

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

Claude Formanek<sup>1,2</sup>, Asad Jawa<sup>1</sup>, Jonathan Shock<sup>2</sup>, Arnu Pretorius<sup>1</sup>

<sup>1</sup> InstaDeep, <sup>2</sup> University of Cape Town



UNIVERSITY OF CAPE TOWN  
IYUNIVESITHI YASEKAPA • UNIVERSITEIT VAN KAAPSTAD

**TLDR:** Offline MARL is a nascent research field that has, to date, been hampered by the lack of standardised datasets to measure progress. To address this, we provide a diverse collection of high-quality multi-agent datasets with baselines.



Being able to harness the power of large datasets for developing autonomous systems could unlock enormous value for real-world applications. Many important industrial systems are multi-agent in nature and, in industry, large amounts of logged data from distributed system processes are stored. Offline MARL provides a promising paradigm for building effective decentralised online controllers from such datasets.

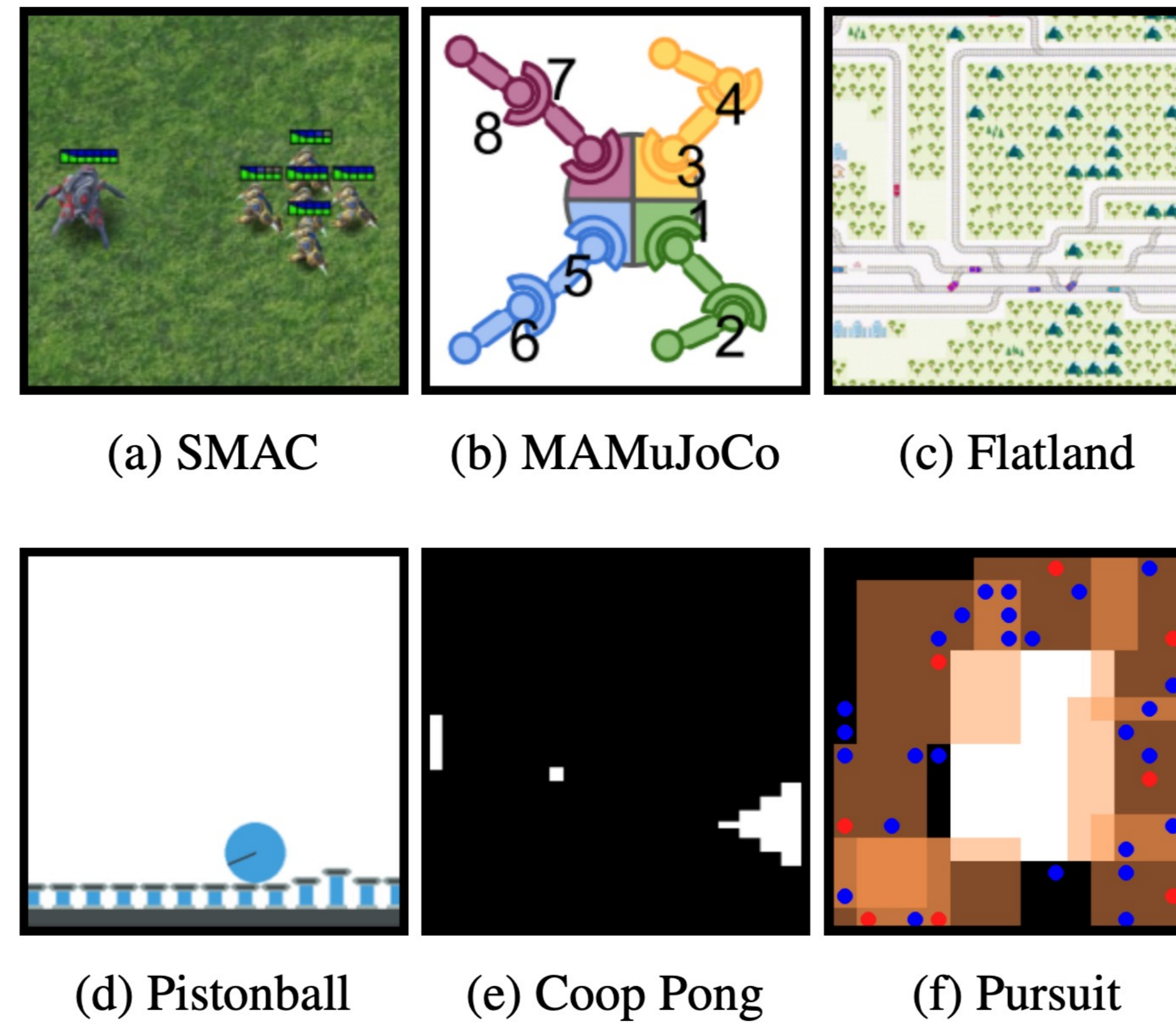
However, offline MARL is still in its infancy, and, therefore, lacks standardised benchmarks, baselines and evaluation protocols typically found in more mature subfields of RL. This deficiency makes it difficult for the community to sensibly measure progress. In this work, we aim to fill this gap by releasing off-the-grid MARL (OG-MARL): a framework for generating offline MARL datasets and algorithms. We release an initial set of datasets and baselines for cooperative offline MARL. Our datasets provide settings that are characteristic of real-world systems, including complex dynamics, non-stationarity, partial observability, suboptimality and sparse rewards, and are generated from popular online MARL benchmarks.

