
Off-the-Grid MARL: Datasets with Baselines for Offline Multi-Agent Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Being able to harness the power of large datasets for developing cooperative multi-
2 agent controllers promises to unlock enormous value for real-world applications.
3 Many important industrial systems are multi-agent in nature and are difficult to
4 model using bespoke simulators. However, in industry, distributed processes can
5 often be recorded during operation, and large quantities of demonstrative data
6 stored. Offline multi-agent reinforcement learning (MARL) provides a promis-
7 ing paradigm for building effective decentralised controllers from such datasets.
8 However, offline MARL is still in its infancy and therefore lacks standardised
9 benchmark datasets and baselines typically found in more mature subfields of
10 reinforcement learning (RL). These deficiencies make it difficult for the community
11 to sensibly measure progress. In this work, we aim to fill this gap by releasing
12 *off-the-grid MARL (OG-MARL)*: a growing repository of high-quality datasets with
13 baselines for cooperative offline MARL research. Our datasets provide settings that
14 are characteristic of real-world systems, including complex environment dynamics,
15 heterogeneous agents, non-stationarity, many agents, partial observability, subopti-
16 mality, sparse rewards and demonstrated coordination. For each setting, we provide
17 a range of different dataset types (e.g. Good, Medium, Poor, and Replay) and
18 profile the composition of experiences for each dataset. We hope that OG-MARL
19 will serve the community as a reliable source of datasets and help drive progress,
20 while also providing an accessible entry point for researchers new to the field.

21 1 Introduction

22 RL algorithms typically require extensive online interactions with an environment to be able to learn
23 robust policies (Yu, 2018). This limits the extent to which previously-recorded experience may be
24 leveraged for RL applications, forcing practitioners to instead rely heavily on optimised environment
25 simulators that are able to run quickly and in parallel on modern compute hardware.

26 In a simulation, it is not atypical to be able to generate years of operating behaviour of a specific
27 system (Berner et al., 2019; Vinyals et al., 2019). However, achieving this level of online data
28 generation throughput in real-world systems, where a realistic simulator is not readily available, can
29 be challenging or near impossible. More recently, the field of offline RL has offered a solution to
30 this challenge by bridging the gap between RL and supervised learning. In offline RL, the aim is to
31 develop algorithms that are able to leverage large existing datasets of sequential decision-making to
32 learn optimal control strategies that can be deployed online (Levine et al., 2020). Many researchers

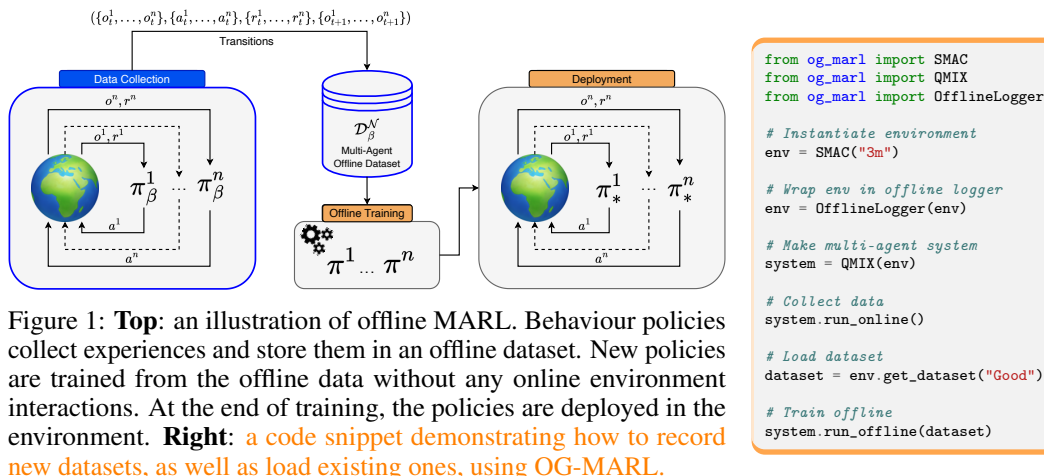


Figure 1: **Top**: an illustration of offline MARL. Behaviour policies collect experiences and store them in an offline dataset. New policies are trained from the offline data without any online environment interactions. At the end of training, the policies are deployed in the environment. **Right**: a code snippet demonstrating how to record new datasets, as well as load existing ones, using OG-MARL.

33 believe that offline RL could help unlock the full potential of RL when applied to the real world,
 34 where success has been limited (Dulac-Arnold et al., 2021).

35 Although the field of offline RL has experienced a surge in research interest in recent years (Prudencio
 36 et al., 2023), the focus on offline approaches specific to the multi-agent setting has remained relatively
 37 neglected, despite the fact that many real-world problems are naturally formulated as multi-agent
 38 systems (e.g. traffic management (Zhang et al., 2019), a fleet of ride-sharing vehicles (Sykora et al.,
 39 2020), a network of trains (Mohanty et al., 2020) or electricity grid management (Khattar and Jin,
 40 2022)). Moreover, systems that require multiple agents (programmed and/or human) to execute
 41 coordinated strategies to perform optimally, arguably have a higher barrier to entry when it comes to
 42 creating bespoke simulators to model their online operating behaviour.

43 Offline RL research in the single agent setting has benefited greatly from publicly available datasets
 44 and benchmarks such as D4RL (Fu et al., 2020) and RL Unplugged (Gulcehre et al., 2020). Without
 45 such offerings in the multi-agent setting to help standardise research efforts and evaluation, it remains
 46 challenging to gauge the state of the field and reproduce results from previous work. Ultimately, to
 47 develop new ideas that drive the field forward, standardised sets of tasks and baselines are required.

48 In this paper, we present OG-MARL, a rich set of datasets specifically curated for cooperative offline
 49 MARL. We generated diverse datasets on a range of popular cooperative MARL environments. For
 50 each environment, we provide different types of behaviour resulting in *Good*, *Medium* and *Poor*
 51 datasets as well as *Replay* datasets (a mixture of the previous three). We developed and applied a
 52 quality assurance methodology to validate our datasets to ensure that they contain a diverse spread
 53 of experiences. Together with our datasets, we provide initial baseline results using state-of-the-art
 54 offline MARL algorithms.

55 The OG-MARL code and datasets are publicly available through our website.¹ Additionally, we
 56 invite the community to contribute their own datasets to the growing repository on OG-MARL and
 57 use our website as a platform for storing and distributing datasets for the benefit of the research
 58 community. We hope the lessons contained in our methodology for generating and validating datasets
 59 help future researchers to produce high-quality offline MARL datasets and help drive progress.

