

---

# Chargax: A JAX Accelerated EV Charging Simulator

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Deep Reinforcement Learning can play a key role in addressing sustainable energy  
2 challenges. For instance, many grid systems are heavily congested, highlighting  
3 the urgent need to enhance operational efficiency. However, reinforcement learning  
4 approaches have traditionally been slow due to the high sample complexity and ex-  
5 pensive simulation requirements. While recent works have effectively used GPUs  
6 to accelerate data generation by converting environments to JAX, these works  
7 have largely focussed on classical toy problems. This paper introduces Chargax,  
8 a JAX-based environment for realistic simulation of electric vehicle charging sta-  
9 tions designed for accelerated training of RL agents. We validate our environment  
10 in a variety of scenarios based on real data, comparing reinforcement learning  
11 agents against baselines. Chargax delivers substantial computational performance  
12 improvements of over 100x-1000x over existing environments. Additionally, Char-  
13 gax’ modular architecture enables the representation of diverse real-world charging  
14 station configurations.<sup>1</sup>

## 15 1 Introduction

16 Deep Reinforcement Learning (RL) can approxi-  
17 mate optimal policies for difficult decision prob-  
18 lems that are impossible to solve with traditional  
19 mathematical methods. Such problems occur fre-  
20 quently in sustainable energy challenges such as  
21 operation of windfarms (Fernandez-Gauna et al.,  
22 2022), electric vehicle charging (Rehman et al.,  
23 2024), and nuclear fusion reactors (Seo et al.,  
24 2024). While RL has achieved successful solu-  
25 tions to these challenges, further development of  
26 RL algorithms hinges on the availability of real-  
27 istic simulation environments and benchmarks  
28 (Ponse et al., 2024).

29 Unfortunately, reinforcement learning is noto-  
30 riously sample-inefficient (Yarats et al., 2020;  
31 Kaiser et al., 2024). It often requires many envi-  
32 ronments samples which are slow and possibly  
33 expensive to generate. These simulations have  
34 often been running on the CPU – disallowing RL  
35 researchers from truly harvesting the potential  
36 scale-up of GPUs that other machine learning fields have been enjoying (Scarfe et al., 2025). To this

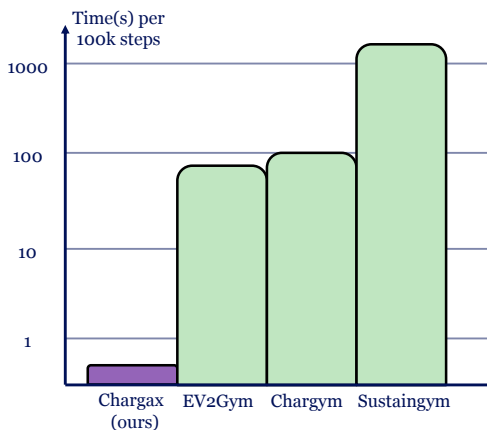


Figure 1: Comparison between Chargax and prior EV Gym Environments in seconds to complete 100k training steps using PPO. See Table 2 for a more complete overview.

<sup>1</sup>Available on GitHub at <https://github.com/anonymous-submission/anonymous-submission>

37 end, the development of RL environments using JAX (Bradbury et al., 2018) has recently gained  
 38 increasing attention (Freeman et al., 2021; Lange, 2022; Pignatelli et al., 2024; Bonnet et al., 2024).  
 39 However, current implementations remain largely confined to simplified toy problems, highlighting a  
 40 significant gap in real-world applications utilizing JAX.

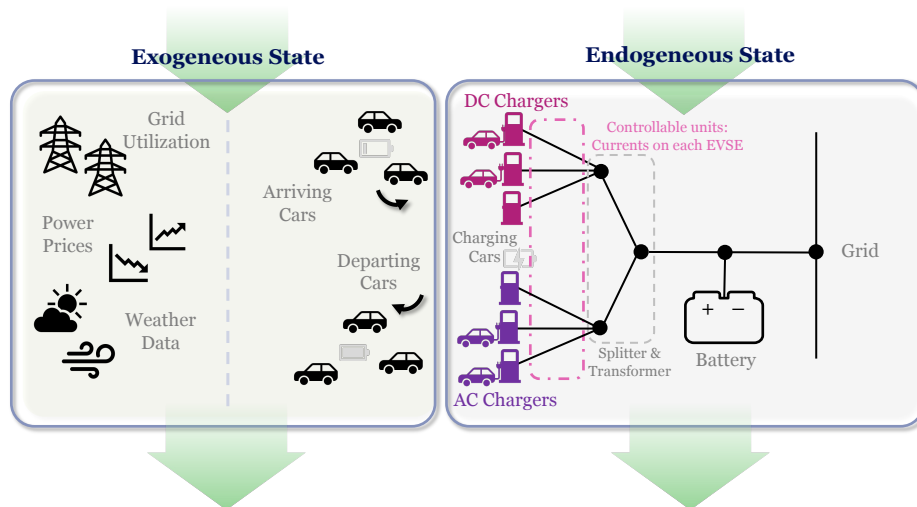


Figure 2: An overview of the Chargax environment. The *endogenous state* describes the state variables that are influenced directly by the agent. The *exogenous state* evolves via, agent-independent, predefined time series data.

41 **Contribution** In this work, we aim to bridge this gap by introducing, to the best of our knowledge,  
 42 the first reinforcement learning environment for EV charging implemented in JAX.

- 43 • Our environment, *Chargax*, achieves a significant speedup of 100x-1000x compared to exist-  
 44 ing environments for EV charging (Yeh et al., 2024; Orfanoudakis et al., 2024; Karatzinis  
 45 et al., 2022). This lowers training times from hours or even days to mere minutes – allowing  
 46 for orders of magnitude more experiments (see Figure 1).
- 47 • *Chargax* extends the generalisability of existing frameworks. As highlighted in a recent  
 48 survey (Alaee et al., 2023), optimising electric vehicle charging strategies involves a diverse  
 49 set of potential objectives. We demonstrate that many of these objectives can be addressed  
 50 within a single simulation framework by ensuring sufficient flexibility.
- 51 • *Chargax* can function as a high-performance test bed for reinforcement learning benchmark-  
 52 ing on real-world applications. Empirically, we demonstrate how RL agents are able to  
 53 outperform baselines and allow for flexible goals such as user satisfaction. We open source  
 54 *Chargax*<sup>1</sup> for the wider community to experiment with.

55 *Chargax* is equipped with predefined datasets, reward functions, and charging station architectures for  
 56 various scenarios. Moreover, all components are fully customizable, enabling researchers to tailor the  
 57 environment to specific requirements, thereby facilitating efficient and adaptable RL-based solutions  
 58 for EV charging optimization.

## 59 2 Related Work

60 Prior work in EV charging includes the gym environments *Sustaingym* (Yeh et al., 2024) (based on  
 61 (Lee et al., 2020b)), *Chargym* (Karatzinis et al., 2022), and the more recently released *EV2Gym*  
 62 (Orfanoudakis et al., 2024). Compared to the works of Yeh et al. 2024; Lee et al. 2020b and Karatzinis  
 63 et al. 2022 our framework provides additional flexibility for the architecture of the charging station,  
 64 scenario selection, and the customer and car profiles. Compared to (Orfanoudakis et al., 2024),  
 65 which also prioritises flexibility, our approach features a more streamlined state and architecture  
 66 representation. To the best of our knowledge, *Chargax* is the only Gym-like environment that includes

67 local car and price data across multiple regions. Furthermore, Chargax is orders of magnitude faster  
68 and in turn allows for large scale experiments on the GPU (See Figure 1). Apart from these Gym-like  
69 simulators, there exist a history of EV charging simulators (Saxena, 2013; Rigas et al., 2018; Balogun  
70 et al., 2023; Cañigüeral, 2023).

71 In recent years, many classical Gym environments have been reimplemented in JAX. We direct the  
72 reader to the following non-exhaustive list (Freeman et al., 2021; Lange, 2022; Nikulin et al., 2023;  
73 Rutherford et al., 2023; Koyamada et al., 2023; Pignatelli et al., 2024; Bonnet et al., 2024). These  
74 implementations have largely been reimplementations of classical toy problems, highlighting a gap in  
75 environments modelling real-world problems.

### 76 3 Preliminaries

#### 77 Markov Decision Process

78 Formally, an environment is represented as a Markov Decision Process (MDP; Sutton & Barto 2018)  
79 defined by a tuple  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, p_0, p, r, \gamma)$ . Here,  $\mathcal{S}$  is a state space,  $\mathcal{A}$  is a action space,  $p_0 \in \Delta(\mathcal{S})^2$   
80 is the initial state distribution,  $p(\cdot|s, a) \in \Delta(\mathcal{S})$  is the probabilistic transition function,  $r(s, a, s')$   
81 denotes the reward function and  $\gamma \in [0, 1)$  is the discount factor. In the next section (4), we provide  
82 a detailed discussion of the motivation behind the choices for each MDP component and formally  
83 define these quantities used within the framework.