## Framework for Offline MARL

```
1 from og_marl.environments.smac import SMAC
2 from og_marl.offline_tools import MAOfflineEnvLogger
3 from og_marl.systems.executor_base import RandomExecutor
4
5 # Instantiate environment
6 env = SMAC("3m")
7
8 # Executor that chooses random actions for each agent
9 executor = RandomExecutor()
10
11 # Wrap env in offline logger
12 wrapped_env = MAOfflineEnvLogger(
13     environment=env
14 )
15
16 # Reset environment
17 timestep = wrapped_env.reset()
18
19 # Get some random actions from executor
20 all_agent_actions = executor.select_actions()
21
22 # Step the environment
23 next_timestep = wrapped_env.step(
24     actions=all_agent_actions
25 )
```

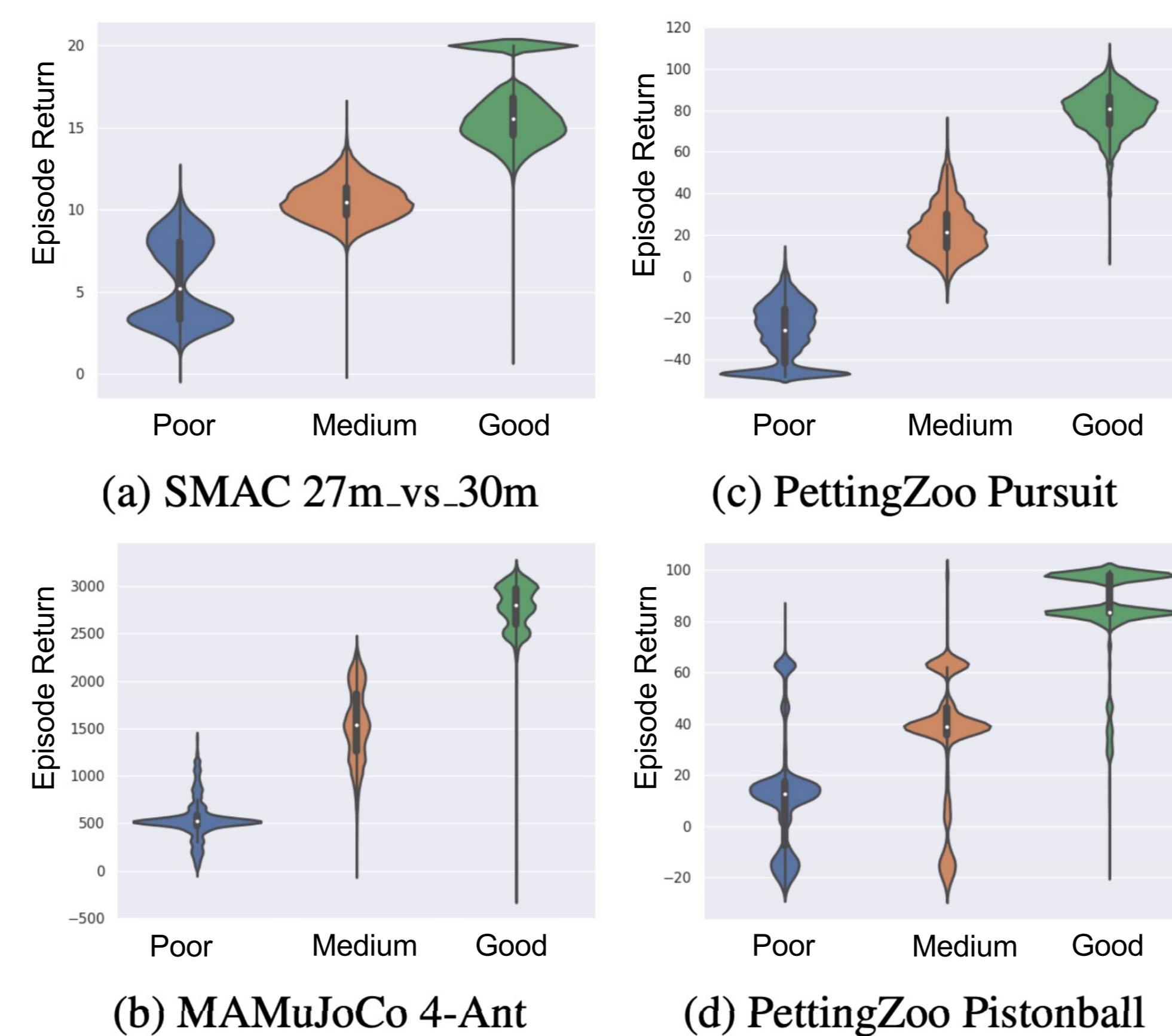
- Utilities for generating and analysing offline MARL datasets.
- Standardised re-implementation of popular offline MARL algorithms.
- Framework for developing novel offline MARL algorithms.
- Website to store and distribute datasets.

