used to filter out the influence
 486 of the stochastic electricity consumption of the rest of the building. This
 487 scaling expresses a hypothetical scenario where a small PV installation is
 488 available solely for EV charging.

489 Also, 98% of the charging transactions in the data set occur between 7:00
 490 and 20:00. Therefore, 7:00 and 20:00 are set as the start and end times of each
 491 episode respectively. The algorithms are tested using three control time steps:
 492 $L_c = 5$, $L_c = 15$ and $L_c = 60$ minutes. These control time steps have been
 493 selected by considering the 15 minutes time step for the imbalance electricity
 494 market, the 60 minutes time step of the day-ahead electricity market, and in
 495 order to approach real-time operation, a time step of 5 minutes.

496 To evaluate the performance of the proposed control framework the sim-
 497 ulation results are compared with those from the business as usual (BAU)
 498 and “perfect information optimum” (PIO) strategies. Recall that the BAU
 499 strategy is equivalent to passive charging, where each EV is charged imme-
 500 diately when it is plugged in to the charging station at its maximum power
 501 iP^{lim} until the required energy iE_t^{req} reaches zero. The PIO strategy as-
 502 sumes complete knowledge for the entire day of all EV arrival and departure
 503 times, required energy, maximum power and the PV generation and electric-
 504 ity prices. The problem is formulated as a constrained nonlinear optimisa-
 505 tion problem as shown in (37) and solved using the sequential least-squares
 506 quadratic programming algorithm [24, 27]. Due to limited computational
 507 resources, $L_c = 15$ minutes is used. The BAU and PIO are selected as the
 508 baselines because; (i) the BAU is the strategy that is used in most charg-
 509 ing stations, and (ii) the PIO provides a theoretical baseline considering a
 510 scenario in which all the information on the different system variables is
 511 available. These baselines provide a best- and worst-case comparison.

$$\begin{aligned}
\min_{\mathbf{P}^r} C^{day} &= \sum_{t=0}^{T-1} C_t^{elec} \\
\text{s.t. } C_t^{elec} &= \begin{cases} C_t^{cons} = \sum_{i=1}^{N_{EV}} (iP_t^r - P_t^{pv}) \times (L_c/60) \times \lambda_t^{cons}, & \text{if } \sum_{i=1}^{N_{EV}} iP_t^r - P_t^{pv} \geq 0, \\ C_t^{inj} = \sum_{i=1}^{N_{EV}} (P_t^{pv} - iP_t^r) \times (L_c/60) \times \lambda_t^{inj}, & \text{otherwise,} \end{cases} \\
0 \leq iP_t^r &\leq iP_t^{lim} \quad \forall i \in [1, N_{ev}], \forall t \in [0, T-1], \\
\sum_{i=1}^{N_{ev}} \sum_{t=iT^{arr}}^{iT^{dep}} iP_t^r \times (L_c/60) &= iE^{req}.
\end{aligned} \tag{37}$$

512

513 The neural network used to represent the actor-critic in the PPO algo-
514 rithm consists of a first hidden layer with 128 nodes, shared between the
515 actor and the critic. Both networks contain two hidden layers with 64 nodes
516 each. The number of layers and nodes are obtained based on [22] and [26].
517 The `tanh` activation function is used. A representation of the actor-critic
518 network is shown in Fig. 4.

519 4.2. Simulation results

520 The performance of the proposed control strategy is evaluated by consid-
521 ering two simulation experiments. The first experiment evaluates the perfor-
522 mance of the control strategy for the different MDP formulations while the
523 second experiment investigates the scalability of the control framework. The
524 net electricity cost per day (an optimisation horizon of one day is used) is
525 considered as the key performance indicator.

526 4.2.1. Experiment 1: performance evaluation of MDPs

527 This experiment compares the performance of the control framework for
528 the three MDP formulations and investigates the influence of the time step
529 L_c and the RTC on the performance. The MDPs are tested for $L_c=5$, $L_c=15$
530 and $L_c=60$ minutes, and each of these instances is tested without the RTC
531 and with the RTC (for $\kappa = 1.5$ and with $\kappa = 5.0$). For each instance, the
532 training/testing loop is executed for 5×10^7 agent-environment interactions.
533 The remaining MDP hyper-parameters are set as follows: $N_{past} = N_{fut} = 2$,
534 $K_1 = 5.0 \times 10^4$, $K_2 = 1$, $K_3 = 1$ (MDP_{hid}).

