# TAAC: Temporally Abstract Actor-Critic for Continuous Control

*Haonan Yu*, *Wei Xu*, and *Haichao Zhang*
Horizon Robotics

# Motivation: temporal abstraction (TA)

Temporally correlated actions

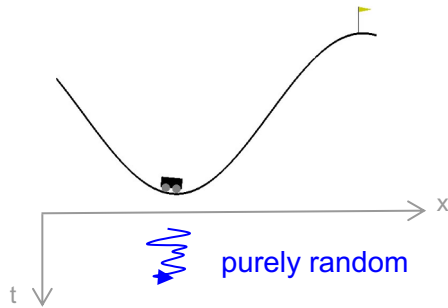# Motivation: temporal abstraction (TA)

Temporally correlated actions

1. Temporally persistent exploration

# Motivation: temporal abstraction (TA)
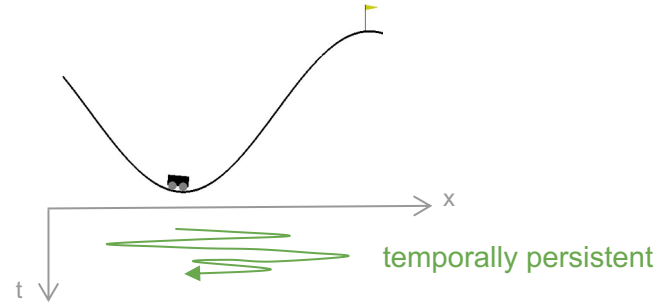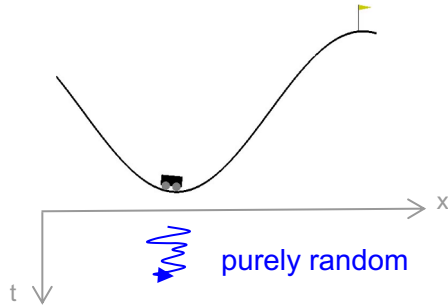
Temporally correlated actions

1. Temporally persistent exploration



purely random

# Motivation: temporal abstraction (TA)

Temporally correlated actions

1. Temporally persistent exploration
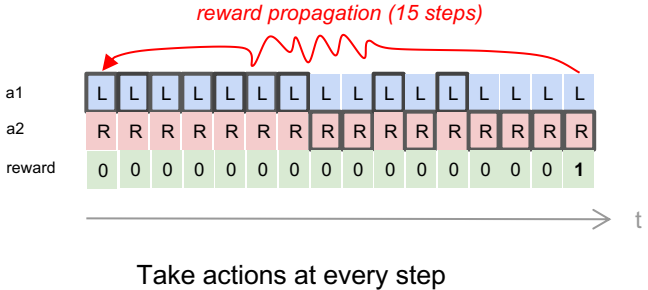


Moving out of the valley more easily!

# Motivation: temporal abstraction (TA)

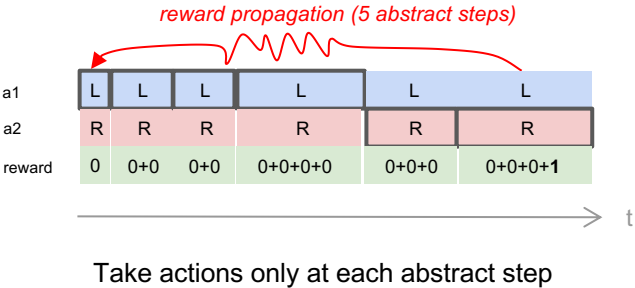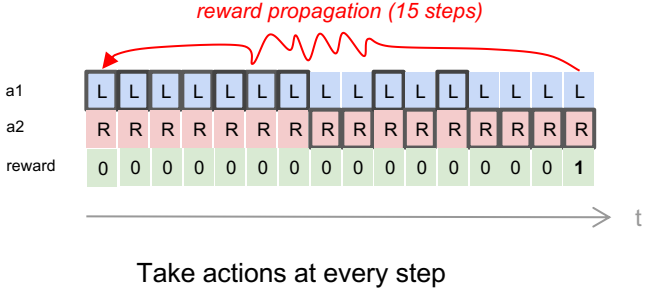2. Better credit assignment with delayed reward (shorter horizon)

# Motivation: temporal abstraction (TA)

2. Better credit assignment with delayed reward (shorter horizon)



Take actions at every step

# Motivation: temporal abstraction (TA)

## 2. Better credit assignment with delayed reward (shorter horizon)



Take actions at every step

Take actions only at each abstract step

# Open-loop *vs.* closed-loop repetition

Action repetition: perhaps the simplest temporal abstraction technique.