60 2 Related Work

61 **Datasets.** In the single-agent RL setting, D4RL (Fu et al., 2020) and RL Unplugged (Gulcehre
 62 et al., 2020) have been important contributions, providing a comprehensive set of offline datasets for
 63 benchmarking offline RL algorithms. While not originally included, D4RL was later extended by Lu
 64 et al. (2022) to incorporate datasets with pixel-based observations, which they highlight as a notable

¹<https://sites.google.com/view/og-marl>

65 deficiency of other datasets. The ease of access to high-quality datasets provided by D4RL and RL
 66 Unplugged has enabled the field of offline RL to make rapid progress over the past years (Kostrikov
 67 et al., 2021; Ghasemipour et al., 2022; Nakamoto et al., 2023). However, these repositories lack
 68 datasets for MARL, which we believe, alongside additional technical difficulties such as large joint
 69 action spaces (Yang et al., 2021), has resulted in slower progress in the field.

70 **Offline Multi-Agent Reinforcement Learning.** To date, there has been a limited number of papers
 71 published on cooperative offline MARL, resulting in benchmarks, datasets and algorithms that do
 72 not adhere to any unified standard, making comparisons between works difficult. In brief, Zhang
 73 et al. (2021) carried out an in-depth theoretical analysis of finite-sample offline MARL. Jiang and
 74 Lu (2021) proposed a decentralised multi-agent version of the popular offline RL algorithm BCQ
 75 (Fujimoto et al., 2019) and evaluated it on their own datasets of a multi-agent version of MuJoCo
 76 (MAMuJoCo) (Peng et al., 2021). Yang et al. (2021) highlighted how extrapolation error accumulates
 77 rapidly in the number of agents and propose a new method they call *Implicit Constraint Q-Learning*
 78 (ICQ) to address this. The authors evaluate their method on their own datasets collected using the
 79 popular *StarCraft Multi-Agent Challenge* (SMAC) (Samvelyan et al., 2019). Pan et al. (2022) showed
 80 that *Conservative Q-Learning* (CQL) (Kumar et al., 2020), a very successful offline RL method,
 81 does not transfer well to the multi-agent setting since the multi-agent policy gradients are prone to
 82 uncoordinated local optima. To overcome this, the authors proposed a zeroth-order optimization
 83 method to better optimize the conservative value functions, and evaluate their method on their own
 84 datasets of a handful of SMAC scenarios, the two agent HalfCheetah scenario from MAMuJoCo and
 85 some simple Multi Particle Environments (MPE) (Lowe et al., 2017). Meng et al. (2021) propose a
 86 *multi-agent decision transformer* (MADT) architecture, which builds on the *decision transformer*
 87 (DT) (Chen et al., 2021), and demonstrated how it can be used for offline pre-training and online
 88 fine-tuning in MARL by evaluating their method on their own SMAC datasets. Barde et al. (2023)
 89 explored a model-based approach for offline MARL and evaluated their method on MAMuJoCo.

90 **Datasets and baselines for Offline MARL.** In all of the aforementioned works, the authors generate
 91 their own datasets for their experiments and provide only a limited amount of information about the
 92 composition of their datasets (e.g. spread of episode returns and/or visualisations of the behaviour
 93 policy). Furthermore, each paper proposes a novel algorithm and typically compares their proposal to
 94 a set of baselines specifically implemented for their work. The lack of commonly shared benchmark
 95 datasets and baselines among previous papers has made it difficult to compare the relative strengths
 96 and weaknesses of these contributions and is one of the key motivations for our work.

97 Finally, we note works that have already made use of the pre-release version of OG-MARL. Formanek
 98 et al. (2023) investigated selective “reincarnation” in the multi-agent setting and Zhu et al. (2023)
 99 explored using diffusion models to learn policies in offline MARL. Both these works made use of
 100 OG-MARL datasets for their experiments, which allows for easier reproducibility and more sound
 101 comparison with future work using OG-MARL.

102 3 Preliminaries

103 **Multi-Agent Reinforcement Learning.** There are three main formulations of MARL tasks: com-
 104 petitive, cooperative and mixed. The focus of this work is on the cooperative setting. Cooperative
 105 MARL can be formulated as a *decentralised partially observable Markov decision process* (Dec-
 106 POMDP) (Bernstein et al., 2002). A Dec-POMDP consists of a tuple $\mathcal{M} = (\mathcal{N}, \mathcal{S}, \{\mathcal{A}^i\}, \{\mathcal{O}^i\}, P,$
 107 $E, \rho_0, r, \gamma)$, where $\mathcal{N} \equiv \{1, \dots, n\}$ is the set of n agents in the system and $s \in \mathcal{S}$ describes the full
 108 state of the system. The initial state distribution is given by ρ_0 . Each agent $i \in \mathcal{N}$ receives only partial
 109 information from the environment in the form of a local observation o_t^i , given according to an emission
 110 function $E(o_t^i | s_t, i)$. At each timestep t , each agent chooses an action $a_t^i \in \mathcal{A}^i$ to form a joint action
 111 $\mathbf{a}_t \in \mathcal{A} \equiv \prod_i^N \mathcal{A}^i$. Due to partial observability, each agent typically maintains an observation history
 112 $o_{0:t}^i = (o_0^i, \dots, o_t^i)$, or implicit memory, on which it conditions its policy $\mu^i(a_t^i | o_{0:t}^i)$, when choosing
 113 an action. The environment then transitions to a new state in response to the joint action selected in
 114 the current state, according to the state transition function $P(s_{t+1} | s_t, \mathbf{a}_t)$ and provides a shared scalar

115 reward to each agent according to a reward function $r(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$. We define an agent’s return
 116 as its discounted cumulative rewards over the T episode timesteps, $G = \sum_{t=0}^T \gamma^t r_t$, where $\gamma \in (0, 1]$
 117 is the discount factor. The goal of MARL in a Dec-POMDP is to find a joint policy $(\pi^1, \dots, \pi^n) \equiv \pi$
 118 such that the return of each agent i , following π^i , is maximised with respect to the other agents’
 119 policies, $\pi^{-i} \equiv (\pi \setminus \pi^i)$. That is, we aim to find π such that $\forall i : \pi^i \in \arg \max_{\hat{\pi}^i} \mathbb{E} [G \mid \hat{\pi}^i, \pi^{-i}]$

120 **Offline Reinforcement Learning.** An offline RL algorithm is trained on a static, previously collected
 121 dataset \mathcal{D}_β of transitions (o_t, a_t, r_t, o_{t+1}) from some (potentially unknown) behaviour policy π_β ,
 122 without any further online interactions. There are several well-known challenges in the offline RL
 123 setting which have been explored, predominantly in the single-agent literature. The primary issues
 124 are related to different manifestations of data distribution mismatch between the offline data and the
 125 induced online data. [Levine et al. \(2020\)](#) provide a detailed survey of the problems and solutions in
 126 offline RL.

127 **Offline Multi-Agent Reinforcement Learning.** In the multi-agent setting, offline MARL algorithms
 128 are designed to learn an optimal *joint* policy $(\pi^1, \dots, \pi^n) \equiv \pi$, from a static dataset \mathcal{D}_β^N of previously
 129 collected multi-agent transitions $(\{o_t^1, \dots, o_t^n\}, \{a_t^1, \dots, a_t^n\}, \{r_t^1, \dots, r_t^n\}, \{o_{t+1}^1, \dots, o_{t+1}^n\})$, gen-
 130 erated by a set of interacting behaviour policies $(\pi_\beta^1, \dots, \pi_\beta^n) \equiv \pi_\beta$.