#### 84 JAX

85 JAX is a Python library aimed at accelerator-orientated programming with a NumPy interface (Brad-  
86 bury et al., 2018). It offers function transformations to perform, for example, just-in-time-compilation,  
87 vectorization, and differentiation. Although JAX is a common foundation for deep learning frame-  
88 works (Heek et al., 2024; Kidger & Garcia, 2021), its just-in-time compilation transformation allows  
89 users to easily run plain Python code on accelerators such as GPUs and TPUs. Although JAX imposes  
90 some constraints on how these functions should be constructed, it enables complete environment  
91 transition functions to operate on the GPU. This allows many more operations and environments to  
92 run in parallel and eliminates data transfers between the CPU and GPU for gradient descent updates,  
93 both of which can potentially decrease the computational time requirements of reinforcement learning  
94 experiments significantly (Lu, 2024; Hessel et al., 2021).

### 95 4 Environment Design

96 In many real-world control environments not all state variables are directly affected by the actions of  
97 the agent. Instead, some of the state variables transition into their next state via an (agent-independent)  
98 function (often time series). These functions often rely on some external data source and therefore  
99 these variables describe exactly the entry points for data integration that can be flexibly interchanged  
100 within Chargax. Although this data distinction is often implicitly present (Ponse et al., 2024), we will  
101 formalise this separation explicitly in Chargax to make clear which parts of the state can flexibly be  
102 interchanged.

103 Consequently, we split the environment state in an *endogenous* and an *exogenous* state space. The  
104 endogenous state space refers to the typical state variables that are influenced by the agents' actions  
105 during the transition function. In contrast, exogenous state variables transition into their next state via  
106 an (agent-independent) time series. Examples of exogenous state variables are weather variables, or  
107 national electricity prices. Even though these variables are not affected by the agents' actions, they  
108 may influence the agent by providing an additional learning signal and/or alter the reward.

109 An overview of Chargax is shown in Figure 2 and in the following we provide a high-level overview  
110 of the Chargax environment. Full implementation details, including all equations for transition  
111 dynamics and reward functions, are provided in Appendix A.

---

<sup>2</sup> $\Delta(\mathcal{X})$  denotes the set of probability distributions over a set  $\mathcal{X}$

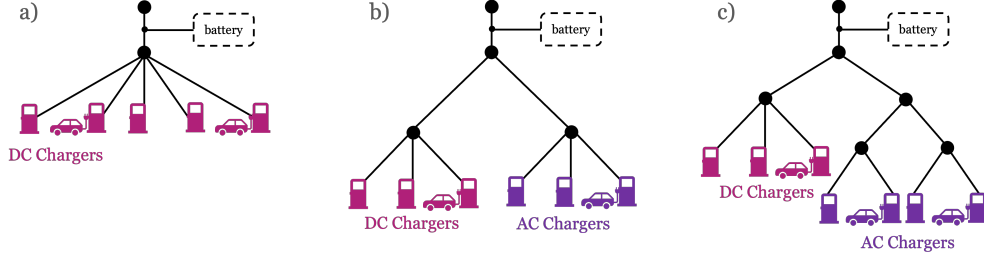


Figure 3: Trees representing different architectures: a) simplest scenario, one type of charger; b) multiple types of chargers, one splitter per charging type; c) multiple types of chargers with multiple splitters per type, imposing additional constraints on the currents. Each node represents a combination of splitters, transformers, cables, and other electrical components.

## 112 EV Station Layout

113 When initialising a Chargax environment, a fixed architectural design for the station is generated or  
 114 provided. This design is fixed and, therefore, not influenced by the transition function. We represent  
 115 this electronic infrastructure of the charging station in the form of a tree (Lee et al., 2021), with  
 116 leaves representing the charging ports (Electric Vehicle Supply Equipment; EVSE; Lee et al. (2020b))  
 117 (see Figure 3). The root node represents the grid connection access, and all other nodes represent a  
 118 combination of splitters, cables, and transformers, and are equipped with a maximum power capacity  
 119 and efficiency coefficient, imposing constraints on the system. In Chargax, we additionally assume a  
 120 fixed voltage  $V$  for each of the EVSEs in the architecture.

121 Chargax supplies methods for generating some charging station architectures. However, custom  
 122 architectures can be built by constructing a tree of simple nodes to mirror existing real-world  
 123 infrastructure.

## 124 Endogenous State Space

125 The endogenous state consists of the state of the various charging ports and their connected cars, and  
 126 the station battery. As each charging port (and the battery) has a fixed voltage level, we allow the  
 127 actual power drawn to be regulated by controlling the current (Orfanoudakis et al., 2024). Losses are  
 128 incorporated through efficiency coefficients at each node (including the charging ports).

129 In addition to the set current at each charging port ( $I_{\text{drawn}}(t) \in [0, I_{\text{max}}]$ ), and whether the port  
 130 is currently occupied ( $\mathbb{1}_{\text{occup}}$ ), the endogenous state contains information for the connected cars.  
 131 This includes their state-of-charge (SoC) and the remaining required power  $\Delta E_{\text{remain}}$ . Additional  
 132 information for each car is supplied exogenously and remains fixed until the car leaves. We will  
 133 expand on this information in the next section. The endogenous state space can optionally be  
 134 expanded with a station battery. This battery is modelled similarly to an EVSE – with a fixed voltage  
 135 and controlled via the set current. The battery allows the agent to store energy to facilitate effective  
 136 discharging strategies. In brief summary, the endogenous state is represented by:

- 137 • For each EVSE:  $I_{\text{drawn}}(t) \in \mathbb{R}_{\geq 0}$ ,  $\mathbb{1}_{\text{occup}}(t) \in \{0, 1\}$ ,  $\Delta E_{\text{remain}}$ ,  $\text{SoC}(t)$
- 138 • Battery:  $I_{\text{battery}}(t)$ ,  $\text{SoC}_{\text{battery}}(t)$

139 Enumerating the existing EVSEs by  $i = 1, \dots, N$ , the total endogenous state space can be expressed  
 140 as  $s_{\text{end}} = (s_{\text{battery}}, s_{c,1}, \dots, s_{c,N})$ . A complete overview of the state space is given in Appendix A.1.

## 141 Exogenous State Space

142 As described previously, the exogenous state variables evolve independently of the agent’s actions.  
 143 As such, the remainder of the variables discussed here are typically sampled from distributions that  
 144 are generated via a provided time series or some predefined function. Currently, Chargax works with  
 145 exogenous state variables for arrival data, user profiles, car profiles, and grid price data.

146 The **arrival data** represents the number of cars that arrive at a given timestep. Typically, this depends  
 147 on the time and location of the charging station. Likewise, the location can also stipulate the typical  
 148 **user profile** of the arriving cars. This profile describes the state of the car that is induced by their  
 149 owner, such as the arrival SoC, desired charging level, and time of departure. **Car profile** variables  
 150 are derived from the physical properties of the cars themselves. These include the maximum capacity  
 151 of the car battery and the maximum charge speed. Lastly, the **grid prices** are an important exogenous  
 152 variable for calculating the profit, which is often a large factor in the reward.

153 Chargax comes equipped with a variety of standard datasets (see Table 1), most of which are based on  
 154 real data. These datasets can be used to sample exogenous variables that resemble realistic scenarios.  
 155 For example, Europe and the US have a different distribution of electric vehicles on the road; in  
 156 turn, the distribution of charging demands is different in both regions. While datasets are provided,  
 157 Chargax is built such that users can use their own data or functions to populate these variables.

## 158 Action Space

At each timestep, the agent controlling the charging station can adjust the power at each EVSE by altering the current (Orfanoudakis et al., 2024), i.e. an action is characterized as

$$a(t) = (\Delta I_i(t))_{i=1}^{N+1} \in \mathbb{R}^{N+1}.$$

159 Here, for the sake of notational convenience, the battery is treated as the  $N + 1$ -th charging pole.  
 160 Notably, the agent cannot accept/decline cars and is assumed to serve arriving cars, as long as there  
 161 are free spots.

## 162 Transition Function

163 At a high level, the transition function consists of four sequential steps, which we detail below. Full  
 164 implementation details can be found in Appendix A.2.

- 165 • **Apply Actions** First, we apply the agent’s to adjust the power drawn by each charging port.  
 166 We limit the maximum power by the capacity of the port, as well as the current maximum  
 167 (dis)charging rate of the car stationed at each charging port.
- 168 • **Charge Stationed Cars** With the newly set power levels, we (dis)charge each car over the  
 169 time interval of a timestep. Here, we assume a constant charging rate over the full interval  
 170  $\Delta t$ .
- 171 • **Departure of Cars** Next, cars fully charged (charge-sensitive users) or with no time remain-  
 172 ing (time-sensitive users) will leave.
- 173 • **Arrival of new Cars** Finally, an amount of new cars will be sampled through our exogenous  
 174 data, along with a *user profile* and *car profile*. The amount of new cars is clipped by the  
 175 number of free spots available and the remaining cars are automatically rejected. Arriving  
 176 cars will park in the first available spot as provided by the provided station architecture.