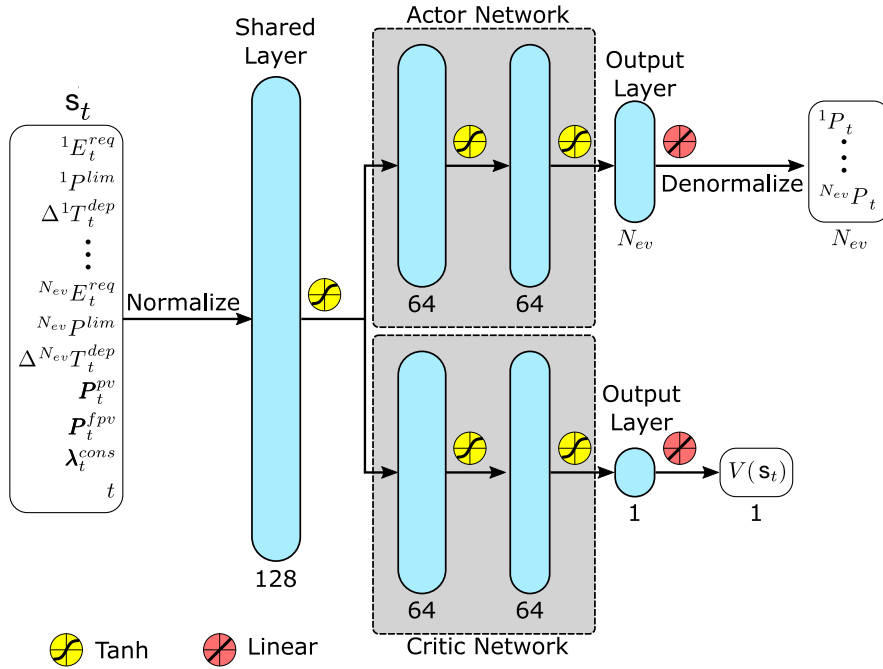
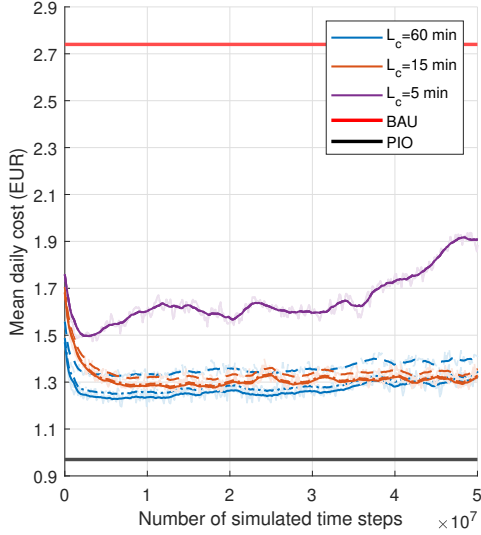
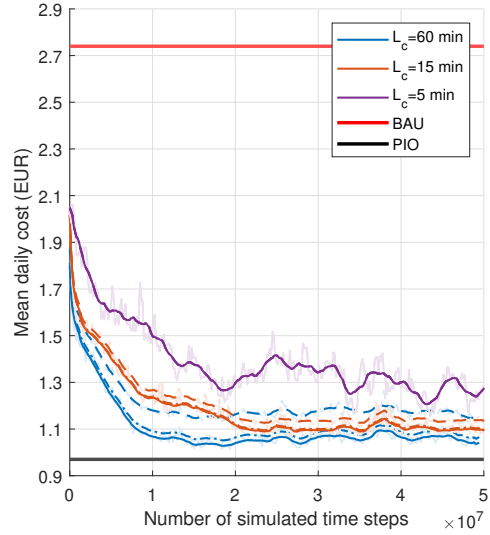


Figure 4: Illustration of the actor-critic network. This network is used for action selection - the policy by the actor - and for estimating the value function - critic.

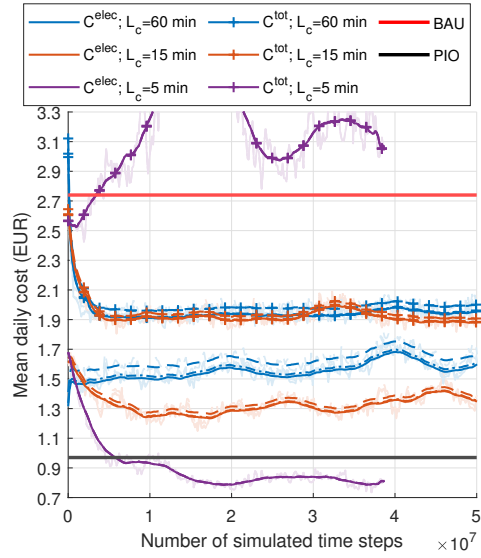
535 The simulation results are presented in Fig. 5. These results are expressed
536 in terms of the total cost $C^{tot} = C^{cons} + C^{inj} + C^{anx}$ and the electricity cost
537 $C^{elec} = C^{cons} + C^{inj}$. For MDP_{base} and MDP_{agg} , $C^{anx} = 0$. MDP_{agg} obtains the
538 lowest electricity cost and has the most stable learning curve while MDP_{base}
539 tends to converge quickly to a sub-optimal local minimum. The resulting
540 electricity cost for MDP_{agg} is 0.2€ lower than that of MDP_{base} (for $L_c =$
541 60). For MDP_{hid} , the trade-off between electricity cost and the fraction of
542 uncharged energy (Fig. 5c and 5d) depends heavily on L_c . When $L_c = 60$
543 minutes, there is a 0.2€ increase in electricity cost compared to MDP_{base} , and
544 2% of the EVs leave with 25% of their required energy not met. It is worth
545 noting that the EVs leave fully charged in the case of MDP_{base} and MDP_{agg}
546 as knowledge of the required energy and departure times of the EVs allows
547 the agent to obtain (near) optimal schedules. Moreover, this knowledge of
548 required energy and departure times allows the backup controller to override
549 the actions of the learning agent and ensure that the EVs are fully charged
550 before departure.



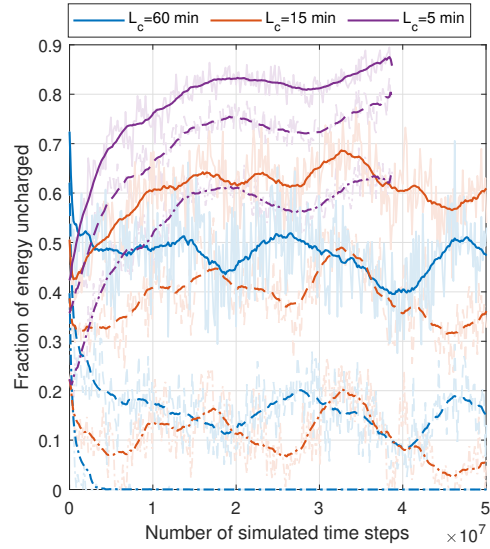
(a) Evolution of C^{elec} for MDP_{base} (solid = RTC with $\kappa = 5.0$, dash-dotted = RTC with $\kappa = 1.5$, dash = no RTC)



(b) Evolution of C^{elec} for MDP_{agg} (solid = RTC with $\kappa = 5.0$, dash-dotted = RTC with $\kappa = 1.5$, dash = no RTC)