131 4 Task Properties

132 In order to design an offline MARL benchmark which is maximally useful to the community, we
 133 carefully considered the properties that the environments and datasets in our benchmark should
 134 satisfy. A major drawback in most prior work has been the limited diversity in the tasks that the
 135 algorithms were evaluated on. [Meng et al. \(2021\)](#) for example only evaluated their algorithm on
 136 SMAC datasets and [Jiang and Lu \(2021\)](#) only evaluated on MAMuJoCo datasets. This makes it
 137 difficult to draw strong conclusions about the generalisability of offline MARL algorithms. Moreover,
 138 these environments fail to test the algorithms along dimensions which may be important for real-world
 139 applications. In this section, we outline the properties we believe are important for evaluating offline
 140 MARL algorithms.

141 **Centralised and Independent Training.** The environments supported in OG-MARL are designed
 142 to test algorithms that use decentralised execution, i.e. at execution time, agents need to choose
 143 actions based on their local observation histories only. However, during training, centralisation (i.e.
 144 sharing information between agents) is permissible, although not required. *Centralised training*
 145 *with decentralised execution* (CTDE) ([Kraemer and Banerjee, 2016](#)) is one of the most popular
 146 MARL paradigms and is well-suited for many real-world applications. Being able to test both
 147 centralised and independent training algorithms is important because it has been shown that neither
 148 paradigm is consistently better than the other ([Lyu et al., 2021](#)). As such, both types of algorithms
 149 can be evaluated using OG-MARL datasets and we also provide baselines for both centralised and
 150 independent training.

151 **Different types of Behaviour Policies.** We generated datasets with several different types of
 152 behaviour policies including policies trained using online MARL with fully independent learners (e.g.
 153 independent DQN and independent TD3), as well as CTDE algorithms (e.g. QMIX and MATD3).
 154 Furthermore, some datasets generated with CTDE algorithms used a state-based critic while others
 155 used a joint-observation critic. It was important for us to consider both of these critic setups as they
 156 are known to result in qualitatively different policies ([Lyu et al., 2022](#)). More specific details of which
 157 algorithms were used to generate which datasets can be found in [Table B.1](#) in the appendix.

158 **Partial Information.** It is common for agents to receive only local information about their envi-
 159 ronment, especially in real-world systems that rely on decentralised components. Thus, some of
 160 the environments in OG-MARL test an algorithm’s ability to leverage agents’ *memory* in order to
 161 choose optimal actions based only on partial information from local observations. This is in contrast
 162 to settings such as MAMuJoCo where prior methods ([Jiang and Lu, 2021](#); [Pan et al., 2022](#)) achieved
 163 reasonable results without instilling agents with any form of memory.

164 **Different Observation Modalities.** In the real world, agent observations come in many different
165 forms. For example, observations may be in the form of a feature vector or a matrix representing a
166 pixel-based visual observation. Lu et al. (2022) highlighted that prior single-agent offline RL datasets
167 failed to test algorithms on high-dimensional pixel-based observations. OG-MARL tests algorithms
168 on a diverse set of observation modalities, including feature vectors and pixel matrices of different
169 sizes.

170 **Continuous and Discrete Action Spaces.** The actions an agent is expected to take can be either
171 discrete or continuous across a diverse range of applications. Moreover, continuous action spaces
172 can often be more challenging for offline MARL algorithms as the larger action spaces make them
173 more prone to extrapolation errors, due to out-of-distribution actions. OG-MARL supports a range
174 of environments with both discrete and continuous actions.

175 **Homogeneous and Heterogeneous Agents.** Real-world systems can either be homogeneous or
176 heterogeneous in terms of the types of agents that comprise the system. In a homogeneous system,
177 it may be significantly simpler to train a single policy and copy it to all agents in the system. On
178 the other hand, in a heterogeneous system, where agents may have significantly different roles and
179 responsibilities, this approach is unlikely to succeed. OG-MARL provides datasets from environments
180 that represent both homogeneous and heterogeneous systems.

181 **Number of Agents.** Practical MARL systems may have a large number of agents. Most prior works
182 to date have evaluated their algorithms on environments with typically fewer than 8 agents (Pan et al.,
183 2022; Yang et al., 2021; Jiang and Lu, 2021). In OG-MARL, we provide datasets with between 2 and
184 27 agents, to better evaluate *large-scale* offline MARL (see Table B.1).

185 **Sparse Rewards.** Sparse rewards are challenging in the single-agent setting, but in the multi-agent
186 setting, it can be even more challenging due to the multi-agent credit assignment problem (Zhou
187 et al., 2020). Prior works focused exclusively on dense reward settings (Pan et al., 2022; Yang et al.,
188 2021). To overcome this, OG-MARL also provides datasets with sparse rewards.

189 **Team and Individual Rewards.** Some environments have team rewards while others can have an
190 additional local reward component. Team rewards exacerbate the multi-agent credit assignment
191 problem, and having a local reward component can help mitigate this. However, local rewards may
192 result in sub-optimality, where agents behave too greedily with respect to their local reward and as a
193 result jeopardize achieving the overall team objective. OG-MARL includes tasks to test algorithms
194 along both of these dimensions.

195 **Procedurally Generated and Stochastic Environments.** Some popular MARL benchmark environ-
196 ments are known to be highly deterministic (Ellis et al., 2022). This limits the extent to which the
197 generalisation capabilities of algorithms can be evaluated. Procedurally generated environments have
198 proved to be a useful tool for evaluating generalisation in single-agent RL (Cobbe et al., 2020). In
199 order to better evaluate generalisation in offline MARL, OG-MARL includes stochastic tasks that
200 make use of procedural generation.

201 **Realistic Multi-Agent Domains.** Almost all prior offline MARL works have evaluated their al-
202 gorithms exclusively on game-like environments such as StarCraft (Yang et al., 2021) and particle
203 simulators (Pan et al., 2022). Although a large subset of open research questions may still be readily
204 investigated in such simulated environments, we argue that in order for offline MARL to become
205 more practically relevant, benchmarks in the research community should begin to closer reflect real-
206 world problems of interest. Therefore, in addition to common game-like benchmark environments,
207 OG-MARL also supports environments which simulate more real-world like problems including
208 energy management and control (Vazquez-Canteli et al., 2020; Wang et al., 2021). While there
209 remains a large gap between these environments and truly real-world settings, it is a step in the right
210 direction to keep pushing the field forward and enable useful contributions in the development of new
211 algorithms and improving our understanding of key difficulties and failure modes.