## 177 Reward Function

178 In RL, the reward functions reflects the notion of optimality, i.e. the desired behaviour. In this section,  
 179 we outline some of the reward functions that are available in Chargax, and how they reflect different  
 180 objectives. We provide additional details in Appendix A.3.

Table 1: Overview of available Profiles in Chargax. Default settings are marked in bold.

Price Profiles	Architectures	Car Distributions	Arrival Frequency	User Profiles
NL	Simple: Single	<b>Europe</b>	Low Traffic	Highway
FR	Charger Type	US	<b>Medium Traffic</b>	Residential
DE	<b>Simple: Multiple</b>	World	High Traffic	Work
<i>Custom</i>	<b>Charger Types</b>	<i>Custom</i>	<i>Custom</i>	<b>Shopping</b>
	<i>Custom</i>			<i>Custom</i>

181 **Profit Maximisation** Profit maximisation lies at the core of most Charging Station Operations  
 182 (Alinejad et al., 2021; Chang et al., 2021; Mirzaei & Kazemi, 2021; Ye et al., 2022). The amount of  
 183 net energy transferred into cars in the interval  $[t, t + \Delta t]$  is denoted by  $\Delta E_{\text{net}}(t)$ . The amount of  
 184 energy fed into the grid as a result from discharging cars is denoted by  $\Delta E_{\rightarrow \text{grid}}(t)$ , and the amount of  
 185 energy that has to be drawn from the net to transfer set levels of energy into cars  $\Delta E_{\text{grid} \rightarrow}(t)$ . Lastly,  
 186 the energy contributed by (dis-)charging the battery  $\Delta E_{\text{b,net}}(t)$  has to be incorporated, resulting in  
 187 the following net energy that is drawn from (or pushed into) the grid

$$\Delta E_{\text{grid,net}} = \Delta E_{\text{grid} \rightarrow}(t) + \Delta E_{\rightarrow \text{grid}}(t) + \Delta E_{\text{b,net}}(t). \quad (1)$$

188 We further assume that the price at which we sell and buy power from car owners is the same, i.e.  
 189  $p_{\text{sell}}$ . This results in the following profit

$$\Pi(t) = \begin{cases} p_{\text{sell}}(t) \cdot \Delta E_{\text{net}}(t) - p_{\text{buy}}(t) \cdot \Delta E_{\text{grid,net}} - c_{\Delta t} & \Delta E_{\text{grid,net}} > 0, \\ p_{\text{sell}}(t) \cdot \Delta E_{\text{net}}(t) - p_{\text{sell,grid}}(t) \cdot \Delta E_{\text{grid,net}} - c_{\Delta t} & \Delta E_{\text{grid,net}} \leq 0. \end{cases} \quad (2)$$

190 Here,  $c_{\Delta t}$  denotes the fixed cost for running the facility per  $\Delta t$ .

191 **Profit Maximisation under constraints** To further steer agents' learnt behaviour in a direction,  
 192 constraints can be induced to penalise certain (undesired) behaviour through penalty terms  $c(t)$ . The  
 193 resulting reward will be the profit minus the linear combination of (possibly) multiple penalty terms

$$r(t) = \Pi(t) - \sum_c \alpha_c c(t). \quad (3)$$

194 Different linear combinations of different penalty terms allow Chargax to be flexible in its optimization  
 195 objective. Chargax comes equipped with various of these penalty terms to better optimize for, for  
 196 example, customer satisfaction, battery degradation, or violating node constraints. We provide a more  
 197 complete list of possible penalty terms along with a formal expression in Appendix A.3. However, we  
 198 emphasise that these are mere suggestions, and that these rewards are not comprehensive in reflecting  
 199 the full landscape of Charging Station Optimisation challenges, and we encourage users to customise  
 200 their reward function within the provided framework.

## 201 5 Experiments

202 In this section, we demonstrate the use of Chargax across different included scenarios. Additionally,  
 203 we highlight performance improvements of Chargax compared to previous EV charging simulations.  
 204 Full details of the used model and configuration parameters, along with additional experimental  
 205 results, can be found in Appendix B and D respectively.

206 In Figure 4a, we have trained a standard PPO agent based on PureJaxRL (Lu et al., 2022). We trained  
 207 on our included *shopping* scenario in varying amounts of traffic using a 16 charger station (10 DC, 6  
 208 AC). We observe how our PPO agent increases its profit over a standard baseline. The baseline is set  
 209 to always charge to its maximum potential within the constraints of the EVSE and the connected car.  
 210 As expected, the potential for profit increases in scenarios with higher amounts of traffic, but this  
 211 increase diminishes as we kept the charging station size the same.

212 Our baseline should yield a high customer satisfaction as customers should be charged within the  
 213 minimum amount of possible time. In contrast, our charging station agent may optimize fully for  
 214 short-term profit without consideration of user satisfaction. This is likely undesirable and may affect  
 215 long-term profits. However, Chargax allows for flexible reward signals that may optimize for this. In  
 216 Figure 4b and 4c, we trained our PPO agent to optimize for profit and user satisfaction at varying  
 217  $\alpha$  levels. Notably in Figure 4b, we can see the agent manages to find preferential policies that  
 218 substantially increase user satisfaction (decrease the amount of kWh that was not charged at departure  
 219 time), while keeping profit levels quite similar.

220 Beyond finding appropriate reward signals, real-world deployment typically involves training an  
 221 agent on historical exogenous data. During deployment, the agent likely encounters data that is  
 222 has not yet observed. Possibly, the entire data set has shifted, for example, due to a rise in energy  
 223 prices year-over-year. Therefore, it is important that system that deal with exogenous time-series  
 224 data can deal with – and test for – this distribution shift (Yeh et al., 2024). As Chargax is flexibly  
 225 designed to allow for any exogenous data, it readily allows to test for these distribution shift problems

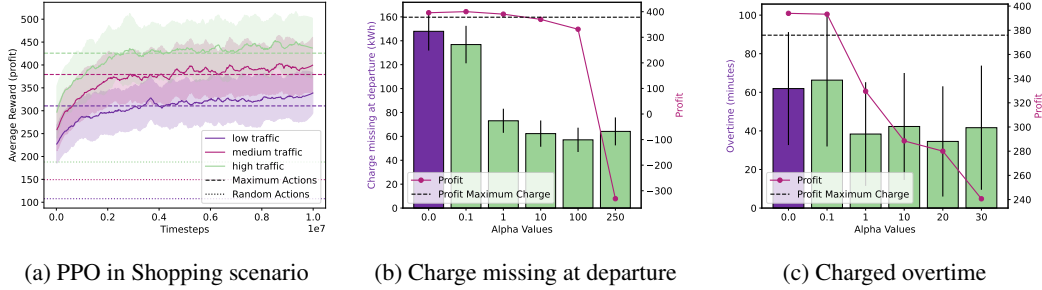


Figure 4: In a) average episode rewards during training a PPO agent in the shopping scenario with different levels of traffic. The RL solution manages to increase profit over the baseline that always charged the maximum possible amount. In b) and c), user satisfaction measured as charge (kWh) missing at time of departure (b), and time exceeded to fully charge cars (c). Higher  $\alpha$ -values weigh the measured variable greater in the reward (Eq. 3). Increasing user satisfaction tends to decrease daily profit. Notably however in b), optimizing for user satisfaction has steered the agent to find policies that reduce the missing charge percentages while retaining a near-identical profit level.

226 – as is displayed in Figure 5, where we have trained and evaluated RL agents on data of different  
 227 price electricity years. Interestingly, although rewards would be assumed to peak when training and  
 228 testing in the same year, employing data from 2021 or 2023 actually yielded higher rewards in 2022  
 229 compared to using the 2022 data directly. The EU region experienced significant energy price surges  
 230 in 2022, likely complicating the training process with the data for this year.

231 Table 2 and Figure 1, showcases the performance of our environment compared to existing EV  
 232 charging simulations that support reinforcement learning through a Gym API. We can see that in a  
 233 typical training scenario, we can decrease learning times by factors exceeding 100. It is important to  
 234 acknowledge that these environments are not identical and might simulate different behaviours (for  
 235 example, SustainGym does not allow discharging). Therefore, this comparison may be considered  
 236 rough. However, the significant differences in scale clearly demonstrate the advantages of using  
 237 Chargax- and JAX-based environments for RL in general. Training cycles can be reduced entire  
 238 working days to well under 5 minutes, allowing for many more iterations of training and testing.

## 239 6 Discussion & Conclusion

240 This work presented Chargax, an EV charging simulator built in JAX. Chargax aims to bridge the gap  
 241 between toy problems and real-world implementations, accelerating simulations while maintaining  
 242 practical relevance. However, it remains a simulator, constrained by simplifying assumptions,  
 243 requiring future work to further close the gap between simulation and deployment.

