

## Time Your Rewards:

Learning Temporally Consistent Rewards from a Single Video Demonstration Huaxiaoyue (Yuki) Wang\*, Will Huey\*, Anne Wu, Yoav Artzi, Sanjiban Choudhury

Cornell University

Learn the Task Specified in the Video Demo Through RL

Define a class of reward function based on the video demo



A good reward function needs to:

PORTAL

**DEnforce** temporal constraints

**Avoid reward** hacking



SDTW+ reward helps agent • achieve higher success rate • train stably



DTW+: DTW with Cumulative Reward Bonus

0.2

0.0+

0.00

[1] S. Sontakke, J. Zhang, S. Arnold, K. Pertsch, E. Bıyık, D. Sadigh, C. Finn, and L. Itti. Robo-clip: One demonstration is enough to learn robot policies. Advances in Neural Information Processing Systems, 36, 2024.



DTW+

SDTW+

RoboCLIP

## **O** 0 S 0.2 0.0 -0.25 2.00 0.00 2.00 0.50 1.75 0.25 1.50 1.75 1e6 1e6 Environment Steps Environment Steps

Please check out our paper for more results.

## Future Work

- Longer horizon, periodic tasks
  - More domains (e.g. cooking)
- Learning from human video demos



- Vision encoder that bridges visual embodiment gap
- More theoretical analysis on reward hacking given different sequence matching rewards