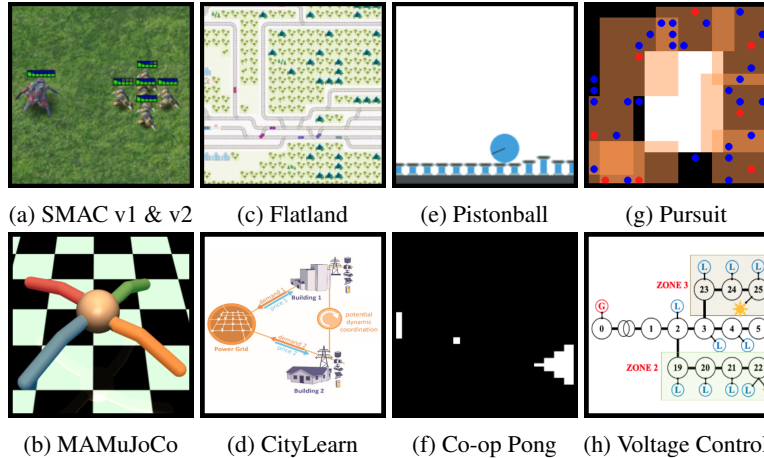


Figure 2: MARL environments for which we provide datasets in OG-MARL.

212 5 Environments

213 **SMAC v1** (*hetero- and homogeneous agents, local observations*). SMAC is the most popular
 214 cooperative offline MARL environment used in the literature (Gorsane et al., 2022). SMAC focuses
 215 on the micromanagement challenge in StarCraft 2 where each unit is controlled by an independent
 216 agent that must learn to cooperate and coordinate based on local (partial) observations. SMAC played
 217 an important role in moving the MARL research community beyond grid-world problems and has
 218 also been very popular in the offline MARL literature (Yang et al., 2021; Meng et al., 2021; Pan et al.,
 219 2022). Thus, it was important for OG-MARL to support a range of SMAC scenarios.

220 **SMAC v2** (*procedural generation, local observations*). Recently some deficiencies in SMAC have
 221 been brought to light. Most importantly, SMAC is highly deterministic, and agents can therefore
 222 learn to *memorise* the best policy by conditioning on the environment timestep only. To address this,
 223 SMACv2 (Ellis et al., 2022) was recently released and includes non-deterministic scenarios, thus
 224 providing a more challenging benchmark for MARL algorithms. In OG-MARL, we publicly release
 225 the first set of SMACv2 datasets.

226 **MAMuJoCo** (*hetero- and homogeneous agents, continuous actions*). The MuJoCo environment
 227 (Todorov et al., 2012) has been an important benchmark that helped drive research in continuous con-
 228 trol. More recently, MuJoCo has been adapted for the multi-agent setting by introducing independent
 229 agents that control different subsets of the whole MuJoCo robot (MAMuJoCo) (Peng et al., 2021).
 230 MAMuJoCo is an important benchmark because there are a limited number of continuous action
 231 space environments available to the MARL research community. MAMuJoCo has also been widely
 232 adopted in the offline MARL literature (Jiang and Lu, 2021; Pan et al., 2022). Thus, in OG-MARL
 233 we provide the largest openly available collection of offline datasets on scenarios in MAMuJoCo
 234 (Pan et al. (2022)), for example, only provided a single dataset on 2-Agent HalfCheetah).

235 **PettingZoo** (*pixel observations, discrete and continuous actions*). OpenAI’s Gym (Brockman
 236 et al., 2016) has been widely used as a benchmark for single agent RL. PettingZoo is a gym-like
 237 environment-suite for MARL (Terry et al., 2021) and provides a diverse collection of environments.
 238 In OG-MARL, we provide a general-purpose environment wrapper which can be used to generate
 239 new datasets for any PettingZoo environment. Additionally, we provide initial datasets on three
 240 PettingZoo environments including *PistonBall*, *Co-op Pong* and *Pursuit* (Gupta et al., 2017). We
 241 chose these environments because they have visual (pixel-based) observations of varying sizes; an
 242 important dimension along which prior works have failed to evaluate their algorithms.

243 **Flatland** (*real-world problem, procedural generation, sparse local rewards*). The train scheduling
 244 problem is a real-world challenge with significant practical relevance. Flatland (Mohanty et al., 2020)
 245 is a simplified 2D simulation of the train scheduling problem that is an appealing benchmark for

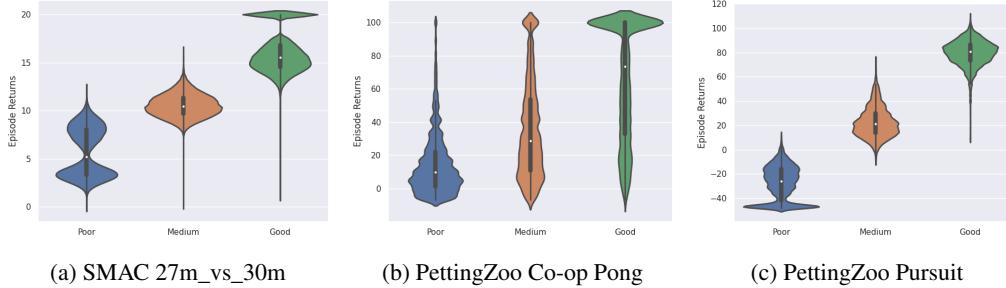


Figure 3: Violin plots of the probability distribution of episode returns for selected datasets in OG-MARL. In blue the Poor datasets, in orange the Medium datasets and in green the Good datasets. Wider sections of the violin plot represent a higher probability of sampling a trajectory with a given episode return, while the thinner sections correspond to a lower probability. The violin plots also include the median, interquartile range and min/max episode return for the datasets.

246 cooperative MARL for several reasons. Firstly, it randomly generates a new train track layout and
 247 timetable at the start of each episode, thus testing the generalisation capabilities of MARL algorithms
 248 to a greater degree than many other environments. Secondly, Flatland has a very sparse and noisy
 249 reward signal, as agents only receive a reward on the final timestep of the episode. Finally, agents
 250 have access to a local reward component. These properties make the Flatland environment a novel,
 251 challenging and realistic benchmark for offline MARL.

252 **Voltage Control and CityLearn** (*real-world problem, continuous actions*). Energy management (Yu
 253 et al., 2021) is another appealing real-world application for MARL, especially given the large potential
 254 efficiency gains and corresponding positive effects on climate change that could be had (Rolnick
 255 et al., 2022). As such, we provide datasets for two challenging MARL environments related to energy
 256 management. Firstly, we provide datasets for the *Active Voltage Control on Power Distribution Net-*
 257 *works* environment (Wang et al., 2021). Secondly, we provide datasets for the CityLearn environment
 258 (Vazquez-Canteli et al., 2020) where the goal is to develop agents for distributed energy resource
 259 management and demand response between a network of buildings with batteries and photovoltaics.

260 6 Datasets

261 To generate the transitions in the datasets, we recorded environment interactions of partially trained
 262 online algorithms, as has been common in prior works for both single-agent (Gulcehre et al., 2020)
 263 and multi-agent settings (Yang et al., 2021; Pan et al., 2022). For discrete action environments, we
 264 used QMIX (Rashid et al., 2018) and independent DQN and for continuous action environments,
 265 we used independent TD3 (Fujimoto et al., 2018) and MATD3 (Lowe et al., 2017; Ackermann et al.,
 266 2019). Additional details about how each dataset was generated are included in Appendix C.