Table 2: Performance comparison between Chargax and other EV charging Gym environments, based on data collected by performing 100k environment steps. We evaluated both taking random actions (assessing the performance of the transition function), and a training a PPO agent. The PPO agent was tested both with a single environment, and in a more typical training scenario with vectorized environments. Here we observe performance improvements of over 100x. The results are obtained on an NVIDIA RTX 4000 Ada GPU and an AMD EPYC 2.8 GHz CPU. For the comparison environments, we used Stable-Baselines3 (Raffin et al., 2021) with CUDA enabled for the PPO implementation.

	Chargax	Ev2Gym	Chargym	Sustaingym
			Speedup	Speedup
Random	1.36	77.95	57x	1554.57
PPO (1)	9.79	170.05	17x	1718.71
PPO (16)	0.65	86.99	<b>134x</b>	<b>1836.00</b>

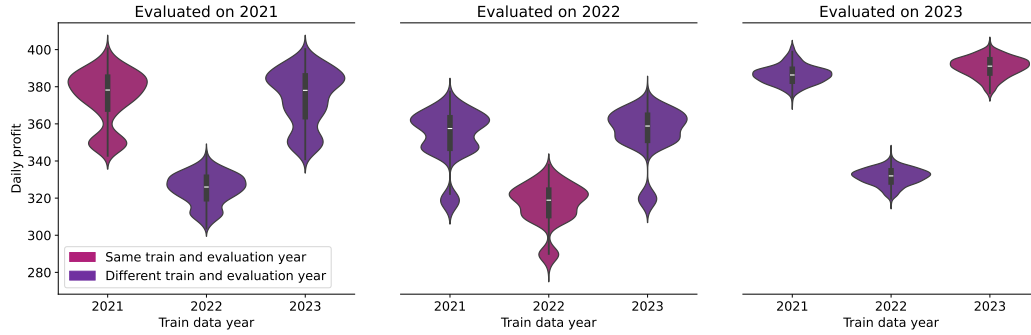


Figure 5: A PPO agent trained and tested on three separate years of Dutch electricity prices. For each experiment, a particular year’s data was used for training, while testing on a fixed year. Substantial price increases in the year 2022 results in suboptimal training when using this year’s data – even when evaluating on this same year.

244 Our model assumes an isolated power network for the EV charging station, avoiding shared trans-  
 245 formers that could introduce uncontrollable constraints. Expanding the model to include additional  
 246 control variables, such as dynamic pricing strategies or vehicle allocation mechanisms, would increase  
 247 its realism. Furthermore, accounting for temperature dependence in the system, or incorporating  
 248 government-imposed regulatory constraints could make it more reflective of real-world charging sta-  
 249 tions. Furthermore, a natural addition for future work would be to incorporate local energy production  
 250 systems (such as solar panels) and weather data.

251 In its current state, Chargax achieves training time reductions of over 100x, compared to existing  
 252 simulators. Usual training durations of (multiple) working days can be completed in Chargax in  
 253 well under 5 minutes, allowing for many more additional training and testing runs. We have built  
 254 Chargax to be flexible, allowing for custom data sources for the exogenous state, and flexible reward  
 255 structures. However, Chargax does provide base datasets and reward penalties to get started. Chargax  
 256 can also be used as a test bed reinforcement learning benchmarking as it is currently only one of few  
 257 JAX-based environments that models a real-world problem. We open source Chargax for the wider  
 258 community to experiment with.<sup>1</sup>

## 259 References

260 Pegah Alaei, Julius Bems, and Amjad Anvari-Moghaddam. A Review of the Latest Trends in  
 261 Technical and Economic Aspects of EV Charging Management. *Energies*, 16(9):3669, January  
 262 2023. ISSN 1996-1073. DOI: 10.3390/en16093669.

263 Mahyar Alinejad, Omid Rezaei, Ahad Kazemi, and Saeed Bagheri. An Optimal Management for  
 264 Charging and Discharging of Electric Vehicles in an Intelligent Parking Lot Considering Vehicle  
 265 Owner’s Random Behaviors. *Journal of Energy Storage*, 35:102245, March 2021. ISSN 2352152X.  
 266 DOI: 10.1016/j.est.2021.102245.

267 Emmanuel Balogun, Elizabeth Buechler, Siddharth Bhela, Simona Onori, and Ram Rajagopal. Ev-  
 268 ecosim: A grid-aware co-simulation platform for the design and optimization of electric vehicle  
 269 charging infrastructure. *IEEE Transactions on Smart Grid*, 2023.

270 Clément Bonnet, Daniel Luo, Donal Byrne, Shikha Surana, Sasha Abramowitz, Paul Duckworth,  
 271 Vincent Coyette, Laurence I. Midgley, Elshadai Tegegn, Tristan Kalloniatis, Omayma Mahjoub,  
 272 Matthew Macfarlane, Andries P. Smit, Nathan Grinsztajn, Raphael Boige, Cemlyn N. Waters, Mo-  
 273 hamed A. Mimouni, Ulrich A. Mbou Sob, Ruan de Kock, Siddarth Singh, Daniel Furelos-Blanco,  
 274 Victor Le, Arnau Pretorius, and Alexandre Laterre. Jumanji: a diverse suite of scalable reinforcement  
 275 learning environments in jax, 2024. URL <https://arxiv.org/abs/2306.09884>.

276 James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal  
 277 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and  
 278 Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL  
 279 <http://github.com/google/jax>.



- 280 M. Cañigüeral. evsim: Electric vehicle charging sessions simulation, 2023. R package version 1.2.0.  
281 [Online]. Available: <https://github.com/mcanigüeral/evsim/>.
- 282 Shuo Chang, Yugang Niu, and Tinggang Jia. Coordinate scheduling of electric vehicles in charging  
283 stations supported by microgrids. *Electric Power Systems Research*, 199:107418, October 2021.  
284 ISSN 03787796. DOI: 10.1016/j.epsr.2021.107418.
- 285 Onur Elma. A dynamic charging strategy with hybrid fast charging station for electric vehicles.  
286 *Energy*, 202:117680, July 2020. ISSN 03605442. DOI: 10.1016/j.energy.2020.117680.
- 287 Fastned. Charge speed, 2025. URL [https://www.fastnedcharging.com/en/  
288 brands-overview](https://www.fastnedcharging.com/en/brands-overview). Accessed: 2025-02-14.
- 289 Borja Fernandez-Gauna, Manuel Graña, Juan-Luis Osa-Amilibia, and Xabier Larrucea. Actor-critic  
290 continuous state reinforcement learning for wind-turbine control robust optimization. *Information  
291 Sciences*, 591:365–380, April 2022. ISSN 00200255. DOI: 10.1016/j.ins.2022.01.047.
- 292 C. Daniel Freeman, Erik Frey, Anton Raichuk, Sertan Girgin, Igor Mordatch, and Olivier Bachem.  
293 Brax - a differentiable physics engine for large scale rigid body simulation, 2021. URL [http:  
294 //github.com/google/brax](http://github.com/google/brax).
- 295 Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas  
296 Steiner, and Marc van Zee. Flax: A neural network library and ecosystem for JAX, 2024. URL  
297 <http://github.com/google/flax>.
- 298 Matteo Hessel, Manuel Kroiss, Aidan Clark, Iurii Kemaev, John Quan, Thomas Keck, Fabio Viola,  
299 and Hado van Hasselt. Podracer architectures for scalable reinforcement learning. *arXiv preprint  
300 arXiv:2104.06272*, 2021.
- 301 Shahid Hussain, Yun-Su Kim, Subhasis Thakur, and John G. Breslin. Optimization of Waiting  
302 Time for Electric Vehicles Using a Fuzzy Inference System. *IEEE Transactions on Intelligent  
303 Transportation Systems*, 23(9):15396–15407, September 2022. ISSN 1558-0016. DOI: 10.1109/  
304 TITS.2022.3140461.
- 305 Lukasz Kaiser, Mohammad Babaeizadeh, Piotr Milos, Blazej Osinski, Roy H. Campbell, Konrad  
306 Czechowski, Dumitru Erhan, Chelsea Finn, Piotr Kozakowski, Sergey Levine, Afroz Mohiuddin,  
307 Ryan Sepassi, George Tucker, and Henryk Michalewski. Model-Based Reinforcement Learning  
308 for Atari, April 2024. URL <http://arxiv.org/abs/1903.00374>. arXiv:1903.00374 [cs, stat].
- 309 Georgios Karatzinis, Christos Korkas, Michalis Terzopoulos, Christos Tsaknakis, Aliko Stefanopoulou,  
310 Iakovos Michailidis, and Elias Kosmatopoulos. Charym: An ev charging station model for  
311 controller benchmarking. In *IFIP International Conference on Artificial Intelligence Applications  
312 and Innovations*, pp. 241–252. Springer, 2022.
- 313 Patrick Kidger and Cristian Garcia. Equinox: neural networks in JAX via callable PyTrees and  
314 filtered transformations. *Differentiable Programming workshop at Neural Information Processing  
315 Systems 2021*, 2021.
- 316 Sotetsu Koyamada, Shinri Okano, Soichiro Nishimori, Yu Murata, Keigo Habara, Haruka Kita, and  
317 Shin Ishii. Pgx: Hardware-accelerated parallel game simulators for reinforcement learning. In  
318 *Advances in Neural Information Processing Systems*, volume 36, pp. 45716–45743, 2023.
- 319 Robert Tjarko Lange. gymmax: A JAX-based reinforcement learning environment library, 2022. URL  
320 <http://github.com/RobertTLange/gymmax>.
- 321 Munsu Lee, Jinhyeong Park, Sun-Ik Na, Hyung Sik Choi, Byeong-Sik Bu, and Jonghoon  
322 Kim. An Analysis of Battery Degradation in the Integrated Energy Storage System with Solar  
323 Photovoltaic Generation. *Electronics*, 9(4):701, April 2020a. ISSN 2079-9292. DOI:  
324 10.3390/electronics9040701.
- 325 Zachary J. Lee, Sunash Sharma, Daniel Johansson, and Steven H. Low. ACN-Sim: An Open-Source  
326 Simulator for Data-Driven Electric Vehicle Charging Research. <https://arxiv.org/abs/2012.02809v2>,  
327 December 2020b.