(c) Evolution of C^{elec} and C^{tot} for MDP_{hid} (solid = RTC with $\kappa = 5.0$, dash-dotted = RTC with $\kappa = 1.5$, dash = no RTC)



(d) Evolution of 99- (solid), 98- (dashed) and 96- percentile (dash-dotted) of fraction of energy uncharged

Figure 5: Learning curves for the main experiment comparing the three MDPs, three values of L_c and investigating the impact of the RTC on performance

551 When comparing the influence of the value of L_c when no RTC is used,
 552 for $L_c = 15$ minutes MDP_{hid} obtains the best performance as it shifts charg-
 553 ing priority towards minimizing electricity cost rather than minimizing the
 554 fraction of uncharged energy as previously mentioned. For $L_c=60$ minutes
 555 with no RTC, the algorithm performs worse compared to the case for $L_c=15$
 556 minutes with or without the RTC for MDP_{base} and MDP_{agg}. This is due to
 557 the length of the control time step, which makes the control agent unable
 558 to learn the fluctuations in the PV electricity generation. Recall that, using
 559 the RTC allows to mitigate against the rapidly fluctuating PV generation
 560 especially when learning at a low time resolution. At $L_c = 5$ minutes, the
 561 high time resolution results in unstable learning for MDP_{base} and MDP_{hid}.
 562 In MDP_{hid}, the total cost and the fraction of uncharged energy increase dur-
 563 ing training, while the electricity cost drops below that of the PIO. The RL
 564 algorithm hence has lost the ability to learn an effective policy. On the other
 565 hand, the lower dimension of the state-action space in MDP_{agg} allows the
 566 RL algorithm to maintain its ability to learn (albeit slower than $L_c = 15$) an
 567 effective policy at $L_c = 5$ minutes, obtaining an electricity cost that is 0.1€
 568 higher compared to that obtained when $L_c = 15$ minutes. It may be possible
 569 that the electricity cost decreases further for MDP_{agg} and $L_c = 5$, possibly
 570 even dropping below that obtained for $L_c = 15$, after the measured 5×10^7
 571 iterations. However, the CPU time required for testing this hypothesis would
 572 be impractical, with 5×10^7 iterations already requiring ≈ 24 hours.

573 The influence of the RTC is most noticeable at $L_c = 60$ minutes, resulting
 574 in a decrease in the electricity cost by on average 0.08€ ($\kappa = 1.5$) and 0.10€
 575 ($\kappa = 5.0$) compared to the MDPs with $L_c = 60$ minutes without RTC. For
 576 $L_c = 15$ and $L_c = 5$ minutes the performance gain from the RTC is lower since
 577 those instances already learn on a high time resolution. The lowest electricity
 578 cost, $C^{elec} = 1.03\text{€}$, is obtained by MDP_{agg} with RTC and $\kappa = 5.0$. This
 579 value is only 0.05€ above the electricity cost obtained by the PIO. Moreover,
 580 MDP_{agg} with $L_c = 60$, with RTC and $\kappa = 5.0$ converges approximately two
 581 times faster and obtains an electricity cost of 0.1€ lower than MDP_{agg} with
 582 $L_c = 15$ and without the RTC, which is the second best performing instance.

583 Figures 6 to 8 show examples of charging schedules for three sample
 584 days in the test set: a sunny day (Fig. 6), a day with variable sunshine
 585 (Fig. 7), and an overcast day (Fig. 8). The figures clearly show how the
 586 proposed framework moves the charging of EVs to later moments in the day
 587 when more PV-generated power is available and grid consumption prices are
 588 typically lower. This is a substantial improvement compared to the worst-

589 case scenarios of the BAU strategy. The figures also show that, because of
 590 its prior knowledge, the PIO strategy is able to spread the actual charging
 591 over the entire time the EV is plugged in, resulting in less volatile charging
 592 power levels. It however does not always find the most optimal solution, as
 593 shown in Figure 7, where it is outperformed by the MDP_{agg} that is able to
 594 avoid a small fraction of charging from the grid in the morning.

595 The behavior of the trained MDP instances on the sunny day is very
 596 similar, with all instances deferring charging to the middle of the day. With
 597 E^{req} and T^{dep} unknown to the agent in MDP_{hid} , it favours charging the EVs
 598 sooner - at a higher electricity cost - compared to MDP_{base} to avoid an EV
 599 leaving with its battery not fully charged. However, as shown on the figure,
 600 EV_7 (gray) still leaves with an uncharged energy fraction of 0.117 due to its
 601 unusually short session around noon. The difference between MDP_{base} and
 602 MDP_{agg} is the order in which they charge the EVs. The heuristic dispatch
 603 in MDP_{agg} prioritizes charging EV_2 (orange) in the morning since it has the
 604 lowest flexibility (due to its early departure time). In contrast, MDP_{base} does
 605 not exhibit any logical charging priority.

606 On the test day with variable sunshine (Fig. 7), the benefit of using a
 607 smaller L_c or using the RTC is clearly visible. MDP_{agg} with $L_c = 60$ and
 608 with the RTC obtains a perfect schedule ($C^{cons} = 0$) by being able to learn
 609 an effective policy on a broad timescale and adapting the charging power to
 610 the rapidly varying PV generation using the RTC. MDP_{agg} without the RTC
 611 and $L_c = 5$ obtains a near perfect schedule. MDP_{hid} fails to fully charge
 612 an EV with an unusually short session occurring around noon; EV_2 (yellow)
 613 leaves with an uncharged energy fraction of 0.212.