267 **Diverse Data Distributions.** It is well known from the single-agent offline RL literature that the
 268 quality of experience in offline datasets can play a large role in the final performance of offline RL
 269 algorithms (Fu et al., 2020). In OG-MARL, we include a range of dataset distributions including
 270 Good, Medium, Poor and Replay datasets in order to benchmark offline MARL algorithms on a
 271 range of different dataset qualities. The dataset types are characterised by the quality of the joint
 272 policy that generated the trajectories in the dataset, which is the same approach taken in previous
 273 works (Meng et al., 2021; Yang et al., 2021; Pan et al., 2022). To ensure that all of our datasets have
 274 sufficient coverage of the state and action spaces, while also containing minimal repetition i.e. not
 275 being too narrowly focused around a single strategy, we used 3 independently trained joint policies
 276 to generate each dataset, and additionally added a small amount of exploration noise to the policies.
 277 The boundaries for the different categories were assigned independently for each environment and
 278 were related to the maximum attainable return in the environment. Additional details about how the
 279 different datasets were curated can be found in Appendix C.

Table 1: Results on the *Pursuit* and *Co-op Pong* datasets. The mean episode return with one standard deviation across all seeds is given. In each row the best mean episode return is in bold.

Scenario	Dataset	BC	QMIX	QMIX+BCQ	QMIX+CQL	MAICQ
Co-op Pong	Good	31.2±3.5	0.6±3.5	1.9±1.1	90.0±4.7	75.4±3.9
	Medium	21.6±4.8	10.6±17.6	20.3±12.2	64.9±15.0	84.6±0.9
	Poor	1.0±0.9	14.4±16.0	30.2±20.7	52.7±8.5	74.8±7.8
Pursuit	Good	78.3±1.8	6.7±19.0	66.9±14.0	54.4±6.3	92.7±3.7
	Medium	15.0±1.6	-24.4±20.2	16.6±10.7	20.6±10.3	35.3±3.0
	Poor	-18.5±1.6	-43.7±5.6	-0.7±4.0	-19.6±3.3	-4.1±0.7

280 **Statistical characterisation of datasets.** It is common in both the single-agent and multi-agent
 281 offline RL literature for researchers to curate offline datasets by unrolling episodes using an RL policy
 282 that was trained to a desired *mean* episode return. However, authors seldom report the distribution
 283 of episode returns induced by the policy. Reporting only the mean episode return of the behaviour
 284 policy can be misleading (Agarwal et al., 2021). To address this, we provide violin plots to visualise
 285 the distribution of expected episode returns. A violin plot is a powerful tool for visualising numerical
 286 distributions as they visualise the density of the distribution as well as several summary statistics
 287 such as the minimum, maximum and interquartile range of the data. These properties make the violin
 288 plot very useful for understanding the distribution of episode returns in the offline datasets, assisting
 289 with interpreting offline MARL results. Figure 3 provides a sample of the violin plots for different
 290 scenarios (the remainder of the plots can be found in the appendix). In each figure, the difference
 291 in shape and position of the three violins (blue, orange and green) illustrates the difference in the
 292 datasets with respect to the expected episode return. Additionally, we provide a table with the mean
 293 and standard deviation of the episode returns for each of the datasets in Table C.1, similar to Meng
 294 et al. (2021).

295 7 Baselines

296 In this section, we present the initial baselines that we provide with OG-MARL. This serves two
 297 purposes: *i*) to validate the quality of our datasets and *ii*) to enable the community to use these initial
 298 results for development and performance comparisons in future work. In the main text, we present
 299 results on two PettingZoo environments (*Pursuit* and *Co-op Pong*), since these environments and
 300 their corresponding datasets are a novel benchmark for offline MARL. Furthermore, it is the first set
 301 of environments with pixel-based observations to be used to evaluate offline MARL algorithms. We
 302 include all additional baseline results in Appendix D (Table D.4 and Table D.5).

303 **Baseline Algorithms.** State-of-the-art algorithms were implemented from seminal offline MARL
 304 work. For discrete action environments we implemented *Behaviour Cloning* (BC), QMIX (Rashid
 305 et al., 2018), QMIX with *Batch Constrained Q-Learning* (Fujimoto et al., 2019) (QMIX+BCQ),
 306 QMIX with *Conservative Q-Learning* (Kumar et al., 2020) (QMIX+CQL) and MAICQ (Yang et al.,
 307 2021). For continuous action environments, Behaviour Cloning (BC), Independent TD3 (ITD3), ITD3
 308 with *Behaviour Cloning* regularisation (Fujimoto and Gu, 2021) (ITD3+BC), ITD3 with *Conservative*
 309 *Q-Learning* (ITD3+CQL) and OMAR (Pan et al., 2022) were implemented. Appendix D provides
 310 additional implementation details on the baseline algorithms.

311 **Experimental Setup.** On *Pursuit* and *Co-op Pong*, all of the algorithms were trained offline for 50000
 312 training steps with a fixed batch size of 32. At the end of training, we evaluated the performance of
 313 the algorithms by unrolling the final joint policy in the environment for 100 episodes and recording
 314 the mean episode return over the episodes. We repeated this procedure for 10 independent seeds as
 315 per the recommendation by Gorsane et al. (2022). We kept the online evaluation budget (Kurenkov
 316 and Kolesnikov, 2022) fixed for all algorithms by only tuning hyper-parameters on *Co-op Pong*
 317 and keeping them fixed for *Pursuit*. Controlling for the online evaluation budget is important when
 318 comparing offline algorithms because online evaluation may be expensive, slow or dangerous in

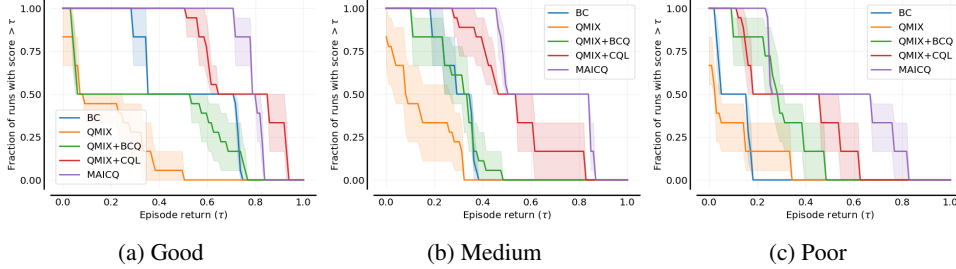


Figure 4: Performance profiles (Agarwal et al., 2021) aggregated across all seeds on *Pursuit* and *Co-op Pong*. Shaded regions show pointwise 95% confidence bands based on percentile bootstrap with stratified sampling.

319 real-world problems, making online hyper-parameter fine-tuning infeasible. See Appendix D for a
 320 further discussion on hyper-parameter tuning in OG-MARL.