- 328 Zachary J. Lee, George Lee, Ted Lee, Cheng Jin, Rand Lee, Zhi Low, Daniel Chang, Christine Ortega,  
329 and Steven H. Low. Adaptive Charging Networks: A Framework for Smart Electric Vehicle  
330 Charging. *IEEE Transactions on Smart Grid*, 12(5):4339–4350, September 2021. ISSN 1949-3061.  
331 DOI: 10.1109/TSG.2021.3074437.
- 332 Yang Li, Meng Han, Zhen Yang, and Guoqing Li. Coordinating Flexible Demand Response and  
333 Renewable Uncertainties for Scheduling of Community Integrated Energy Systems With an Electric  
334 Vehicle Charging Station: A Bi-Level Approach. *IEEE Transactions on Sustainable Energy*, 12(4):  
335 2321–2331, October 2021. ISSN 1949-3037. DOI: 10.1109/TSTE.2021.3090463.
- 336 Chris Lu. luchris429/purejaxrl, September 2024. URL [https://github.com/luchris429/  
337 purejaxrl](https://github.com/luchris429/purejaxrl). original-date: 2023-02-25T15:38:11Z.
- 338 Chris Lu, Jakub Kuba, Alistair Letcher, Luke Metz, Christian Schroeder de Witt, and Jakob Foerster.  
339 Discovered policy optimisation. *Advances in Neural Information Processing Systems*, 35:16455–  
340 16468, 2022.
- 341 Mohammad Javad Mirzaei and Ahad Kazemi. A two-step approach to optimal management of electric  
342 vehicle parking lots. *Sustainable Energy Technologies and Assessments*, 46:101258, August 2021.  
343 ISSN 22131388. DOI: 10.1016/j.seta.2021.101258.
- 344 Alexander Nikulin, Vladislav Kurenkov, Ilya Zisman, Viacheslav Sinii, Artem Agarkov, and Sergey  
345 Kolesnikov. XLand-minigrid: Scalable meta-reinforcement learning environments in JAX. In  
346 *Intrinsically-Motivated and Open-Ended Learning Workshop, NeurIPS2023*, 2023. URL [https:  
347 //openreview.net/forum?id=xALDC4aHGz](https://openreview.net/forum?id=xALDC4aHGz).
- 348 Stavros Orfanoudakis, Cesar Diaz-Londono, Yunus E. Yilmaz, Peter Palensky, and Pedro P. Vergara.  
349 EV2Gym: A Flexible V2G Simulator for EV Smart Charging Research and Benchmarking, April  
350 2024.
- 351 Eduardo Pignatelli, Jarek Liesen, Robert Tjarko Lange, Chris Lu, Pablo Samuel Castro, and Laura  
352 Toni. Navix: Scaling minigrid environments with jax. *arXiv preprint arXiv:2407.19396*, 2024.
- 353 Koen Ponse, Felix Kleuker, Márton Fejér, Álvaro Serra-Gómez, Aske Plaat, and Thomas Moerland.  
354 Reinforcement learning for sustainable energy: A survey. *arXiv preprint arXiv:2407.18597*, 2024.
- 355 Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah  
356 Dormann. Stable-Baselines3: Reliable Reinforcement Learning Implementations. *Journal of  
357 Machine Learning Research*, 22(268):1–8, 2021. ISSN 1533-7928. URL [http://jmlr.org/  
358 papers/v22/20-1364.html](http://jmlr.org/papers/v22/20-1364.html).
- 359 Anis ur Rehman, Haris M Khalid, and SM Muyeen. Grid-integrated solutions for sustainable ev  
360 charging: a comparative study of renewable energy and battery storage systems. *Frontiers in  
361 Energy Research*, 12:1403883, 2024.
- 362 Emmanouil S Rigas, Sotiris Karapostolakis, Nick Bassiliades, and Sarvapali D Ramchurn. Evlibsim:  
363 A tool for the simulation of electric vehicles’ charging stations using the evlib library. *Simulation  
364 Modelling Practice and Theory*, 87:99–119, 2018.
- 365 Alexander Rutherford, Benjamin Ellis, Matteo Gallici, Jonathan Cook, Andrei Lupu, Gardar Ing-  
366 varsson, Timon Willi, Akbir Khan, Christian Schroeder de Witt, Alexandra Souly, Saptarashmi  
367 Bandyopadhyay, Mikayel Samvelyan, Minqi Jiang, Robert Tjarko Lange, Shimon Whiteson,  
368 Bruno Lacerda, Nick Hawes, Tim Rocktaschel, Chris Lu, and Jakob Nicolaus Foerster. Jaxmarl:  
369 Multi-agent rl environments in jax. *arXiv preprint arXiv:2311.10090*, 2023.
- 370 S. Saxena. Vehicle-to-grid simulator, version 00, November 2013. [Online]. Available: [https:  
371 //www.osti.gov/biblio/1437011](https://www.osti.gov/biblio/1437011).
- 372 Tim Scarfe, Jakob Foerster, and Chris Lu. Imagenet moment for reinforcement learning?, feb 2025.  
373 URL [https://www.dropbox.com/s/cl/fi/yqjszhntfr00bhjh6t565/JAKOB.pdf?rlkey=  
374 scvny4bnwj8th42fjv8zsfu2y&e=1&d1=0](https://www.dropbox.com/s/cl/fi/yqjszhntfr00bhjh6t565/JAKOB.pdf?rlkey=scvny4bnwj8th42fjv8zsfu2y&e=1&d1=0). Machine Learning Streetwork podcast episode.

- 375 Jaemin Seo, SangKyeun Kim, Azarakhsh Jalalvand, Rory Conlin, Andrew Rothstein, Joseph Abbate,  
376 Keith Erickson, Josiah Wai, Ricardo Shousha, and Egemen Kolemen. Avoiding fusion plasma  
377 tearing instability with deep reinforcement learning. *Nature*, 626(8000):746–751, 2024.
- 378 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. Adaptive  
379 Computation and Machine Learning Series. The MIT Press, Cambridge, Massachusetts, second  
380 edition, 2018. ISBN 978-0-262-03924-6.
- 381 Fynn Welzel, Carl-Friedrich Klinck, Yannick Pohlmann, and Mats Bednarczyk. Grid and  
382 user-optimized planning of charging processes of an electric vehicle fleet using a quantita-  
383 tive optimization model. *Applied Energy*, 290:116717, May 2021. ISSN 03062619. DOI:  
384 10.1016/j.apenergy.2021.116717.
- 385 Denis Yarats, Amy Zhang, Ilya Kostrikov, Brandon Amos, Joelle Pineau, and Rob Fergus. Improving  
386 Sample Efficiency in Model-Free Reinforcement Learning from Images, July 2020. URL <http://arxiv.org/abs/1910.01741>. arXiv:1910.01741 [cs].  
387
- 388 Zuzhao Ye, Yuanqi Gao, and Nanpeng Yu. Learning to Operate an Electric Vehicle Charging Station  
389 Considering Vehicle-Grid Integration. *IEEE Transactions on Smart Grid*, 13(4):3038–3048, July  
390 2022. ISSN 1949-3061. DOI: 10.1109/TSG.2022.3165479.
- 391 Christopher Yeh, Victor Li, Rajeev Datta, Julio Arroyo, Nicolas Christianson, Chi Zhang, Yize Chen,  
392 Mohammad Mehdi Hosseini, Azarang Golmohammadi, Yuanyuan Shi, et al. Sustaingym: Rein-  
393 forcement learning environments for sustainable energy systems. *Advances in Neural Information*  
394 *Processing Systems*, 36, 2024.
- 395 Yongmin Zhang, Pengcheng You, and Lin Cai. Optimal Charging Scheduling by Pricing for EV  
396 Charging Station With Dual Charging Modes. *IEEE Transactions on Intelligent Transportation*  
397 *Systems*, 20(9):3386–3396, September 2019. ISSN 1558-0016. DOI: 10.1109/TITS.2018.2876287.