# Open-loop *vs.* closed-loop repetition

Action repetition: perhaps the simplest temporal abstraction technique.
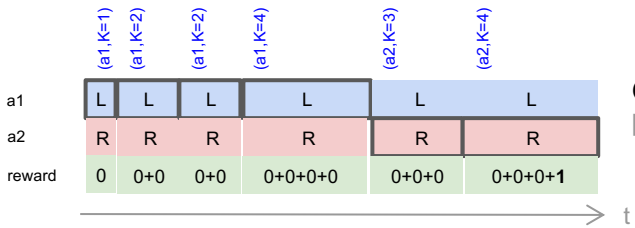
***Key questions:*** *what action to repeat & how long to repeat it?*

# Open-loop *vs.* closed-loop repetition

Action repetition: perhaps the simplest temporal abstraction technique.

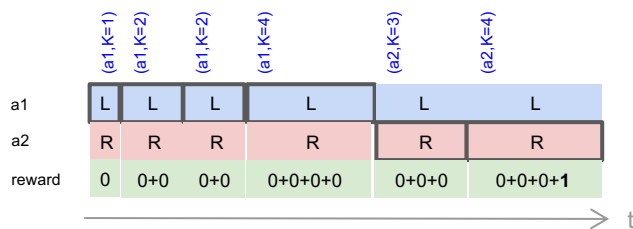***Key questions:*** *what action to repeat & how long to repeat it?*



Open-loop methods output an action and its duration at once
[Sharma et al., 2017] [Biedenkapp et al., 2021] [Dabney et al., 2021]
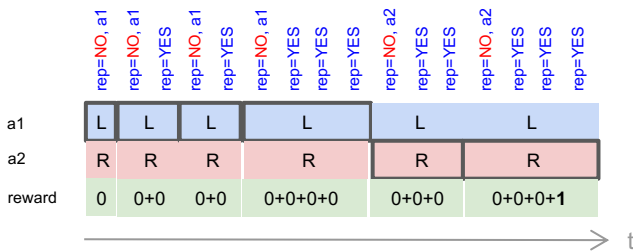
# Open-loop *vs.* closed-loop repetition

Action repetition: perhaps the simplest temporal abstraction technique.

***Key questions:*** *what action to repeat & how long to repeat it?*



Open-loop methods output an action and its duration at once
[Sharma et al., 2017] [Biedenkapp et al., 2021] [Dabney et al., 2021]

Closed-loop methods decide "act-or-repeat" at every step
[Neunert et al., 2020] [Chen et al., 2021]

12

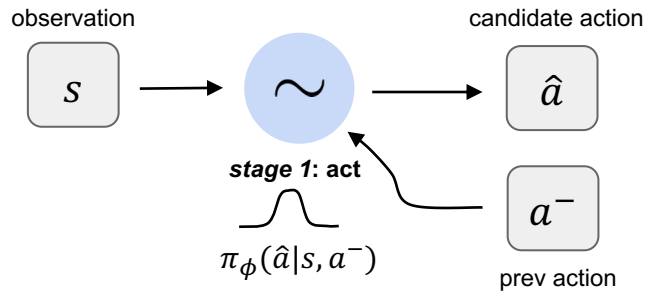# Temporally abstract actor-critic (TAAC)

TAAC incorporates closed-loop action repetition into off-policy actor-critic for continuous control

# Temporally abstract actor-critic (TAAC)

TAAC incorporates closed-loop action repetition into off-policy actor-critic for continuous control

*TAAC's two-stage policy*



observation

candidate action

$s$

$\hat{a}$

*stage 1*: act

$a^-$

$\pi_\phi(\hat{a}|s, a^-)$

prev action

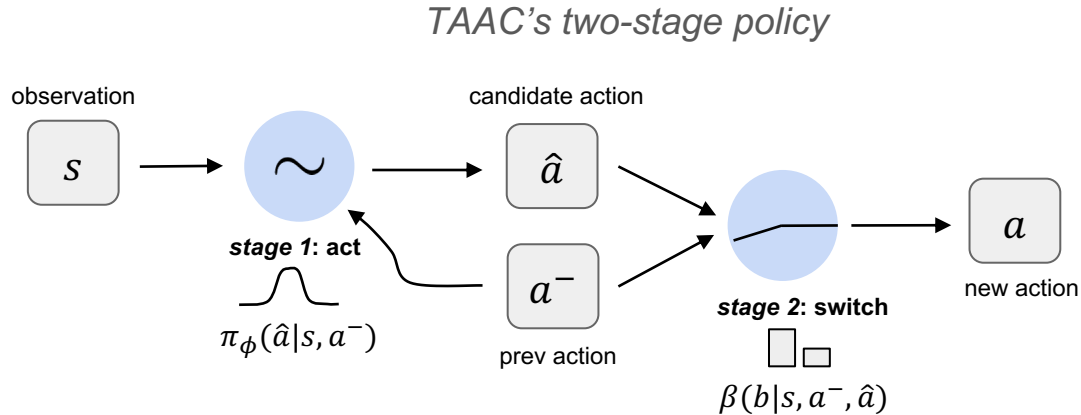# Temporally abstract actor-critic (TAAC)