321 **Results.** In Table 1 we provide the unnormalised mean episode returns for each of the discrete action
 322 algorithms on the different datasets for *Pursuit* and *Co-op Pong*.

323 **Aggregated Results.** In addition to the tabulated results we also provide *aggregated* results as per
 324 the recommendation by Gorsane et al. (2022). In Figure 4 we plot the performance profiles (Agarwal
 325 et al., 2021) of the discrete action algorithms by aggregating across all seeds and the two environments,
 326 *Pursuit* and *Co-op Pong*. To facilitate aggregation across environments, where the possible episode
 327 returns can be very different, we adopt the normalisation procedure from Fu et al. (2020). On the
 328 Good datasets, we found that MAICQ and QMIX+CQL both outperformed behaviour cloning (BC).
 329 QMIX+BCQ did not outperform BC and vanilla QMIX performed very poorly. On the Medium
 330 datasets, MAICQ and QMIX+CQL once again performed the best, significantly outperforming BC.
 331 QMIX+BCQ marginally outperformed BC and vanilla QMIX failed. Finally, on the Poor datasets,
 332 MAICQ, QMIX+CQL and QMIX+BCQ all outperformed BC but MAICQ was the best by some
 333 margin. These results on PettingZoo environments, with pixel observations, further substantiate that
 334 MAICQ is the current state-of-the-art offline MARL algorithm in discrete action settings.

335 8 Discussion

336 **Limitations and future work.** The primary limitation of this work is that it focuses on the cooperative
 337 setting. Additionally, the datasets used in OG-MARL were exclusively generated by online MARL
 338 policies. Future work could explore the inclusion of datasets from alternate sources, such as hand-
 339 designed or human controllers, which may exhibit distinct properties (Fu et al., 2020). Moreover,
 340 an exciting research direction considers the offline RL problem as a sequence modeling task (Chen
 341 et al., 2021; Meng et al., 2021), and we aim to incorporate such models as additional baselines in
 342 OG-MARL in future iterations.

343 **Potential Negative Societal Impacts.** While the potential positive impacts of efficient decentralized
 344 controllers powered by offline MARL are promising, it is essential to acknowledge and address the
 345 potential negative societal impacts (Whittlestone et al., 2021). Deploying a model trained using
 346 offline MARL in real-world applications requires careful consideration of safety measures (Gu et al.,
 347 2022; Xu et al., 2022). Practitioners should exercise caution to ensure the safe and responsible
 348 implementation of such models.

349 **Conclusion.** In this work, we highlighted the importance of offline MARL as a research direction
 350 for applying RL to real-world problems. We specifically focused on the lack of a standard set of
 351 benchmark datasets, which is a significant obstacle to progress. To address this issue, we presented
 352 a set of relevant and diverse datasets for offline MARL. We profiled our datasets by visualising
 353 the distribution of episode returns in violin plots and tabulated mean and standard deviations. We
 354 validated our datasets by providing a set of initial baseline results with state-of-the-art offline MARL

355 algorithms. Finally, we open-sourced all of our software tooling for generating new datasets and
356 provided a website with our code, as well as for hosting and sharing the datasets. It is our hope that
357 the research community will adopt and contribute towards OG-MARL as a framework for offline
358 MARL research and that it helps to drive progress in this nascent field.

359 **References**

- 360 J. Ackermann, V. Gabler, T. Osa, and M. Sugiyama. Reducing overestimation bias in multi-agent
361 domains using double centralized critics. *ArXiv Preprint*, 2019. 7
- 362 R. Agarwal, D. Schuurmans, and M. Norouzi. An optimistic perspective on offline reinforcement
363 learning. *ArXiv Preprint*, 2019. 23
- 364 R. Agarwal, M. Schwarzzer, P. S. Castro, A. C. Courville, and M. Bellemare. Deep reinforcement
365 learning at the edge of the statistical precipice. *Advances in Neural Information Processing Systems*,
366 2021. 8, 9, 26, 27
- 367 P. Barde, J. Foerster, D. Nowrouzezahrai, and A. Zhang. A model-based solution to the offline
368 multi-agent reinforcement learning coordination problem. *ArXiv Preprint*, 2023. 3
- 369 C. Berner, G. Brockman, B. Chan, V. Cheung, P. Debiak, C. Dennison, D. Farhi, Q. Fischer,
370 S. Hashme, C. Hesse, R. Józefowicz, S. Gray, C. Olsson, J. Pachocki, M. Petrov, H. P.
371 de Oliveira Pinto, J. Raiman, T. Salimans, J. Schlatter, J. Schneider, S. Sidor, I. Sutskever, J. Tang,
372 F. Wolski, and S. Zhang. Dota 2 with large scale deep reinforcement learning. *ArXiv Preprint*,
373 2019. 1
- 374 D. S. Bernstein, R. Givan, N. Immerman, and S. Zilberstein. The complexity of decentralized control
375 of markov decision processes. *Mathematics of operations research*, 2002. 3
- 376 G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. Openai
377 gym. *ArXiv Preprint*, 2016. 6
- 378 L. Chen, K. Lu, A. Rajeswaran, K. Lee, A. Grover, M. Laskin, P. Abbeel, A. Srinivas, and I. Mor-
379 datch. Decision transformer: reinforcement learning via sequence modeling. *Advances in Neural*
380 *Information Processing Systems*, 2021. 3, 9
- 381 K. Cobbe, C. Hesse, J. Hilton, and J. Schulman. Leveraging procedural generation to benchmark
382 reinforcement learning. *International Conference on Machine Learning*, 2020. 5
- 383 G. Dulac-Arnold, N. Levine, D. J. Mankowitz, J. Li, C. Paduraru, S. Gowal, and T. Hester. Challenges
384 of real-world reinforcement learning: definitions, benchmarks and analysis. *Springer Machine*
385 *Learning*, 2021. 2
- 386 B. Ellis, S. Moalla, M. Samvelyan, M. Sun, A. Mahajan, J. N. Foerster, and S. Whiteson. Smacv2:
387 An improved benchmark for cooperative multi-agent reinforcement learning. *ArXiv Preprint*, 2022.
388 5, 6
- 389 C. Formanek, C. R. Tilbury, J. P. Shock, K. ab Tessera, and A. Pretorius. Reduce, reuse, recy-
390 cle: Selective reincarnation in multi-agent reinforcement learning. *Workshop on Reincarnating*
391 *Reinforcement Learning at ICLR*, 2023. 3, 17
- 392 J. Fu, A. Kumar, O. Nachum, G. Tucker, and S. Levine. D4rl: Datasets for deep data-driven
393 reinforcement learning. *ArXiv Preprint*, 2020. 2, 7, 9
- 394 S. Fujimoto and S. S. Gu. A minimalist approach to offline reinforcement learning. *Advances in*
395 *Neural Information Processing Systems*, 2021. 8, 23, 24
- 396 S. Fujimoto, H. Hoof, and D. Meger. Addressing function approximation error in actor-critic methods.
397 *International Conference on Machine Learning*, 2018. 7
- 398 S. Fujimoto, D. Meger, and D. Precup. Off-policy deep reinforcement learning without exploration.
399 *International Conference on Machine Learning*, 2019. 3, 8, 23, 24
- 400 T. Gebu, J. Morgenstern, B. Vecchione, J. W. Vaughan, H. Wallach, H. D. I. au2, and K. Crawford.
401 Datasheets for datasets. *ArXiv Preprint*, 2021. 16