## 398 A Environmental Details

### 399 A.1 State Spaces

400 **Exogenous state space** Apart from price data, examples of exogenous state variables include  
 401 power demand of the grid, weather data, or marginal operating emissions rate (MOER, Yeh et al.,  
 402 2024), all of which could influence the maximisation objective but evolve according to some (agent  
 403 independent) time series. It is important to note, that while the environment requires auxiliary data  
 404 for most built-in reward functions, e.g. it is impossible to maximise profit without having access to  
 405 prices, these exogenous state variables may be treated unobservable for the agent. On the contrary,  
 406 one may add data to the exogenous state space, that is not required for any reward calculation, but  
 407 may serve as additional learning signal, for instance day-ahead power prices.

408 Apart from the above examples arrival data, user profiles, and car profiles are part of the exogenous  
 409 state space.

- 410 • **Arrival Data** At each timestep  $t$ , a number of cars  $M(t)$  is characterized as a sample from  
 411 an arrival distribution  $M(t) \sim \mathcal{D}_{\text{arrival}}(t)$ .
- 412 • **Car Profiles** Arriving cars are characterised by their physical properties. This encompasses  
 413 the charging speed  $\hat{r}$  as a function of the SoC. As in (Lee et al., 2020b) we assume a  
 414 piece-wise linear function

$$\hat{r}_{\tau, \bar{r}}(\text{SoC}) = \begin{cases} \bar{r}, & \text{SoC} \leq \tau \\ (1 - \text{SoC}) \frac{\bar{r}}{1 - \tau}, & \text{SoC} > \tau. \end{cases}$$

415 Due to lack of data, we assume that the discharging speed can be obtained by vertically  
 416 flipping the charging curve at  $\text{SoC} = 0.5$ . While we assume, that we have a different  
 417 maximal charging speed for different charger types – by default AC and DC charger – and  
 418 have hence different max charging rates ( $\bar{r} = (\bar{r}_{\text{AC}}, \bar{r}_{\text{DC}})$ ), we assume that both charging  
 419 speed curves use the same  $\tau$ . Lastly, each car has a maximum battery capacity  $C$ , which is  
 420 important for calculating State of Charges. These car profiles are sampled from a pre-defined  
 421 car distribution  $\mathcal{D}_{\text{car}}(t)$ , see also Table 1.

- 422 • **User Profiles** Additionally to the physical properties, the charging demand is a result from  
 423 the habits of the car owner, encompassing a duration of stay  $\Delta t_{\text{remain}}$ , the number of units  
 424 of power to be charged  $\Delta E$ , the SoC upon arrival  $\text{SoC}_0$  and the user preference  $u$ , indicating  
 425 whether a user is time-sensitive (will leave iff  $\Delta t_{\text{remain}} = 0$ ), or charge sensitive (will leave  
 426 iff  $\Delta E_{\text{remain}} = 0$ ). The user profiles are sampled from a distribution  $\mathcal{D}_{\text{user}}(t)$ , see also  
 427 Table 1.

428 **Endogenous state space** The endogenous state consists of the state of the various charging ports  
 429 and their connected cars, and the station battery. For each charging port, we assume a fixed voltage  
 430 and allow the actual power drawn to be regulated by controlling the current  $I_{\text{drawn}}(t) \in [0, I_{\text{max}}]$   
 431 (Orfanoudakis et al., 2024). We assume that the voltage value already encodes the phases, i.e. it  
 432 represents the product  $V \cdot \sqrt{\phi}$  in Orfanoudakis et al. (2024), eliminating the need for the phase as an  
 433 additional variable. To incorporate losses during the (dis)charging process, each EVSE is equipped  
 434 with an efficiency coefficient for charging and discharging. As a charging port may not always be  
 435 occupied, we add a final Boolean to the state  $\mathbb{1}_{\text{occup}}$ , indicating the presence of a car.

436 To properly facilitate discharging, the charging station is equipped with a battery. Similarly to EVSEs,  
 437 the battery will have a fixed voltage  $V_{\text{battery}}$ , with the power flow controlled by the current  $I_{\text{battery}}(t)$ .  
 438 To specify the physical properties of the battery, it also has a maximum capacity  $C$ , the maximal  
 439 charging rate for a car  $\bar{r}$  and  $\tau \in (0, 1)$ . Additionally, we will equip the state with the current SoC of  
 440 the battery

$$s_{\text{battery}} = (I_{\text{battery}}(t), \text{SoC}_{\text{battery}}(t), \hat{r}_{\text{battery}}(t)).$$

441 **Car State** Additionally, the state of each charging port contains information for the connected cars,  
 442 the so-called car state, representing the car that is charging at this port (all zeros if no car is present).  
 443 As this car state consists of exogenous and endogenous variables, it is listed separately. This includes  
 444 the car’s state-of-charge ( $\text{SoC} \in [0, 1]$ ), the remaining required power  $\Delta E_{\text{remain}} \in \mathbb{R}_{\geq 0}$ , the number  
 445 of timesteps the car remains  $\Delta t_{\text{remain}} \in \mathbb{N}$ , and the maximal charging power currently allowed by

446 the car  $\hat{r}(t) \in \mathbb{R}_{\geq 0}$ . The latter one is heavily depended on the State of Charge  $\text{SoC}(t) \in [0, 1]$  of  
 447 the car battery (W6lzel et al., 2021; Fastned, 2025), which is also part of the car-state. The car-state  
 448 also contains information about the physical properties of the car. These are the maximum battery  
 449 capacity  $C$ , the maximum charging rate for a car  $\bar{r}$ , and  $\tau \in (0, 1)$  – the transition point from the  
 450 bulk stage to the absorption stage of the charging process (Lee et al., 2020b). Finally, the car-state  
 451 includes a user preference indicator  $u$ .

452 In brief summary, the state of each charging port is represented by:

- 453 • Current power drawn  $I_{\text{drawn}}(t) \in \mathbb{R}_{\geq 0}$ , occupancy indicator  $\mathbb{1}_{\text{occup}}(t) \in \{0, 1\}$ ;
- 454 • Car-state  $(\Delta E_{\text{remain}}(t), \Delta t_{\text{remain}}(t), \hat{r}(t), \text{SoC}(t), C, \bar{r}, \tau, u)$ .

## 455 A.2 Transition Function

456 The transition function consists of four major steps: (i) Apply Actions, i.e. adapt charging levels at  
 457 each EVSE, (ii) charge stationed cars, (iii) departure of cars, and (iv) arrival of new cars.

458 **Apply Actions** As a first step, the actions taken by the agent are applied to adjust the power drawn  
 459 by each charging pole, specifically

$$I_{\text{drawn},i}(t) = \begin{cases} \min(I_{\text{drawn},i}(t - \Delta t) + a_i(t), \hat{r}(t), I_{\text{max} \rightarrow, i}) & I_{\text{drawn},i}(t - \Delta t) + a_i(t) \geq 0 \\ -\min(-I_{\text{drawn},i}(t - \Delta t) - a_i(t), \hat{r}(t), I_{\text{max} \leftarrow, i}) & \text{else.} \end{cases}$$

460 Hereby constraints on the maximum power drawn imposed by the architecture are enforced by  
 461 assuring that for each subtree  $H$  in the architecture, the constraints

$$\frac{1}{\eta_H} \sum_{h \in \text{leaves}(H)} I_{\text{drawn},h}(t) \leq I_H, \quad (4)$$

462 are satisfied. If the drawn currents violate these constraints, the currents at each leaf are rescaled to  
 463 satisfy the constraints, modelling the potential behaviour of some safety infrastructure on top of the  
 464 controller.

465 **Charge Stationed Cars** After having adjusted the power levels at each charging pole, the charging  
 466 is processed for the time interval, where a constant charging rate over the full interval  $\Delta t$  is assumed.  
 467 The car states are adjusted in the following way:

$$\begin{aligned} \Delta E_{\text{remain},i}(t + \Delta t) &= \Delta E_{\text{remain},i}(t) - \Delta t \cdot V_i \cdot I_{\text{drawn},i}(t) \\ \text{SoC}(t + \Delta t) &= \text{SoC}(t) + \frac{\Delta t \cdot V_i \cdot I_{\text{drawn},i}(t)}{C_i} \\ \hat{r}(t + \Delta t) &= \hat{r}_{\tau_i, \bar{r}_i}(\text{SoC}(t + \Delta t)). \end{aligned}$$

468 Notably, the physical attributes of the car in the car state, i.e. the maximum battery capacity, the  
 469 maximal charging rate and  $\tau$  do not change. As charging has been proceed, we assume that time  
 470 moves on, i.e.  $t \mapsto t + \Delta t$  and  $\Delta t_{\text{remain},i}(t + \Delta t) = \Delta t_{\text{remain},i}(t) - \Delta t$ .