TAAC incorporates closed-loop action repetition into off-policy actor-critic for continuous control

*TAAC's two-stage policy*

# Temporally abstract actor-critic (TAAC)

TAAC incorporates closed-loop action repetition into off-policy actor-critic for continuous control

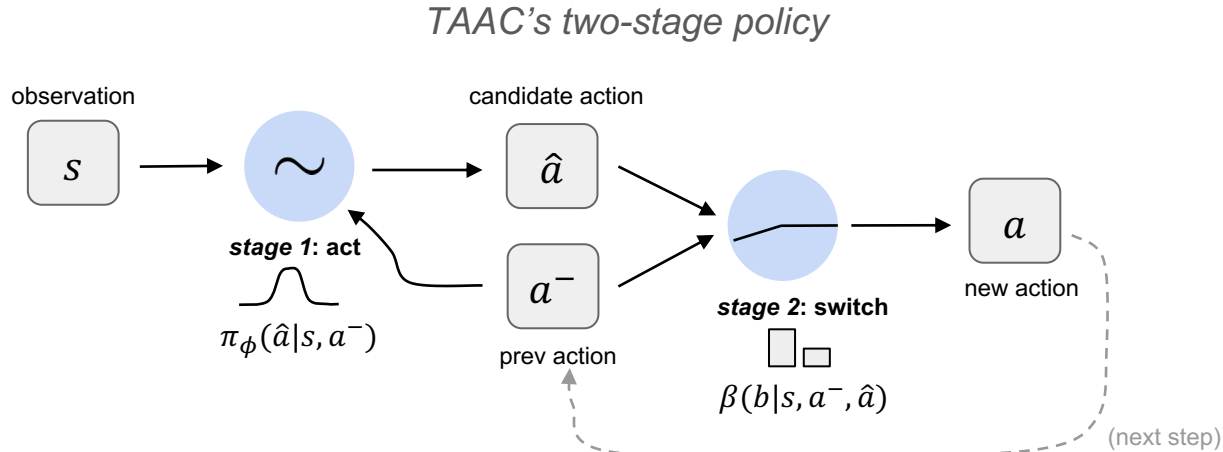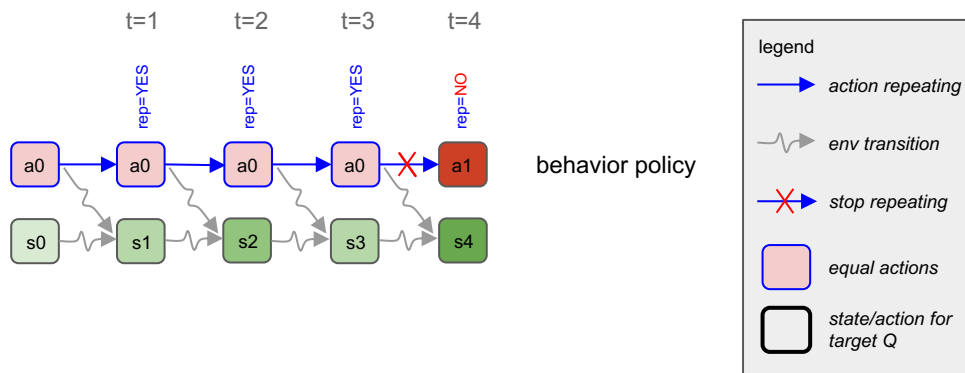*TAAC's two-stage policy*