- 402 K. Ghasemipour, S. S. Gu, and O. Nachum. Why so pessimistic? estimating uncertainties for offline
403 rl through ensembles, and why their independence matters. *Advances in Neural Information*
404 *Processing Systems*, 2022. 3
- 405 R. Gorsane, O. Mahjoub, R. J. de Kock, R. Dubb, S. Singh, and A. Pretorius. Towards a standard-
406 ised performance evaluation protocol for cooperative MARL. *Advances in Neural Information*
407 *Processing Systems*, 2022. 6, 8, 9
- 408 S. Gu, L. Yang, Y. Du, G. Chen, F. Walter, J. Wang, Y. Yang, and A. Knoll. A review of safe
409 reinforcement learning: Methods, theory and applications. *ArXiv Preprint*, 2022. 9
- 410 C. Gulcehre, Z. Wang, A. Novikov, T. Paine, S. Gómez, K. Zolna, R. Agarwal, J. S. Merel, D. J.
411 Mankowitz, C. Paduraru, et al. RL unplugged: A suite of benchmarks for offline reinforcement
412 learning. *Advances in Neural Information Processing Systems*, 2020. 2, 7
- 413 J. K. Gupta, M. Egorov, and M. Kochenderfer. Cooperative multi-agent control using deep reinforce-
414 ment learning. *International Conference on Autonomous Agents and Multiagent Systems*, 2017.
415 6
- 416 J. Hu, S. Jiang, S. A. Harding, H. Wu, and S.-w. Liao. Rethinking the implementation tricks and
417 monotonicity constraint in cooperative multi-agent reinforcement learning, 2021. 23
- 418 J. Jiang and Z. Lu. Offline decentralized multi-agent reinforcement learning. *ArXiv Preprint*, 2021.
419 3, 4, 5, 6
- 420 V. Khattar and M. Jin. Winning the citylearn challenge: Adaptive optimization with evolutionary
421 search under trajectory-based guidance. *ArXiv Preprint*, 2022. 2
- 422 I. Kostrikov, A. Nair, and S. Levine. Offline reinforcement learning with implicit q-learning. *Deep*
423 *RL Workshop at NeurIPS*, 2021. 3
- 424 L. Kraemer and B. Banerjee. Multi-agent reinforcement learning as a rehearsal for decentralized
425 planning. *Elsevier Neurocomputing*, 2016. 4
- 426 A. Kumar, J. Fu, G. Tucker, and S. Levine. Stabilizing off-policy q-learning via bootstrapping error
427 reduction. *Neural Information Processing Systems*, 2019. 23
- 428 A. Kumar, A. Zhou, G. Tucker, and S. Levine. Conservative q-learning for offline reinforcement
429 learning. *Advances in Neural Information Processing Systems*, 2020. 3, 8, 23, 24
- 430 V. Kurenkov and S. Kolesnikov. Showing your offline reinforcement learning work: Online evaluation
431 budget matters. *International Conference on Machine Learning*, 2022. 8, 24
- 432 S. Levine, A. Kumar, G. Tucker, and J. Fu. Offline reinforcement learning: Tutorial, review, and
433 perspectives on open problems. *ArXiv Preprint*, 2020. 1, 4
- 434 R. Lowe, Y. I. Wu, A. Tamar, J. Harb, O. Pieter Abbeel, and I. Mordatch. Multi-agent actor-critic for
435 mixed cooperative-competitive environments. *Advances in neural information processing systems*,
436 30, 2017. 3, 7
- 437 C. Lu, P. J. Ball, T. G. Rudner, J. Parker-Holder, M. A. Osborne, and Y. W. Teh. Challenges and
438 opportunities in offline reinforcement learning from visual observations. *Decision Awareness in*
439 *Reinforcement Learning Workshop at ICML*, 2022. 2, 5
- 440 X. Lyu, Y. Xiao, B. Daley, and C. Amato. Contrasting centralized and decentralized critics in multi-
441 agent reinforcement learning. *International Conference on Autonomous Agents and Multi-Agent*
442 *Systems*, 2021. 4
- 443 X. Lyu, A. Baisero, Y. Xiao, and C. Amato. A deeper understanding of state-based critics in
444 multi-agent reinforcement learning. *ArXiv Preprint.*, 2022. 4