471 **Departure of Cars** At the end of the period, cars fully charged or with no time remaining will  
 472 leave. Consequently the car-states for the corresponding charging poles are updated

$$s_{c,i}(t) = \begin{cases} (0, \dots, 0) & \Delta t_{\text{remain},i}(t) = 0 \text{ and } u_i = 0 \\ (0, \dots, 0) & \Delta E_{\text{remain},i}(t) = 0 \text{ and } u_i = 1 \\ s_{c,i}(t) & \text{else.} \end{cases}$$

473 **Arrival of new Cars** The amount of arriving cars is sampled  $M(t) \sim \mathcal{D}_{\text{arrival}}(t)$ . We model a first-  
 474 come-first-served policy by clipping  $M(t)$  by the number of available free spots  $N - \sum_{i=1}^N \mathbb{1}_{\text{occup},i}(t)$ .  
 475 For each car  $j = 1, \dots, M(t)$  the car profile, and the user profile are sampled from their respective  
 476 distribution, i.e.  $(\Delta t_{\text{remain},j}, \Delta E_j, \text{SoC}_{0,j}, u_j) \sim \mathcal{D}_{\text{profile}}(t)$  and  $(\bar{r}_j, \tau_j, C_j) \sim \mathcal{D}_{\text{car}}(t)$ , respec-  
 477 tively.

478 Each car  $j$  is then allocated to a free charging pole  $k$ , which alters the state of charging pole  $k$  based  
 479 on car  $j$ :

$$s_{c,k}(t) = (0, 1, \Delta E_j, \Delta t_{\text{remain},j}, \hat{r}_{\tau_j, \bar{r}_j}(\text{SoC}_{0,j}), C_j, \bar{r}_j, \tau_j, u_j).$$

### 480 A.3 Reward functions

481 The amount of net energy transferred into cars in the interval  $[t, t + \Delta t]$  can be calculated as  
 482  $\Delta E_{\text{net}}(t) = \Delta t \sum_{i=1}^N V_i \cdot I_{\text{drawn},i}(t)$ . Accounting for losses within the electric architecture of the  
 483 charging station, the amount of energy, that is transferred from cars into the grid can be calculated as

$$\Delta E_{\rightarrow \text{grid}}(t) = \Delta t \sum_{i: I_{\text{drawn},i} < 0} \eta_i \cdot V_i \cdot I_{\text{drawn},i}(t) < 0. \quad (5)$$

484 Similarly, the amount of energy that has to be drawn from the net to transfer set levels of energy into  
 485 cars  $\Delta E_{\text{grid} \rightarrow}(t)$ , after incorporating imperfect efficiencies, can be calculated via  $\Delta E_{\text{grid} \rightarrow}(t) =$   
 486  $\Delta t \sum_{i: I_{\text{drawn},i} > 0} \eta_i^{-1} \cdot V_i \cdot I_{\text{drawn},i}(t) > 0$ . Lastly, the energy contributed by (dis-)charging the  
 487 battery  $\Delta E_{\text{b},\text{net}}(t) = \Delta t I_{\text{battery}}(t) V_{\text{battery}}$  has to be incorporated, resulting in the following net  
 488 energy drawn from (or pushed into) the grid

$$\Delta E_{\text{grid},\text{net}} = \Delta E_{\text{grid} \rightarrow}(t) + \Delta E_{\rightarrow \text{grid}}(t) + \Delta E_{\text{b},\text{net}}(t).$$

489 Further that the price at which we sell and buy power from car owners is the same, i.e.  $p_{\text{sell}}$ . This  
 490 results in the following revenue

$$\Pi(t) = \begin{cases} p_{\text{sell}}(t) \cdot \Delta E_{\text{net}}(t) - p_{\text{buy}}(t) \cdot \Delta E_{\text{grid},\text{net}} - c_{\Delta t} & \Delta E_{\text{grid},\text{net}} > 0, \\ p_{\text{sell}}(t) \cdot \Delta E_{\text{net}}(t) - p_{\text{sell,grid}}(t) \cdot \Delta E_{\text{grid},\text{net}} - c_{\Delta t} & \Delta E_{\text{grid},\text{net}} \leq 0. \end{cases}$$

491 Here,  $c_{\Delta t}$  denotes the fixed cost for running the facility per  $\Delta t$ . The general reward  $r(s(t), a(t), s(t +$   
 492  $\Delta t))$ , abbreviated by  $r(t)$  in Chargax consists of the profit minus the linear combination of some  
 493 penalty terms

$$r(t) = \Pi(t) - \sum_c \alpha_c c(t). \quad (6)$$

494 Some examples of included penalty terms are listed below

495 • **Constraint Violations** The hard constraints imposed by the architecture in Eq. 4 could be  
 496 instead included as as soft constraints (Yeh et al., 2024) via the penalty

$$c_{\text{constraint}}(t) = \max_H \min \left( 0, \frac{1}{\eta_H} \sum_{i \in \text{leaves}(H)} I_{\text{drawn},i}(t) - I_H \right).$$

497 • **Satisfaction penalty** Users can experience dissatisfaction in two ways: Time-sensitive users  
 498 have a desired departure time and are assumed to leave at that time, regardless the SoC of  
 499 their car. To avoid customers leaving the charging station with a suboptimal SoC we propose  
 500 to incorporate a satisfaction penalty

$$c_{\text{Satisfaction},0}(t) = \sum_{i: \Delta t_{\text{remain},i}(t)=0, u_i=0} \max(0, \Delta E_{\text{remain},i}(t)).$$

501 The opposite holds for charge sensitive users, as they are expected to leave when there cars  
 502 are charged to the desired level. However, these users can be overly satisfied by charging  
 503 their car to the desired level faster than desired

$$c_{\text{Satisfaction},1}(t) = \sum_{i: \Delta E_i(t)=0, u_i=1} \max(0, -\Delta t_{\text{remain},i}(t)) - \beta \max(0, \Delta t_{\text{remain},i}(t)).$$

504 Here  $\beta$  controls how much the positive satisfaction from leaving earlier should weight in  
 505 comparison to the negative dissatisfaction from having to stay overtime.

506 • **Sustainability** To enforce the agent to charge cars in the most sustainable way possible, a  
 507 penalty term for non-sustainable behaviour may be added. One solution proposed in (Yeh  
 508 et al., 2024) is to employ the MOER  $m(t)$ , capturing the carbon intensity of a unit of energy  
 509 produced at time  $t$

$$c_{\text{sustain}}(t) = m(t) \cdot \Delta E_{\text{grid},\text{net}}(t).$$

510 • **Rejected Customers** In view of congestion management problems (Zhang et al., 2019;  
 511 Hussain et al., 2022), one might be interested in serving the maximum number of cars, i.e.  
 512 reduce the amount of rejected cars, by adding a penalty term for declined cars

$$c_{\text{declined}}(t) = \max \left( M(t) - \left( N - \sum_{i=1}^N \mathbb{1}_{\text{occup},i}(t) \right), 0 \right).$$

513 • **Battery Degradation** Real world batteries suffer from degradation under use (Lee et al.,  
 514 2020a). This can be incorporated by adding a degradation cost to every discharging of the  
 515 Charging station battery, as well as for the cars. For sake of simplicity, we assume the  
 516 additional degradation to be proportional to the discharged energy

$$c_{\text{degrad,battery}}(t) = |\Delta E_{\text{b,net}}(t)| \cdot \mathbb{1}_{\{\Delta E_{\text{b,net}}(t) < 0\}} \text{ and } c_{\text{degrad,cars}}(t) = |\Delta E_{\rightarrow \text{grid}}(t)|.$$

517 • **Grid Stability** (Only applicable in a V2G scenario) If the agent can discharge cars, this can  
 518 be leveraged to stabilize the grid load (Li et al., 2021; Elma, 2020). This could be reflected  
 519 in a penalty term through an exogenous signal of the grid demand  $d_{\text{grid}}(t) \in \mathbb{R}$

$$c_{\text{grid}}(t) = |\Delta E_{\text{net}}(t) - d_{\text{grid}}(t)|.$$

## 520 B Implementation Details

### 521 B.1 Practical Considerations

522 Table 3 contains environment settings used throughout our experiments whenever not stated. Addi-  
 523 tionally, we list some practical considerations in Chargax here.

- 524 • The episode length defaults to the length of data provided for arriving cars. In our bundled  
 525 scenarios, this equals 24 hours. These bundled scenarios provide their data as average  
 526 numbers per timestep. The actual number of cars arriving is then drawn using a Poisson  
 527 distribution.
- 528 • By default, we train in a Chargax environment utilizing a method akin to exploring starts.  
 529 At environment reset, we sample a random day from the given price data and use this day’s  
 530 prices for the episode. The agent observes the current episode day and whether this is a  
 531 weekday or a workday.
- 532 • Throughout our experiments, we have used a discretised action space, setting the (user-  
 533 defined) discretization level to 10. This allows the agent to select increments as 10%, 20%,  
 534 30%, etc., up to 100% of the maximum current for each charging port.

### 535 B.2 Agent configuration

536 Unless otherwise stated, the experiments conducted in Section 5 and Appendix D trained with a PPO  
 537 agent using the hyperparameters listed in Table 3.