# Policy evaluation

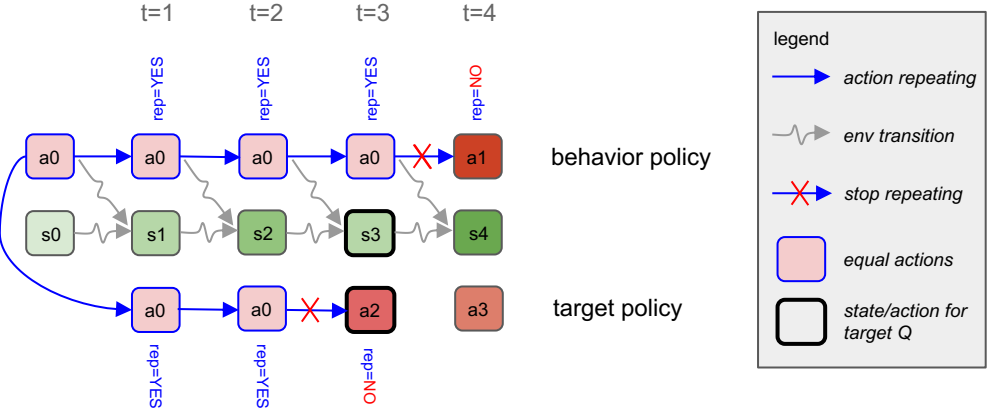A novel compare-through operator for multi-step TD backup

# Policy evaluation

A novel compare-through operator for multi-step TD backup
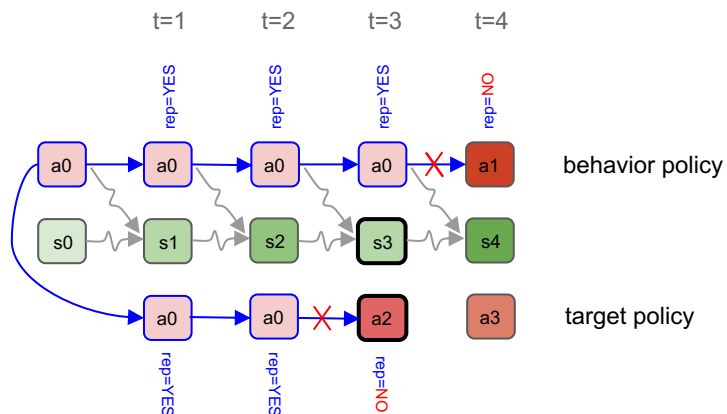
# Policy evaluation

## A novel compare-through operator for multi-step TD backup



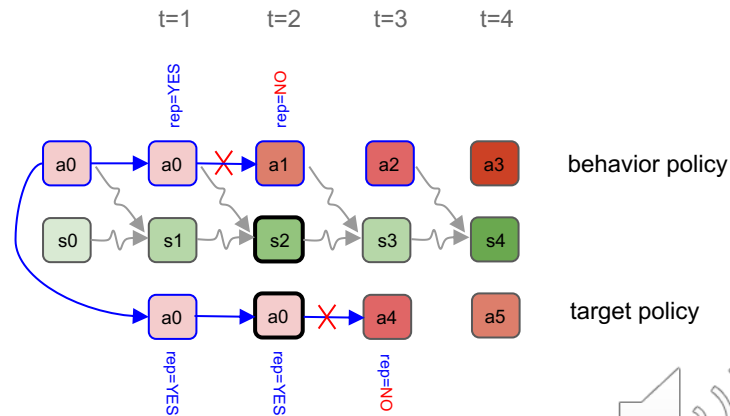Example 1: We use Q(s3,a2) as the target to bootstrap Q(s0,a0) (3-step TD)

# Policy evaluation

A novel compare-through operator for multi-step TD backup



Example 1: We use Q(s3,a2) as the target to bootstrap Q(s0,a0) (3-step TD)

Example 2: We use Q(s2,a0) as the target to bootstrap Q(s0,a0) (2-step TD)

# Policy improvement

A closed-form solution for the switching policy to speed up policy learning: sampling two actions by the exponential Q values

# Policy improvement

A closed-form solution for the switching policy to speed up policy learning: sampling two actions by the exponential Q values

$$\beta(0)^* = \exp\left(\frac{Q(s, a^-)}{\alpha}\right) / Z(s) \qquad \beta(1)^* = \exp\left(\frac{Q(s, \hat{a})}{\alpha}\right) / Z(s)$$

# Policy improvement

A closed-form solution for the switching policy to speed up policy learning: sampling two actions by the exponential Q values

$$\beta(0)^* = \exp\left(\frac{Q(s, a^-)}{\alpha}\right) / Z(s) \qquad \beta(1)^* = \exp\left(\frac{Q(s, \hat{a})}{\alpha}\right) / Z(s)$$
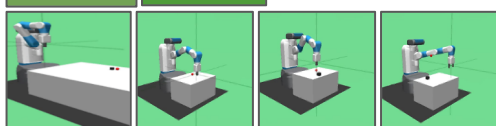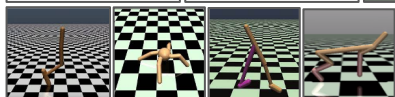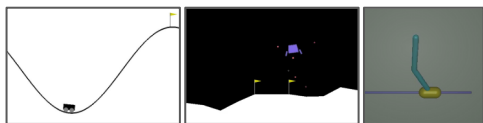
The actor $\pi_\phi(\hat{a}|s, a^-)$ is trained similarly as in DDPG [Lillicrap et al., 2016] and SAC