- 445 L. Meng, M. Wen, Y. Yang, C. Le, X. Li, W. Zhang, Y. Wen, H. Zhang, J. Wang, and B. Xu. Offline
446 pre-trained multi-agent decision transformer: One big sequence model conquers all starcraftii tasks.
447 *ArXiv Preprint*, 2021. [3](#), [4](#), [6](#), [7](#), [8](#), [9](#)
- 448 S. Mohanty, E. Nygren, F. Laurent, M. Schneider, C. Scheller, N. Bhattacharya, J. Watson, A. Egli,
449 C. Eichenberger, C. Baumberger, G. Vienken, I. Sturm, G. Sartoretti, and G. Spigler. Flatland-rl :
450 Multi-agent reinforcement learning on trains. *ArXiv Preprint*, 2020. [2](#), [6](#)
- 451 M. Nakamoto, Y. Zhai, A. Singh, Y. Ma, C. Finn, A. Kumar, and S. Levine. Cal-ql: Calibrated offline
452 rl pre-training for efficient online fine-tuning. *Workshop on Reincarnating Reinforcement Learning*
453 *at ICLR*, 2023. [3](#)
- 454 L. Pan, L. Huang, T. Ma, and H. Xu. Plan better amid conservatism: Offline multi-agent reinforcement
455 learning with actor rectification. *International Conference on Machine Learning*, 2022. [3](#), [4](#), [5](#), [6](#),
456 [7](#), [8](#), [24](#)
- 457 B. Peng, T. Rashid, C. Schroeder de Witt, P.-A. Kamienny, P. Torr, W. Böhmer, and S. Whiteson.
458 Facmac: Factored multi-agent centralised policy gradients. *Advances in Neural Information*
459 *Processing Systems*, 2021. [3](#), [6](#)
- 460 R. F. Prudencio, M. R. O. A. Maximo, and E. L. Colombini. A survey on offline reinforcement
461 learning: Taxonomy, review, and open problems. *IEEE Transactions on Neural Networks and*
462 *Learning Systems*, 2023. [2](#)
- 463 T. Rashid, M. Samvelyan, C. Schroeder, G. Farquhar, J. Foerster, and S. Whiteson. Qmix: Monotonic
464 value function factorisation for deep multi-agent reinforcement learning. *International Conference*
465 *on Machine Learning*, 2018. [7](#), [8](#), [23](#)
- 466 D. Rolnick, P. L. Donti, L. H. Kaack, K. Kochanski, A. Lacoste, K. Sankaran, A. S. Ross, N. Milojevic-
467 Dupont, N. Jaques, A. Waldman-Brown, A. S. Luccioni, T. Maharaj, E. D. Sherwin, S. K. Mukkav-
468 illi, K. P. Kording, C. P. Gomes, A. Y. Ng, D. Hassabis, J. C. Platt, F. Creutzig, J. Chayes, and
469 Y. Bengio. Tackling climate change with machine learning. *ACM Computing Surveys*, 2022. [7](#)
- 470 M. Samvelyan, T. Rashid, C. Schroeder de Witt, G. Farquhar, N. Nardelli, T. G. Rudner, C.-M.
471 Hung, P. H. Torr, J. Foerster, and S. Whiteson. The starcraft multi-agent challenge. *International*
472 *Conference on Autonomous Agents and MultiAgent Systems*, 2019. [3](#), [16](#)
- 473 R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, 2018. [23](#)
- 474 Q. Sykora, M. Ren, and R. Urtasun. Multi-agent routing value iteration network. *International*
475 *Conference on Machine Learning*, 2020. [2](#)
- 476 J. Terry, B. Black, N. Grammel, M. Jayakumar, A. Hari, R. Sullivan, L. S. Santos, C. Dieffendahl,
477 C. Horsch, R. Perez-Vicente, et al. Pettingzoo: Gym for multi-agent reinforcement learning.
478 *Advances in Neural Information Processing Systems*, 2021. [6](#)
- 479 E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. *IEEE/RSJ*
480 *International Conference on Intelligent Robots and Systems*, 2012. [6](#)
- 481 J. R. Vazquez-Canteli, S. Dey, G. Henze, and Z. Nagy. Citylearn: Standardizing research in multi-
482 agent reinforcement learning for demand response and urban energy management. *ArXiv Preprint*,
483 2020. [5](#), [7](#)
- 484 O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell,
485 T. Ewalds, P. Georgiev, J. Oh, D. Horgan, M. Kroiss, I. Danihelka, A. Huang, L. Sifre, T. Cai,
486 J. P. Agapiou, M. Jaderberg, A. S. Vezhnevets, R. Leblond, T. Pohlen, V. Dalibard, D. Budden,
487 Y. Sulsky, J. Molloy, T. L. Paine, C. Gulcehre, Z. Wang, T. Pfaff, Y. Wu, R. Ring, D. Yogatama,
488 D. Wunsch, K. McKinney, O. Smith, T. Schaul, T. Lillicrap, K. Kavukcuoglu, D. Hassabis, C. Apps,
489 and D. Silver. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*,
490 2019. [1](#)

- 491 J. Wang, W. Xu, Y. Gu, W. Song, and T. C. Green. Multi-agent reinforcement learning for active
492 voltage control on power distribution networks. *Advances in Neural Information Processing*
493 *Systems*, 2021. 5, 7
- 494 J. Whittlestone, K. Arulkumaran, and M. Crosby. The societal implications of deep reinforcement
495 learning. *Journal of Artificial Intelligence Research*, 2021. 9
- 496 H. Xu, X. Zhan, and X. Zhu. Constraints penalized q-learning for safe offline reinforcement learning.
497 *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022. 9
- 498 Y. Yang, X. Ma, L. Chenghao, Z. Zheng, Q. Zhang, G. Huang, J. Yang, and Q. Zhao. Believe what
499 you see: Implicit constraint approach for offline multi-agent reinforcement learning. *Advances in*
500 *Neural Information Processing Systems*, 2021. 3, 5, 6, 7, 8, 24
- 501 L. Yu, S. Qin, M. Zhang, C. Shen, T. Jiang, and X. Guan. A review of deep reinforcement learning
502 for smart building energy management. *IEEE Internet of Things Journal*, 2021. 7
- 503 Y. Yu. Towards sample efficient reinforcement learning. *International Joint Conference on Artificial*
504 *Intelligence*, 2018. 1
- 505 H. Zhang, S. Feng, C. Liu, Y. Ding, Y. Zhu, Z. Zhou, W. Zhang, Y. Yu, H. Jin, and Z. Li. CityFlow:
506 A multi-agent reinforcement learning environment for large scale city traffic scenario. *ACM*
507 *International World Wide Web Conference*, 2019. 2
- 508 K. Zhang, Z. Yang, H. Liu, T. Zhang, and T. Başar. Finite-sample analysis for decentralized batch
509 multiagent reinforcement learning with networked agents. *IEEE Transactions on Automatic*
510 *Control*, 2021. 3
- 511 M. Zhou, Z. Liu, P. Sui, Y. Li, and Y. Y. Chung. Learning implicit credit assignment for cooperative
512 multi-agent reinforcement learning. *Arxiv Preprint*, 2020. 5
- 513 Z. Zhu, M. Liu, L. Mao, B. Kang, M. Xu, Y. Yu, S. Ermon, and W. Zhang. Madiff: Offline multi-agent
514 learning with diffusion models. *Arxiv Preprint*, 2023. 3, 17

515 **Checklist**

- 516 1. For all authors...
- 517 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
- 518 contributions and scope? [Yes]
- 519 (b) Did you describe the limitations of your work? [Yes] See [section 8](#).
- 520 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See
- 521 [section 8](#).
- 522 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
- 523 them? [Yes]
- 524 2. If you are including theoretical results...
- 525 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 526 (b) Did you include complete proofs of all theoretical results? [N/A]
- 527 3. If you ran experiments (e.g. for benchmarks)...
- 528 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
- 529 mental results (either in the supplemental material or as a URL)? [Yes] Our datasets
- 530 and code are open-sourced.
- 531 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 532 were chosen)? [Yes] All of the training details are in [section 7](#) and the hyperparameter
- 533 details are in [Appendix D](#).
- 534 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
- 535 ments multiple times)? [Yes] See [Figure 4](#).
- 536 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 537 of GPUs, internal cluster, or cloud provider)? [Yes] See [Appendix D](#).
- 538 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 539 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 540 (b) Did you mention the license of the assets? [Yes] See our datasheet in [Appendix A](#) and
- 541 licence in [Appendix E](#).
- 542 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 543 See our datasheet in [Appendix A](#).
- 544 (d) Did you discuss whether and how consent was obtained from people whose data you’re
- 545 using/curating? [Yes] See our datasheet in [Appendix A](#).
- 546 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 547 information or offensive content? [Yes] See our datasheet in [Appendix A](#).
- 548 5. If you used crowdsourcing or conducted research with human subjects...
- 549 (a) Did you include the full text of instructions given to participants and screenshots, if
- 550 applicable? [N/A]
- 551 (b) Did you describe any potential participant risks, with links to Institutional Review
- 552 Board (IRB) approvals, if applicable? [N/A]
- 553 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 554 spent on participant compensation? [N/A]