Hyperparameter	Value	Environment Parameter	Value
Total timesteps	1e7	Minutes per timestep $\Delta t$	5
Learning rate ( $\alpha$ )	2.5e-4 (annealed)	Discretization factor	10
Discount factor $\gamma$	0.99	Episode length	24 hours
GAE $\lambda$	0.95	Number of Chargers	16
Max grad norm	100.0	Number of DC Chargers	10
Clipping coefficient $\epsilon$	0.2	Sell price to customers ( $p_{\text{sell}}$ )	0.75
Value func clip coefficient	10.0	All reward coefficients $\alpha$ (Eq. 3)	0.0
Entropy coefficient	0.01		
Value function coefficient	0.25		
Vectorized environments	12		
Rollout length (steps)	300		
Number of minibatches	4		
Update epochs	4		
Minibatch size	900		
Batch size	3600		

Table 3: PPO hyperparameters (left) alongside environment settings (right) used throughout our experiments unless otherwise stated.

## 538 C State summary

Table 4: Summary of the state space in Chargax

	symbol	domain	exogenous/ endogenous	variable name
reward data	$p_{\text{sell}}$	$\mathbb{R}_{\geq 0}$	exogenous	Selling price (to Customer) per kWh
	$p_{\text{buy}}$	$\mathbb{R}_{\geq 0}$	exogenous	Buying price per kWh
	$p_{\text{sell,grid}}$	$\mathbb{R}_{\geq 0}$	exogenous	Selling price (to grid) per kWh
	$m$	$\mathbb{R}_{\geq 0}$	exogenous	Marginal Operations Emission Rate
	$d_{\text{grid}}$	$\mathbb{R}$	exogenous	Grid Demand
	$M$	$\mathbb{N}_0$	exogenous	Number of arriving cars
Car state of EVSE i	$\Delta t_{\text{remain},i}$	$\mathbb{N}_0$	exogenous	Remaining time of customer
	$C_i$	$\mathbb{R}_{\geq 0}$	exogenous	Capacity of Car
	$\bar{r}_i$	$\mathbb{R}_{\geq 0}$	exogenous	Maximum charging rate
	$\hat{r}_i$	$\mathbb{R}_{\geq 0}$	exogenous	Maximum charging rate at current SoC
	$\tau_i$	$[0, 1]$	exogenous	
	$u_i$	$\{0, 1\}$	exogenous	User preference
	$\text{SoC}_i$	$[0, 1]$	endogenous	Current SoC
	$\Delta E_{\text{remain},i}$	$\mathbb{R}_{\geq 0}$	endogenous	Remaining Charging demand
State variables of EVSE i	$\mathbb{1}_{\text{occup},i}$	$\{0, 1\}$	endogenous	Occupancy Indicator
	$I_{\text{drawn},i}$	$\mathbb{R}_{> 0}$	endogenous	Current Power drawn at EVSE
Battery state	$I_{\text{battery}}$	$\mathbb{R}_{\geq 0}$	endogenous	Current power drawn at battery
	$\text{SoC}_{\text{battery}}$	$[0, 1]$	endogenous	SoC of Battery
	$\hat{r}_{\text{battery}}$	$\mathbb{R}_{\geq 0}$	endogenous	Maximum charging rate at current SoC

## 539 D Additional Experiments



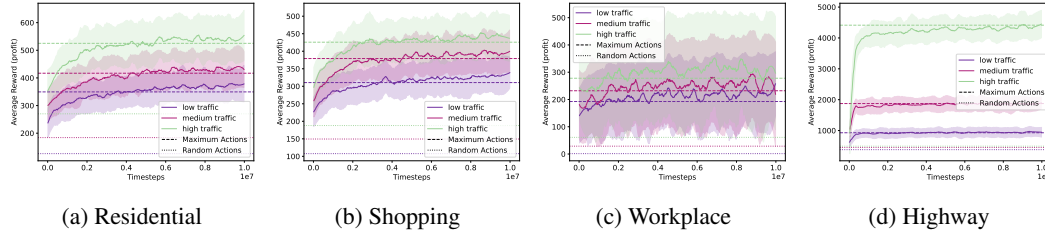


Figure 6: Results on our 4 bundled scenarios using EU cars and 16 chargers (10 DC, 5 AC)

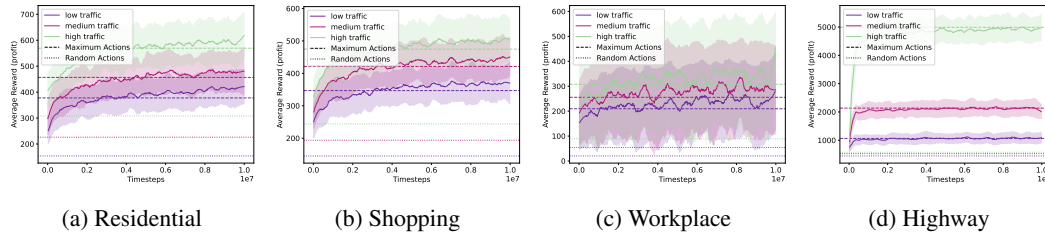


Figure 7: Results on our 4 bundled scenarios using US cars and 16 chargers (10 DC, 5 AC)

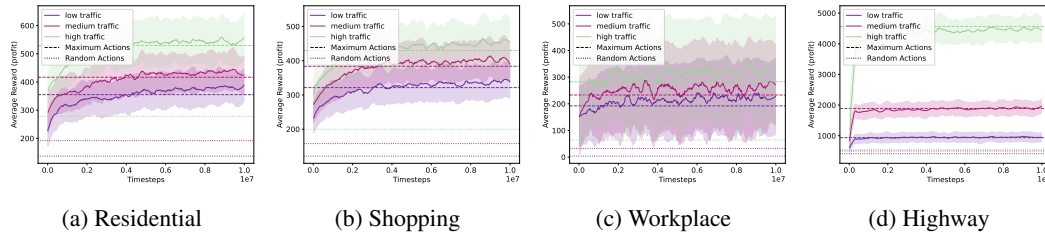


Figure 8: Results on our 4 bundled scenarios using World cars and 16 chargers (10 DC, 5 AC)

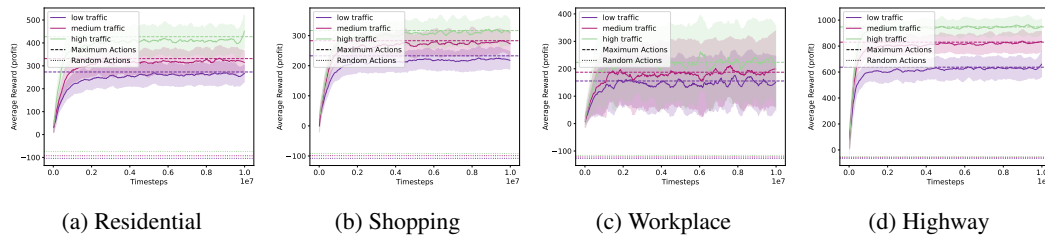


Figure 9: Results on our 4 bundled scenarios using EU cars and 16 AC (11.5kW) chargers

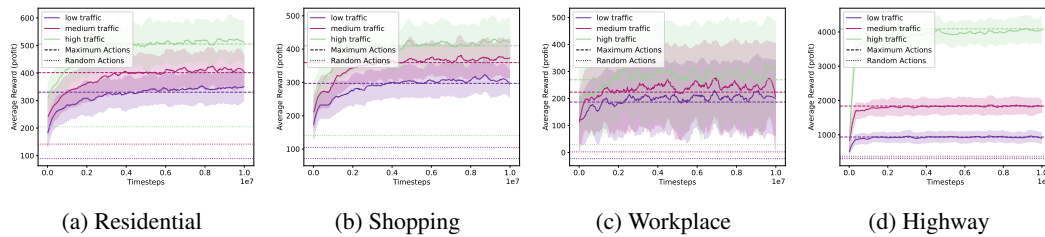


Figure 10: Results on our 4 bundled scenarios using EU cars and 8 AC (11.5kW) and 8 DC (150kW) chargers

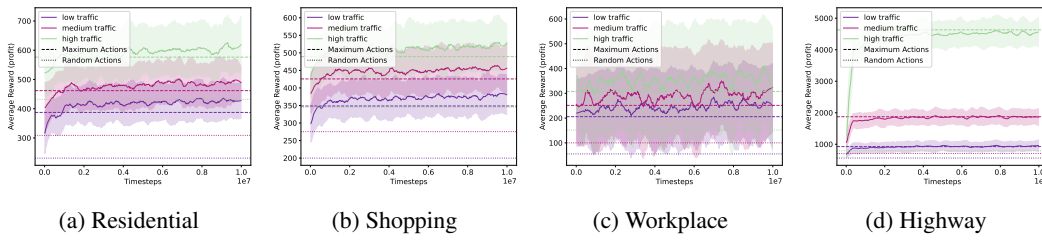


Figure 11: Results on our 4 bundled scenarios using EU cars and 16 DC (150kW) chargers