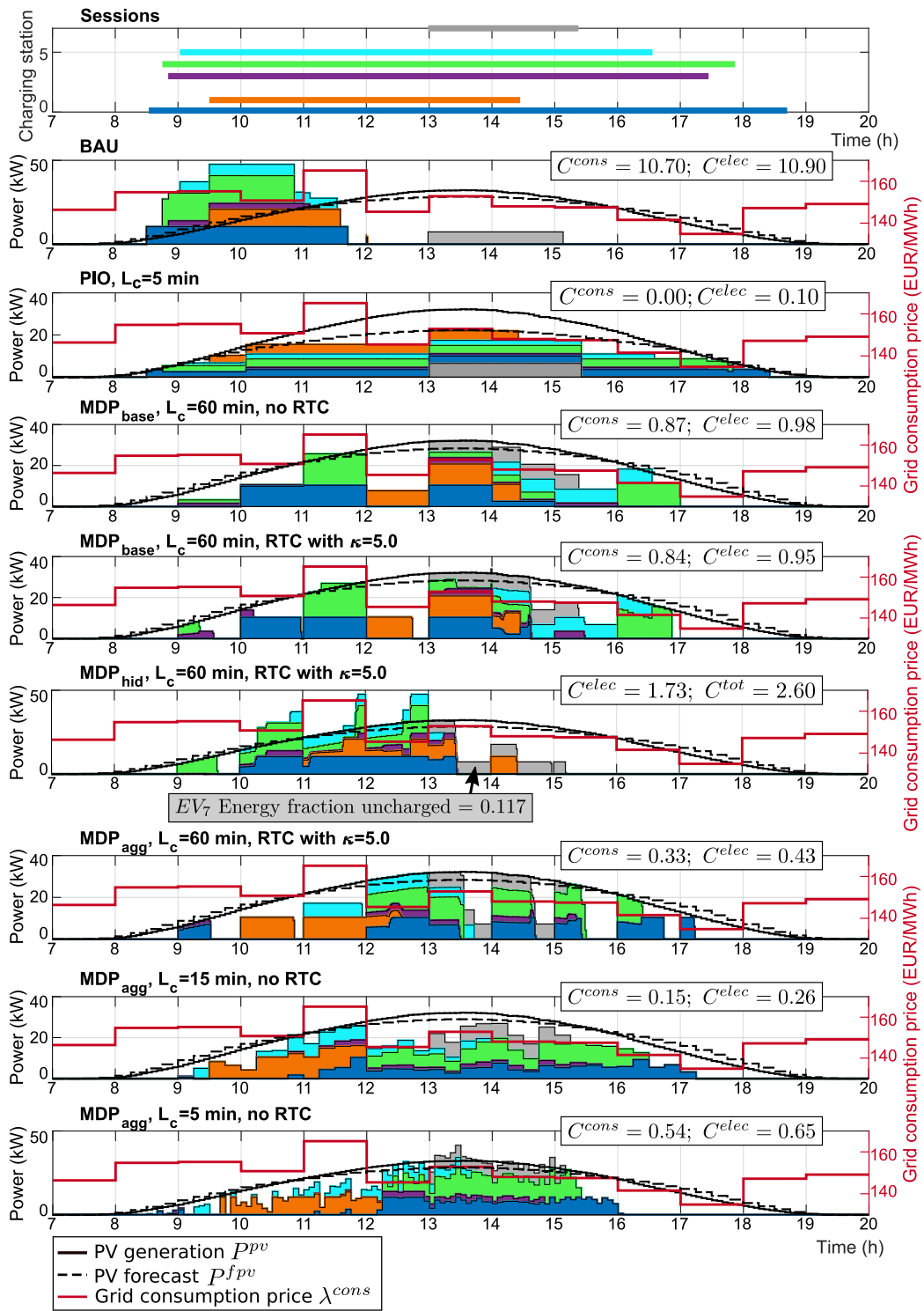


Figure 6: Timeline for a sunny test day for the BAU strategy and several instances of the RL algorithm