[Haarnoja et al., 2018]: $\dfrac{\partial Q(s, \hat{a})}{\partial \phi}\beta^*(1)$ (this is a good approximation to the full gradient)

# Tasks

5 categories of 14 continuous control tasks (13 standard; 1 customized)



| Category | Task | Gym environment name | Observation space | Action space |
|---|---|---|---|---|
| **SimpleControl** | *MountainCarContinuous* <br> *LunarLanderContinuous* <br> *InvertedDoublePendulum* | `MountainCarContinuous-v0` <br> `LunarLanderContinuous-v2` <br> `InvertedDoublePendulum-v2` | $\mathbb{R}^2$ <br> $\mathbb{R}^8$ <br> $\mathbb{R}^{11}$ | $[-1,1]^1$ <br> $[-1,1]^2$ <br> $[-1,1]^1$ |
| **Locomotion** | *Hopper* <br> *Ant* <br> *Walker2d* <br> *HalfCheetah* | `Hopper-v2` <br> `Ant-v2` <br> `Walker2d-v2` <br> `HalfCheetah-v2` | $\mathbb{R}^{11}$ <br> $\mathbb{R}^{111}$ <br> $\mathbb{R}^{17}$ | $[-1,1]^3$ <br> $[-1,1]^8$ <br> $[-1,1]^6$ |
| **Terrain** | *BipedalWalker* <br> *BipedalWalkerHardcore* | `BipedalWalker-v2` <br> `BipedalWalkerHardcore-v2` | $\mathbb{R}^{24}$ | $[-1,1]^4$ |
| **Manipulation** | *FetchReach* <br> *FetchPush* <br> *FetchSlide* <br> *FetchPickAndPlace* | `FetchReach-v1` <br> `FetchPush-v1` <br> `FetchSlide-v1` <br> `FetchPickAndPlace-v1` | $\mathbb{R}^{13}$ <br> $\mathbb{R}^{28}$ | $[-1,1]^4$ |
| **Driving** | *Town01* | `Town01` | "camera": $\mathbb{R}^{128\times64\times3}$, <br> "radar": $\mathbb{R}^{200\times4}$, <br> "collision": $\mathbb{R}^{4\times3}$, <br> "IMU": $\mathbb{R}^7$, <br> "goal": $\mathbb{R}^3$, <br> "velocity": $\mathbb{R}^3$, <br> "navigation": $\mathbb{R}^{8\times3}$ <br> "prev action": $[-1,1]^4$ | $[-1,1]^4$ |

# Experiment results

Comparison methods

SAC [Haarnoja et al., 2018]: flat RL

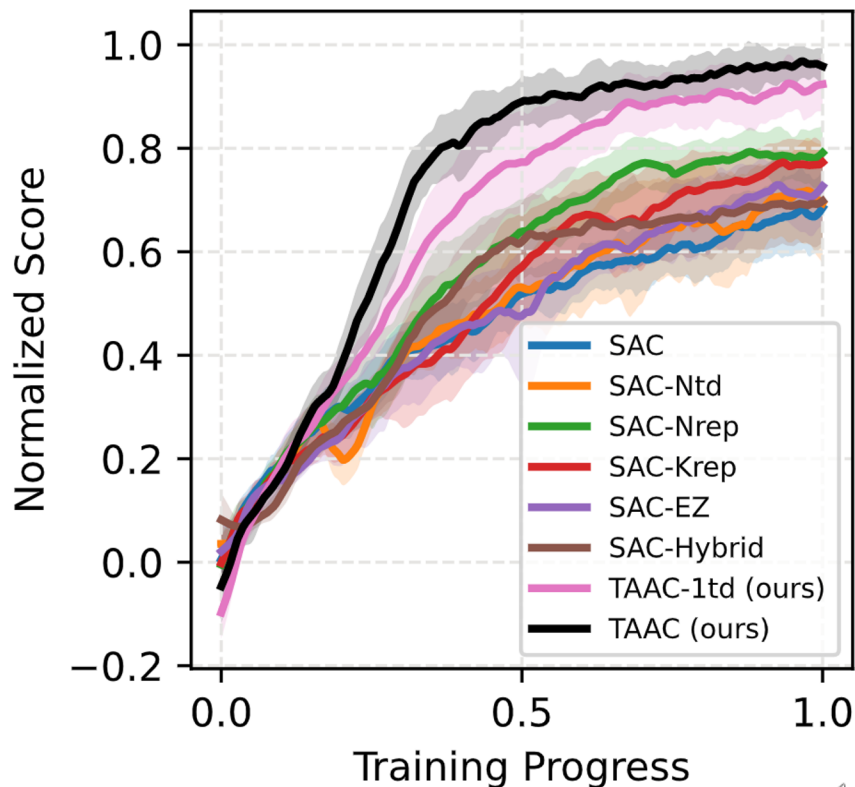SAC-Ntd: SAC + Retrace [Munos et al., 2016] + N-TD