## Diverse Multi-Agent Environments



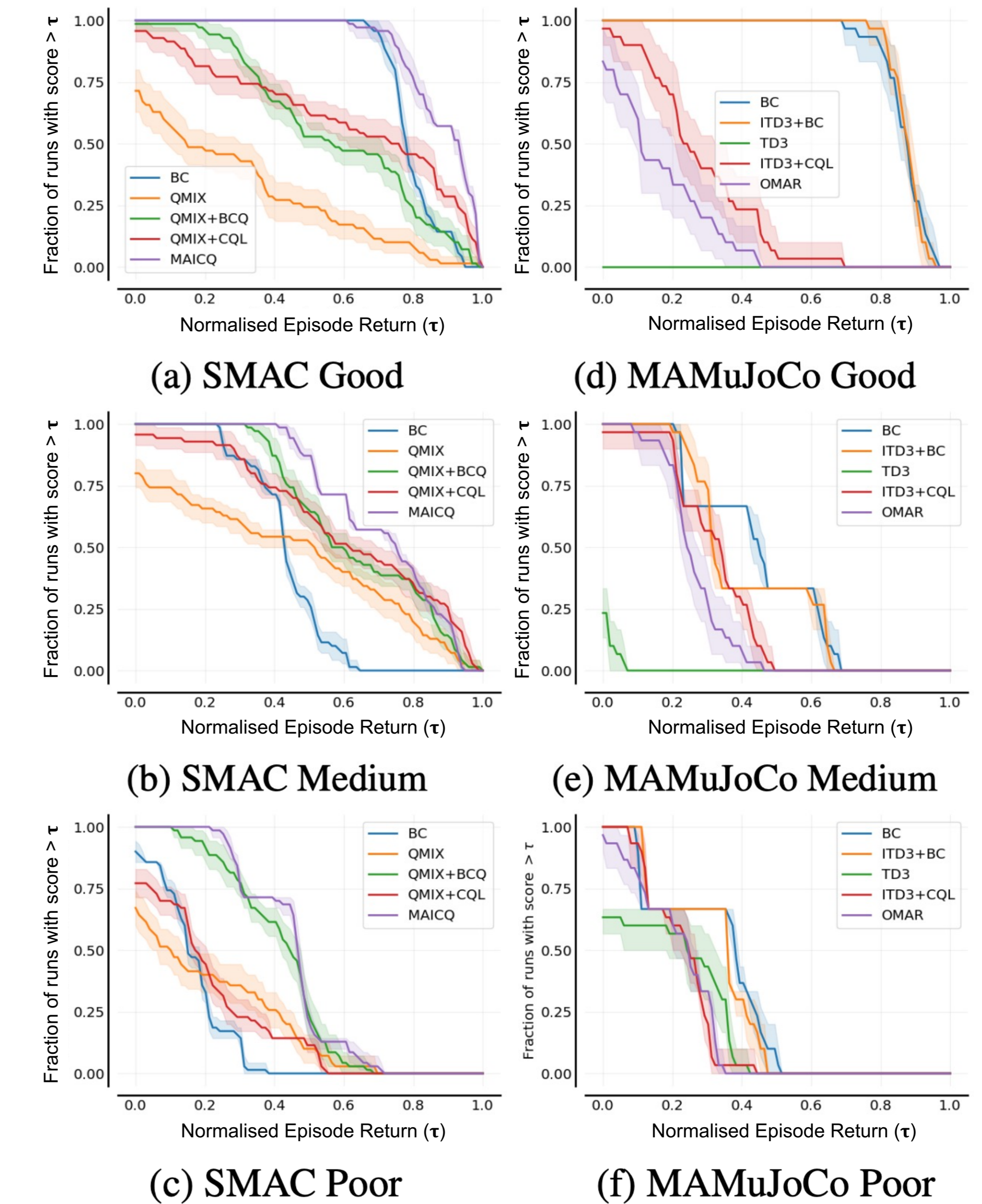
We provide datasets for several popular MARL benchmark environments including the SMAC (Samvelyan et al., 2019), Multi-Agent MuJoCo (Peng et al., 2021), Flatland (Mohanty et al., 2020) and environments from PettingZoo (Terry et al., 2021). Together these environments cover a broad range of task characteristics including: i) discrete and continuous action spaces, ii) vector and pixel-based observations, iii) dense and sparse rewards, iv) a varying number of agents (from 2 to 27 agents), and finally v) heterogeneous and homogeneous agents.

## Dataset Quality Classes



To generate the transitions in the datasets, we recorded environment interactions with online MARL algorithms. For each scenario, we provide three types of datasets. The dataset types are characterised by the quality of the joint policy that generated the trajectories in the dataset. We use violin plots to visualise the distribution of episode returns in the datasets.

## Comprehensive Baseline Results



We provide baselines for all of our datasets. Above we show the performance profiles [Agarwal et al., 2021] for popular offline MARL algorithms [Pan et al., (2022), Yang et al., (2021)]. The performance profiles aggregate across all SMAC and MA MuJoCo scenarios respectively. Roughly speaking, if a curve is above another it means that algorithms performs better than the other, with respect to the episode return at the end of offline training.

## References

- Agarwal, Rishabh, et al., "Deep reinforcement learning at the edge of the statistical precipice", *Advances in Neural Information Processing Systems* (2021).
- Samvelyan, Rashid, et al., "The StarCraft multi-agent challenge", *Autonomous Agents and Multi-Agent Systems* (2019)
- Nygren, et al., "Flatland-RL : Multiagent reinforcement learning on trains", *ArXiv* (2020)
- Terry, Black, et al., "PettingZoo: Gym for multi-agent reinforcement learning", *Advances in Neural Information Processing Systems* (2021)
- Pan, Huang, et al., "Plan better amid conservatism: Offline multi-agent reinforcement learning with actor rectification", *International Conference on Machine Learning* (2022)
- Yang, Ma, et al. "Believe what you see: Implicit constraint approach for offline multi-agent reinforcement learning", *Advances in Neural Information Processing Systems* (2021)



DEEP LEARNING  
INDABA

