



ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

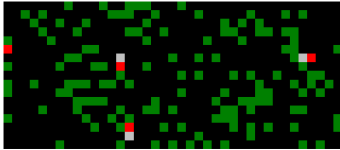
SOLVING THE HARVEST CPR APPROPRIATION PROBLEM WITH POLICY GRADIENT TECHNIQUES

AAS FINAL PROJECT - ACADEMIC YEAR 2020/2021

Alessio Falai
`alessio.falai@studio.unibo.it`
September 10, 2021

Alma Mater Studiorum - University of Bologna

ENVIRONMENT



Full environment



Local observation

- Small 25×7 grid for the single-agent setting and big 39×17 map for multi-agent scenarios
- 9 actions in total: movement + tagging + gifting
- Local observation: RGB image of size $3 \times 20 \times 21$ (20 squares ahead and 10 squares on each side of the agent)

- Social Learning: reshape the reward function of other agents with the goal of promoting cooperation
- Gifting: peer-rewarding strategy in which agents can reward others with a new specialized action
- Gifting mechanisms: each time an agent sends a gift g , its gifting budget is decremented by g
 - Zero-Sum: the budget is infinite, but the agent incurs a penalty $-g$ for every gifting action taken
 - Fixed Budget: the budget is fixed at the start of the episode and when it's empty no more gifting can happen
 - Replenishable Budget: the budget expands as a function of collected environmental rewards

1. Single-agent DQN vs VPG with RLLib
2. Custom VPG vs TRPO vs PPO on Cartpole
3. Custom VPG vs TRPO vs PPO on single-agent Harvest
4. Custom PPO on multi-agent Harvest, with and without Zero-Sum gifting
5. Custom PPO on multi-agent Harvest, with Replenishable and Fixed Budget gifting

THANK YOU FOR YOUR ATTENTION

REFERENCES



A. Lupu and Doina Precup.

Gifting in multi-agent reinforcement learning.

In *AAMAS*, 2020.



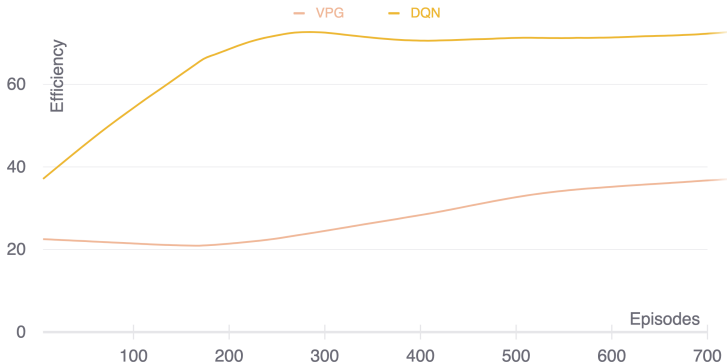
Julien Pérolat, Joel Z. Leibo, Vinícius Flores Zambaldi, Charles Beattie, Karl Tuyls, and Thore Graepel.

A multi-agent reinforcement learning model of common-pool resource appropriation.

CoRR, abs/1707.06600, 2017.

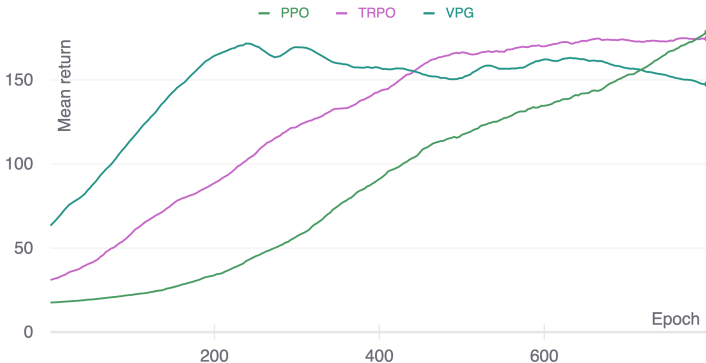
BACKUP FRAMES

SINGLE-AGENT DQN VS VPG WITH RLLIB



Value-based methods seem more suited for the Harvest environment (higher returns)

CUSTOM VPG VS TRPO VS PPO ON CARTPOLE



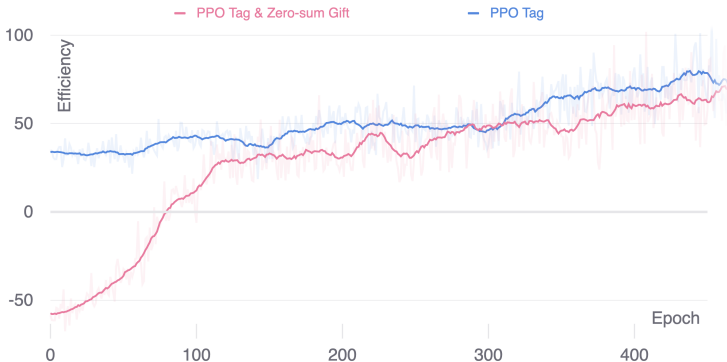
Custom implementations of policy gradient methods are valid, as all agents converge to good returns in the selected test environment

CUSTOM VPG VS TRPO VS PPO ON SINGLE-AGENT HARVEST



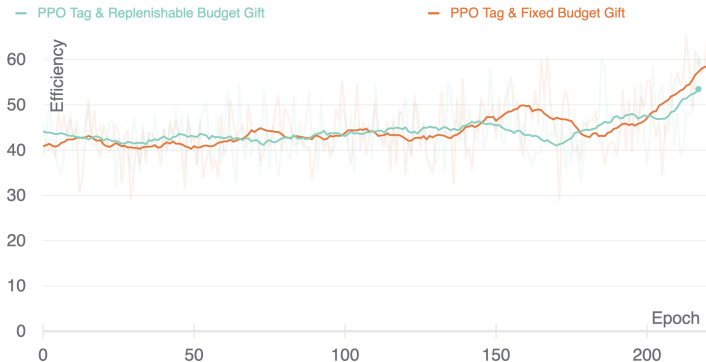
VPG and PPO converge to similar results, while TRPO diverges on the single-agent setting of Harvest

CUSTOM PPO ON MULTI-AGENT HARVEST, WITH AND WITHOUT ZERO-SUM GIFTING



Agents tend to be very generous at the beginning, while later training stages show that enabling or disabling Zero-Sum leads to similar results

CUSTOM PPO ON MULTI-AGENT HARVEST, WITH REPLENISHABLE AND FIXED BUDGET GIFTING



Results show that the Replenishable and Fixed Budget gifting strategies tend to follow similar training curves and converge to comparable results