SAC-Nrep: SAC + Fixed action repetition

SAC-Krep: open-loop action repetition [Sharma et al., 2017; Biedenkapp et al., 2021]

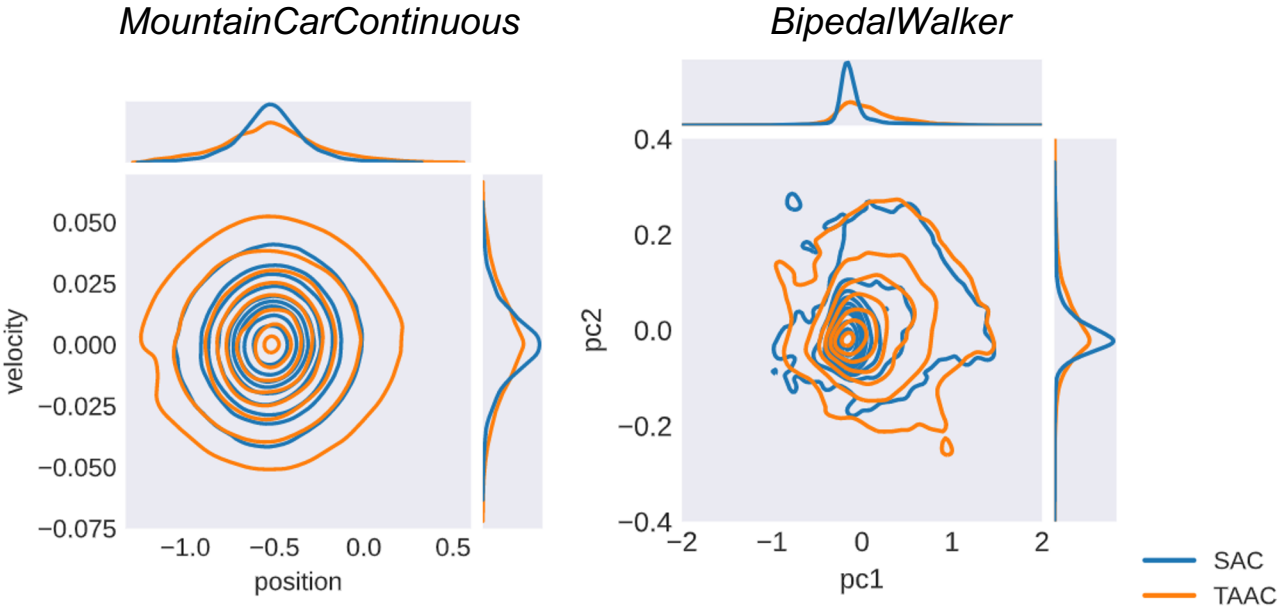SAC-EZ: SAC with EZ-greedy [Dabney et al., 2021]

SAC-Hybrid: closed-loop action repetition formulation from H-MPO [Neunert et al., 2020] with SAC backbone
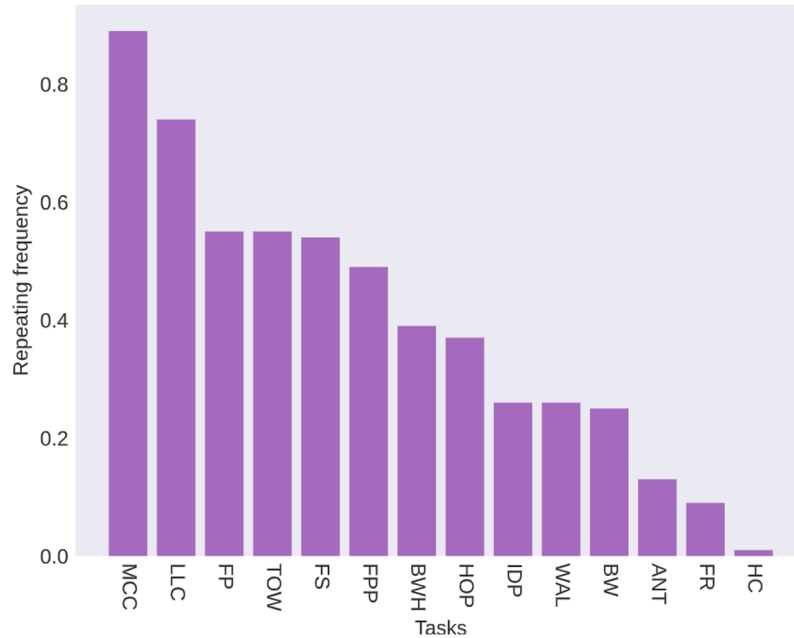
TAAC-1td: TAAC without the compare-through operator