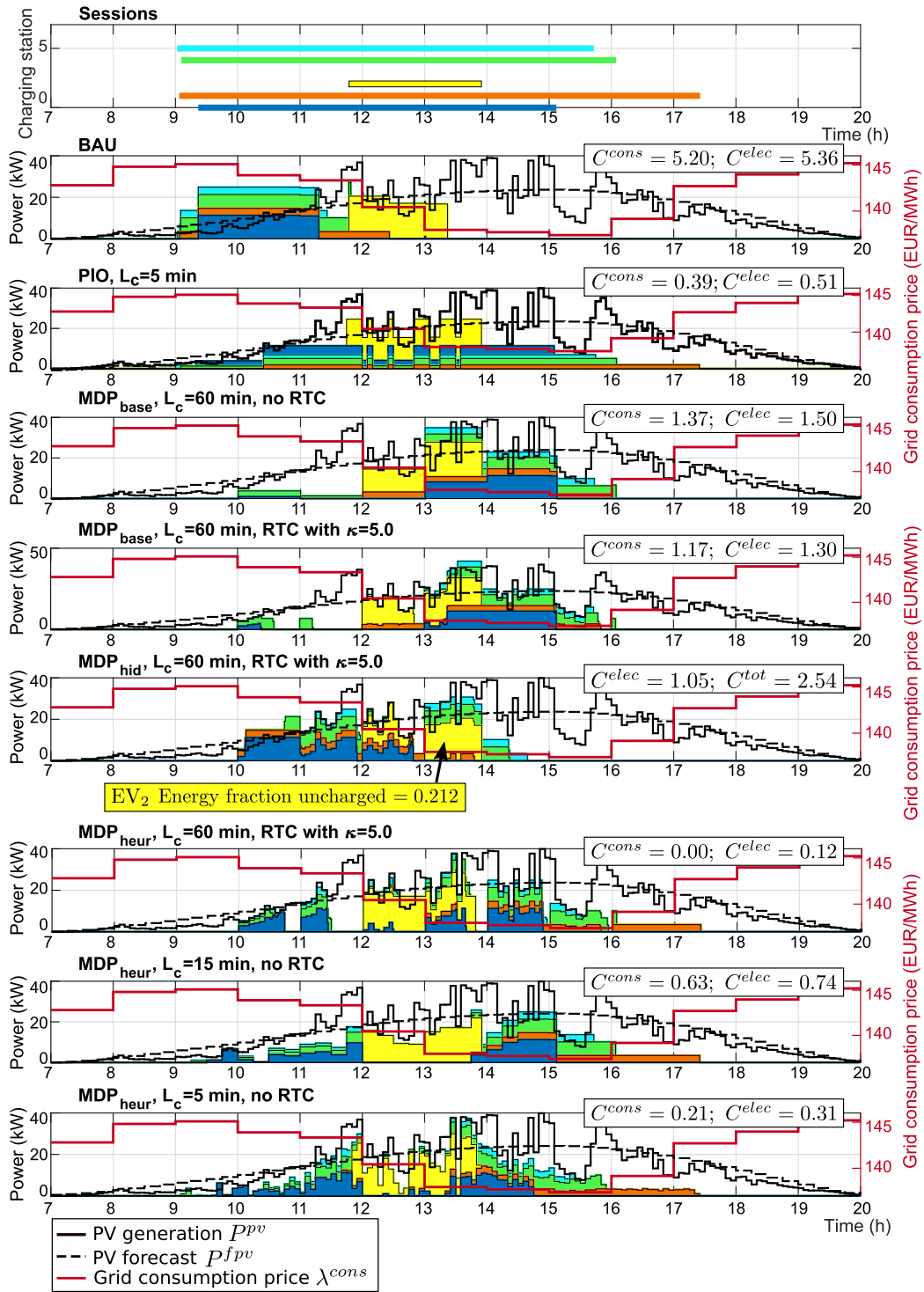


Figure 7: Timeline for a variable sunshine test day for the BAU strategy and several instances of the RL algorithm

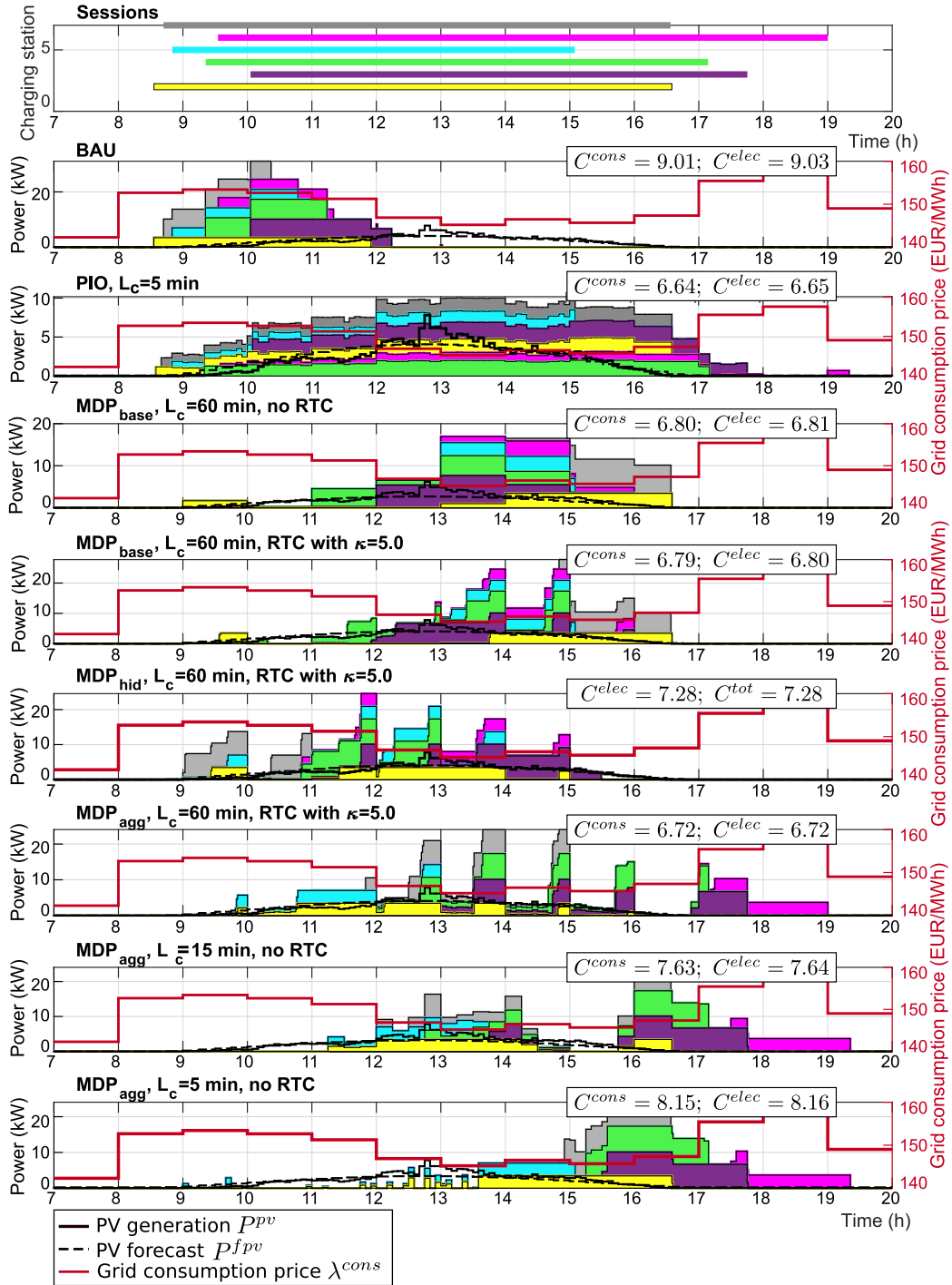


Figure 8: Timeline for an overcast test day for the BAU strategy and several instances of the RL algorithm

614 Finally, the optimal charging strategy on an overcast day (Fig. 8) would
615 be to self-consume all of the PV generation and charge the remaining required
616 energy during periods of low electricity prices. MDP_{agg} with $L_c = 60$ and
617 with the RTC obtains the lowest electricity cost, self-consuming 100% of
618 the PV-generation ($C^{inj} = 0$) and showing a trend towards charging the
619 remaining energy during low λ^{cons} . The major downside of the RTC is that
620 it is more volatile to EV charging power as can be seen in the figure. This
621 volatility can be lowered by using a smaller value of κ , thereby trading the
622 electricity cost for a slightly lower battery degradation. MDP_{hid} prioritizes
623 charging earlier and this time manages to charge all EVs with their respective
624 required energy. The instances with $L_c = 15$ and $L_c = 5$ minutes obtain the
625 highest electricity costs, struggling to follow the PV generation and relying
626 on the backup controller to charge the EVs at the end of their session.