# Experiment results

## Comparison methods

SAC [Haarnoja et al., 2018]: flat RL

SAC-Ntd: SAC + Retrace [Munos et al., 2016] + N-TD

SAC-Nrep: SAC + Fixed action repetition

SAC-Krep: open-loop action repetition [Sharma et al., 2017; Biedenkapp et al., 2021]

SAC-EZ: SAC with EZ-greedy [Dabney et al., 2021]

SAC-Hybrid: closed-loop action repetition formulation from H-MPO [Neunert et al., 2020] with SAC backbone

TAAC-1td: TAAC without the compare-through operator

*(Normalized and averaged over 14 tasks)*

# Exploration behavior analysis



*MountainCarContinuous*

*BipedalWalker*
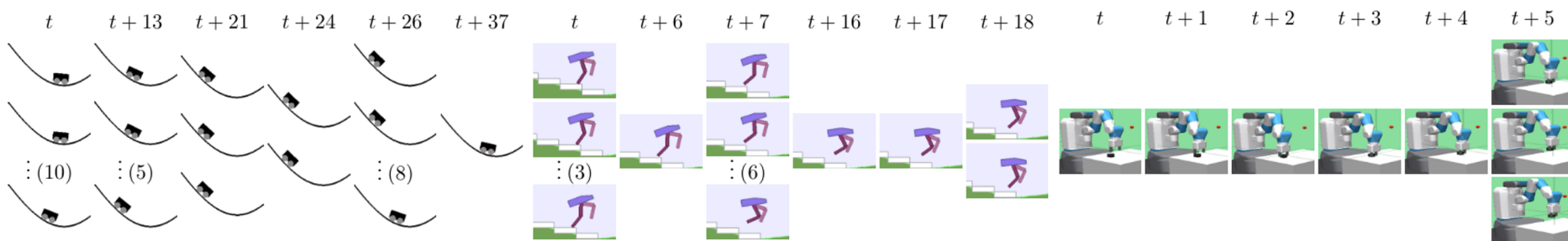
Random policies behaviors

# Action repeating frequency



Evaluating a trained TAAC model
for 100 episodes and calculating the repeating frequency

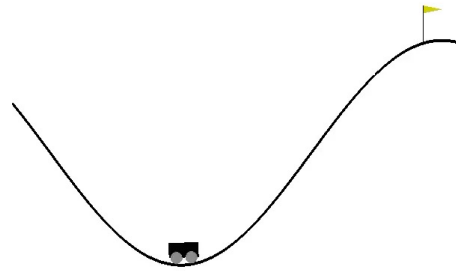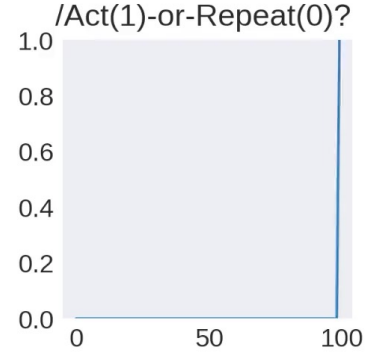# Action repeating patterns



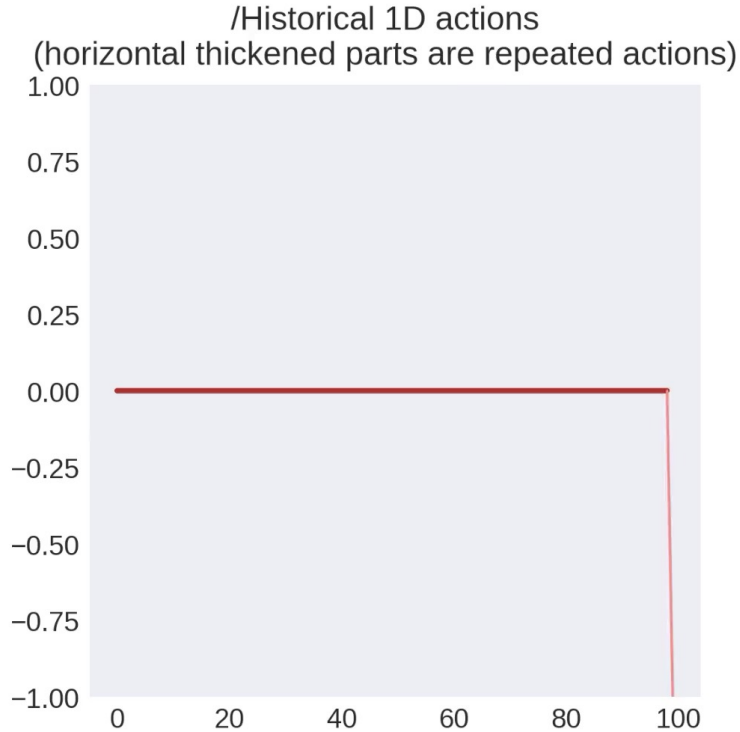*Each column represents the same action*

TAAC learns to skip learning to generate new actions at non-critical states, and save the actor network's representational power for critical states!
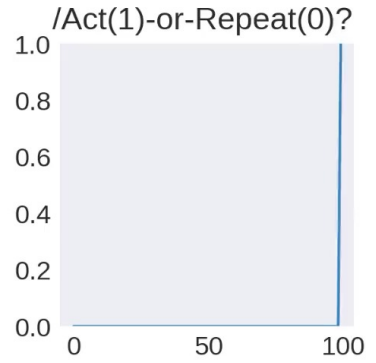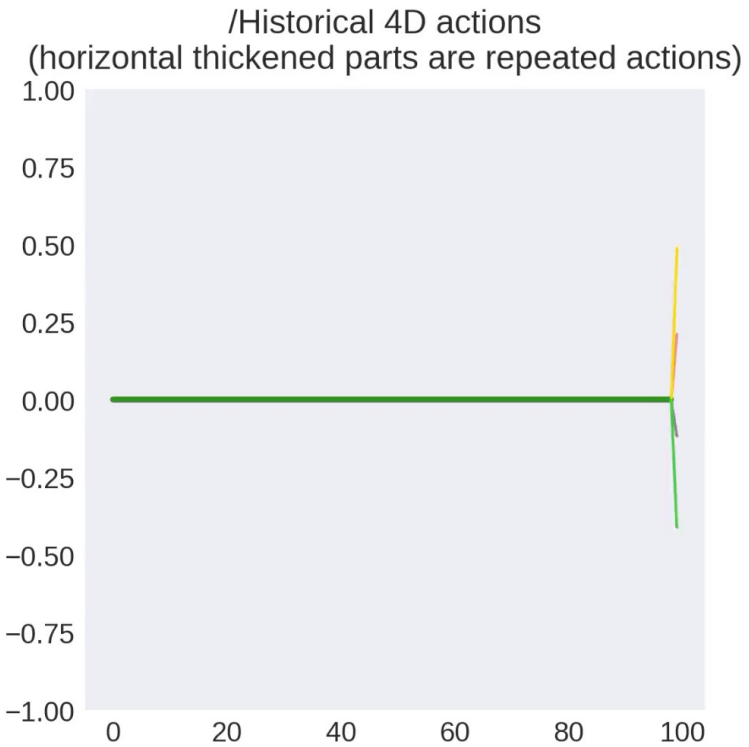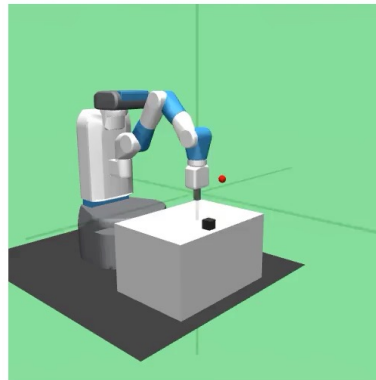
# Demos

# Demo - *MountainCarContinuous*



/Historical 1D actions
(horizontal thickened parts are repeated actions)

/Act(1)-or-Repeat(0)?

# Demo - *BipedalWalker*



/Historical 4D actions
(horizontal thickened parts are repeated actions)

/Act(1)-or-Repeat(0)?

(Played in 0.5x speed for a better view)

# Demo - *FetchPickAndPlace*



/Historical 4D actions
(horizontal thickened parts are repeated actions)

/Act(1)-or-Repeat(0)?

(Played in 0.5x speed for a better view)

# TAAC: Temporally Abstract Actor-Critic for Continuous Control

*Haonan Yu*, *Wei Xu*, and *Haichao Zhang*
Horizon Robotics

Code: *https://github.com/hnyu/taac*