627 The results from this experiment clearly show that the PPO algorithm
628 can achieve an effective charging policy, and a careful design of the MDPs
629 makes a significant difference in the performance of the algorithm. Com-
630 bining the three-step method in MDP_{agg} and the RTC, results in a signif-
631 icantly lower electricity cost (0.2€ lower) compared to the straightforward
632 design MDP_{base} . The resulting daily electricity cost is 1.03€, which is 1.71€
633 (62.5%) lower than the BAU strategy, and approaches the PIO within 0.05€
634 (5%). When E^{req} and T^{dep} are unknown, the choice of K_1 must be done
635 carefully such that it results in the desired trade-off between electricity cost
636 and the amount of uncharged energy. For example, if for 2% of the sessions,
637 the EV leaves with a fraction of uncharged energy greater than 25%, a daily
638 electricity cost improvement over BAU of 1.31€ (48%) is obtained. This
639 value of the electricity cost is 0.4€ (39%) higher compared to the best per-
640 forming instance when the departure time and energy require to fully charge
641 the EV are known.

642 4.2.2. Experiment 2: scalability in terms of fleet size

643 This experiment evaluates the scalability of the proposed control frame-
644 work to larger fleet sizes for the three MDPs. The historical data at Ener-
645 gyVille containing measurements for $N_{ev} = 8$, is used to generate training
646 and testing data for other fleet sizes. The BAU strategy, PIO strategy and
647 the RL algorithm with the three MDPs are tested for $N_{ev} = \{2, 8, 16, 32\}$.
648 The MDP hyper-parameters are $L_c = 60$, $N_{past} = N_{fut} = 2$ and $N_{par} =$
649 $20 \times N_{ev}/8$ (MDP_{agg}). For $N_{ev} = 2$, the actor-critic network is modified to
650 one shared layer with 64 nodes and two layers with 32 nodes each for the

651 actor and the critic networks. Additionally, for $N_{ev} = 2$, the `learning_rate`
 652 PPO hyper-parameter is set to 0.0004 (four times higher than before). Due
 653 to the limited computational resources, the PIO strategy is trained with
 654 $L_c = \{15, 15, 30, 60\}$ [minutes] for $N_{ev} = \{2, 8, 16, 32\}$.

655 The resulting learning curves are shown in Fig. 9 and numerical results,
 656 shown in Table 1, contain for each instance the mean value of the electricity
 657 cost between time steps 9.6×10^6 and 14.4×10^6 . The RL algorithm learns
 658 an effective policy for all three MDPs and for all tested fleet sizes with a
 659 reduction in electricity cost between 46% and 63%.

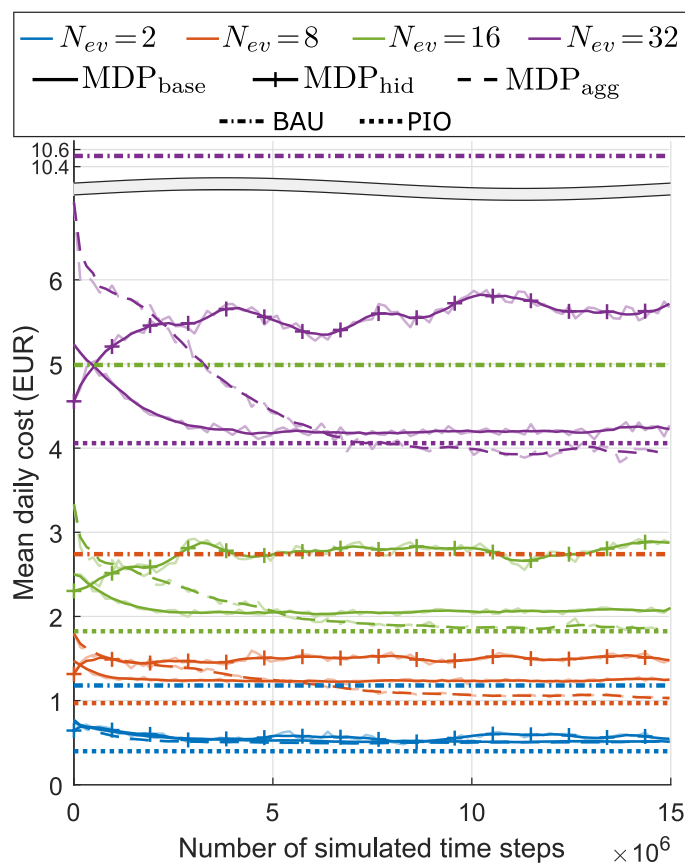


Figure 9: Mean daily $C^{elec}[\text{€}]$ measured on test set for BAU, PIO and the three MDPs for several simulated fleet sizes N_{ev}

660 Compared to MDP_{agg} , MDP_{base} converges faster but obtains a higher final
 661 cost. At $N_{ev} = 2$ both MDPs obtain a similar performance. For larger fleet
 662 sizes, the absolute difference in electricity cost between the MDPs increases,

$N_{ev} =$	2	8	16	32
MDP _{base}	0.515	1.238	2.058	4.208
MDP _{hid}	0.583	1.506	2.740	5.698
MDP _{agg}	0.505	1.051	1.871	3.968
BAU	1.179	2.739	4.988	10.525
PIO	0.398	0.972	1.824	4.059

Table 1: Mean daily $C^{elec}[\text{€}]$ measured on test set for BAU, PIO, and the three MDPs for several simulated fleet sizes N_{ev}

663 for example a difference of 0.187€ for $N_{ev} = 8$ and 0.240€ for $N_{ev} = 32$.
664 The curse of dimensionality makes learning more difficult in high dimensional
665 state-action spaces, and is noticeable for MDP_{base}. Due to the state-action
666 space aggregation, MDP_{agg} is able to mitigate the curse of dimensionality.
667 At $N_{ev} = 32$, the electricity cost for MDP_{base} is less than that obtained by
668 the PIO. This is mainly due to the ability to take an aggregate action for the
669 whole fleet rather than individual actions making it more computationally
670 feasible.

671 The simulation results presented above show that the proposed RL control
672 framework is suitable for coordinating the charging of a fleet of EVs.
673 When knowledge on the departure time of the EV and the energy required
674 to fully charge the EV before its departure are available, the proposed control
675 framework has limited scalability issues especially when the aggregate
676 MDP formulation is used. Even though the proposed control framework
677 outperforms the BAU model when knowledge on departure time and energy
678 required to fully charge the EV are not available, there is no guarantee of the
679 EV being fully charged before departure. Furthermore, it is worth mention-
680 ing that scalability is expected beyond the tested fleet sizes, especially when
681 MDP_{agg} is used as has been shown in [23].

682 In summary, to learn a cost effective control policy to efficiently coordi-
683 nate the charging of large EV fleets, it is necessary to invest in a charging
684 infrastructure that allows obtaining information on the departure time of
685 the EVs and the energy required to fully charge the EVs before departure.
686 The proposed control framework with the aggregate MDP, the real-time con-
687 troller, and a control time step of 60 minutes would be a suitable choice for
688 coordinated charging of large EV fleets.

689 **5. Conclusion**

690 This paper proposes a hybrid control framework that combines the strengths
691 of reinforcement learning and rule-based control for coordinating the charg-
692 ing of an electric vehicle fleet in an office building. Specifically, the framework
693 applies proximal policy optimization, a policy-gradient based reinforcement
694 learning algorithm, to provide a charging schedule with coarse time gran-
695 ularity, which is refined by a rule-based controller to per minute real-time
696 control actions. The control objective is to maximise self-consumption of the
697 local electricity generation and minimize electricity cost. The performance of
698 the proposed framework was evaluated using real-world data from an office
699 building.

700 The simulation results show that the proposed control framework success-
701 fully schedules the charging of an electric vehicle fleet to achieve the control
702 objective. It largely outperforms a business-as-usual strategy and approaches
703 a near-optimal strategy with a 5% performance gap when charging sessions
704 are aggregated before optimization. Simulation results equally show per-
705 formance improvements when information on departure time and required
706 energy is available, and when real-time information on local photo-voltaic
707 electricity generation is used to optimize on a fine time scale.

708 Future work aims to investigate hierarchical reinforcement learning as a
709 replacement for the proposed hybrid method. A downside of the proposed
710 rule-based real-time controller is that it requires expert knowledge and is not
711 generalisable to other settings with a different objective function, such as
712 peak shaving. Hierarchical reinforcement learning may provide a generalis-
713 able and scalable solution for learning on both coarse and fine timescales.
714 Additionally, the proposed framework will be evaluated, and where needed
715 adapted, to better serve the existing range of charging modalities. Specifi-
716 cally the support of fast-charging and vehicle-to-grid charging is of interest
717 to better support the available charging infrastructure at EnergyVille.

718 **Acknowledgement**

719 This research is funded by the Flemish Institute for Technological Re-
720 search (VITO) through a PhD scholarship, and the ROLECS project –
721 Flux50-VLAIO-HBC.2018.0527.

722 **References**

- 723 [1] International Energy Agency, Global EV Outlook 2019, <https://www.>
724 [iea.org/reports/global-ev-outlook-2019](https://www.iea.org/reports/global-ev-outlook-2019), 2019. Accessed: 2020-
725 08-07.
- 726 [2] M. Blonsky, A. Nagarajan, S. Ghosh, K. McKenna, S. Veda, B. Kro-
727 poski, Potential impacts of transportation and building electrification on
728 the grid: A review of electrification projections and their effects on grid
729 infrastructure, operation, and planning, *Current Sustainable/Renewable*
730 *Energy Reports* 6 (2019) 169–176. doi:10.1007/s40518-019-00140-5.
- 731 [3] J. Hu, S. You, M. Lind, J. Østergaard, Coordinated charging of elec-
732 tric vehicles for congestion prevention in the distribution grid, *IEEE*
733 *Transactions on Smart Grid* 5 (2013) 703–711. doi:10.1109/TSG.2013.
734 2279007.
- 735 [4] M. Van Der Kam, W. van Sark, Smart charging of electric vehicles
736 with photovoltaic power and vehicle-to-grid technology in a microgrid; a
737 case study, *Applied energy* 152 (2015) 20–30. doi:10.1016/j.apenergy.
738 2015.04.092.
- 739 [5] N. Sadeghianpourhamami, J. Deleu, C. Develder, Definition and evalu-
740 ation of model-free coordination of electrical vehicle charging with rein-
741 forcement learning, *IEEE Transactions on Smart Grid* 11 (2020) 203–
742 214. doi:10.1109/TSG.2019.2920320.
- 743 [6] A. R. Bhatti, Z. Salam, A rule-based energy management scheme
744 for uninterrupted electric vehicles charging at constant price using
745 photovoltaic-grid system, *Renewable energy* 125 (2018) 384–400. doi:10.
746 1016/j.renene.2018.02.126.
- 747 [7] D. Wang, F. Locment, M. Sechilariu, Modelling, simulation, and man-
748 agement strategy of an electric vehicle charging station based on a DC
749 microgrid, *Applied Sciences* 10 (2020) 2053. doi:[https://doi.org/10.](https://doi.org/10.3390/app10062053)
750 [3390/app10062053](https://doi.org/10.3390/app10062053).
- 751 [8] A. Di Giorgio, F. Liberati, S. Canale, Electric vehicles charging control
752 in a smart grid: A model predictive control approach, *Control Engi-*
753 *neering Practice* 22 (2014) 147–162. doi:10.1016/j.conengprac.2013.
754 10.005.

- 755 [9] S. Bansal, M. N. Zeilinger, C. J. Tomlin, Plug-and-play model predictive
756 control for electric vehicle charging and voltage control in smart grids,
757 in: 53rd IEEE Conference on Decision and Control, IEEE, 2014, pp.
758 5894–5900. doi:10.1109/CDC.2014.7040312.
- 759 [10] B.-R. Choi, W.-P. Lee, D.-J. Won, Optimal charging strategy based
760 on model predictive control in electric vehicle parking lots considering
761 voltage stability, *Energies* 11 (2018) 1812. doi:10.3390/en11071812.
- 762 [11] R. S. Sutton, A. G. Barto, Reinforcement learning: An introduction,
763 MIT press, 2018.
- 764 [12] J. R. Vázquez-Canteli, Z. Nagy, Reinforcement learning for demand
765 response: A review of algorithms and modeling techniques, *Applied En-
766 ergy* 235 (2019) 1072 – 1089. doi:10.1016/j.apenergy.2018.11.002.
- 767 [13] S. Vandael, B. Claessens, D. Ernst, T. Holvoet, G. Deconinck, Rein-
768 forcement learning of heuristic EV fleet charging in a day-ahead elec-
769 tricity market, *IEEE Transactions on Smart Grid* 6 (2015) 1795–1805.
770 doi:10.1109/TSG.2015.2393059.
- 771 [14] A. Chiş, J. Lundén, V. Koivunen, Reinforcement learning-based plug-
772 in electric vehicle charging with forecasted price, *IEEE Transactions
773 on Vehicular Technology* 66 (2016) 3674–3684. doi:10.1109/TVT.2016.
774 2603536.
- 775 [15] S. Wang, S. Bi, Y. A. Zhang, A reinforcement learning approach for
776 EV charging station dynamic pricing and scheduling control, in: 2018
777 IEEE Power Energy Society General Meeting (PESGM), 2018, pp. 1–5.
778 doi:10.1109/PESGM.2018.8586075.
- 779 [16] Z. Wan, H. Li, H. He, D. Prokhorov, Model-free real-time EV charging
780 scheduling based on deep reinforcement learning, *IEEE Transactions on
781 Smart Grid* 10 (2019) 5246–5257. doi:10.1109/TSG.2018.2879572.
- 782 [17] J. Lee, E. Lee, J. Kim, Electric vehicle charging and discharging al-
783 gorithm based on reinforcement learning with data-driven approach
784 in dynamic pricing scheme, *Energies* 13 (2020) 1950. doi:10.3390/
785 en13081950.

- 786 [18] C. J. C. H. Watkins, Learning from delayed rewards (1989). King's
787 College, Cambridge.
- 788 [19] L. Yu, W. Xie, D. Xie, Y. Zou, D. Zhang, Z. Sun, et al., Deep Reinforce-
789 ment Learning for Smart Home Energy Management, IEEE Internet of
790 Things Journal 7 (2020) 2751–2762. doi:10.1109/JIOT.2019.2957289.
- 791 [20] D. Qiu, Y. Ye, D. Papadaskalopoulos, G. Strbac, A deep reinforcement
792 learning method for pricing electric vehicles with discrete charging lev-
793 els, IEEE Transactions on Industry Applications 56 (2020) 5901–5912.
794 doi:10.1109/TIA.2020.2984614.
- 795 [21] H. Li, Z. Wan, H. He, Real-time residential demand response, IEEE
796 Transactions on Smart Grid 11 (2020) 4144–4154. doi:10.1109/TSG.
797 2020.2978061.
- 798 [22] A. Moonens, A. Nowé, Fine-grained control of electric vehicle charging
799 with policy gradient, in: Proceedings of the Adaptive and Learning
800 Agents Workshop 2019 (ALA-19), 2019.
- 801 [23] S. Vandael, B. Claessens, M. Hommelberg, T. Holvoet, G. Deconinck,
802 A scalable three-step approach for demand side management of plug-in
803 hybrid vehicles, IEEE Transactions on Smart Grid 4 (2013) 720–728.
804 doi:10.1109/TSG.2012.2213847.
- 805 [24] P. E. Gill, E. Wong, Sequential quadratic programming methods, in:
806 Mixed integer nonlinear programming, Springer, 2012, pp. 147–224.
- 807 [25] Z. Wan, H. Li, H. He, D. Prokhorov, Model-Free Real-Time EV Charg-
808 ing Scheduling Based on Deep Reinforcement Learning, IEEE Trans-
809 actions on Smart Grid 10 (2019) 5246–5257. doi:10.1109/TSG.2018.
810 2879572.
- 811 [26] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, O. Klimov, Proximal
812 Policy Optimization Algorithms, CoRR abs/1707.06347 (2017). URL:
813 <http://arxiv.org/abs/1707.06347>.
- 814 [27] D. Kraft, A software package for sequential quadratic program-
815 ming, Ein Software-Paket zur sequentiellen quadratischen Optimierung,
816 Forschungsbericht., Technical Report, Institut für Dynamik der Flugsys-
817 teme, Deutsche Forschungs- und Versuchsanstalt für Luft- und Raum-
818 fahrt *DFVLR*, Oberpfaffenhofen, Braunschweig, 1988.