# Temporal Difference Learning for Model Predictive Control

Nicklas Hansen,   Xiaolong Wang*,   Hao Su*
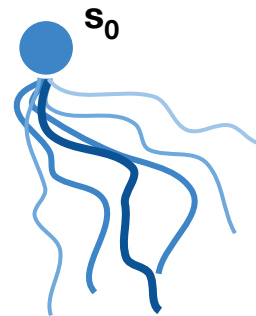
**ICML 2022**

UC San Diego

# Data-Driven Model Predictive Control (*MPC*)

- Plan using a ***learned*** model of the environment

- Objective $\mathbb{E}_{\Gamma \sim \Pi_\theta} \left[ \sum_{t=0}^{\infty} \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right]$ intractable



$\mathbf{s_0}$

(repeat for $\infty$ steps)
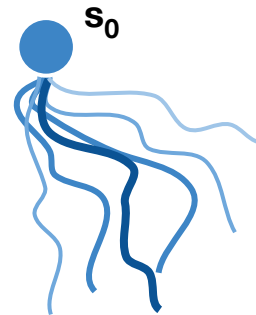
# Data-Driven Model Predictive Control (*MPC*)

- Plan using a ***learned*** model of the environment

- Objective $\mathbb{E}_{\Gamma} \sim \Pi_{\theta} \left[ \sum_{t=0}^{\infty} \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right]$ intractable

- Instead find ***locally optimal*** trajectory $\mathbb{E}_{\Gamma} \sim \Pi_{\theta} \left[ \sum_{t=0}^{H} \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right]$

- **Two major challenges:**

  - Compounding model errors

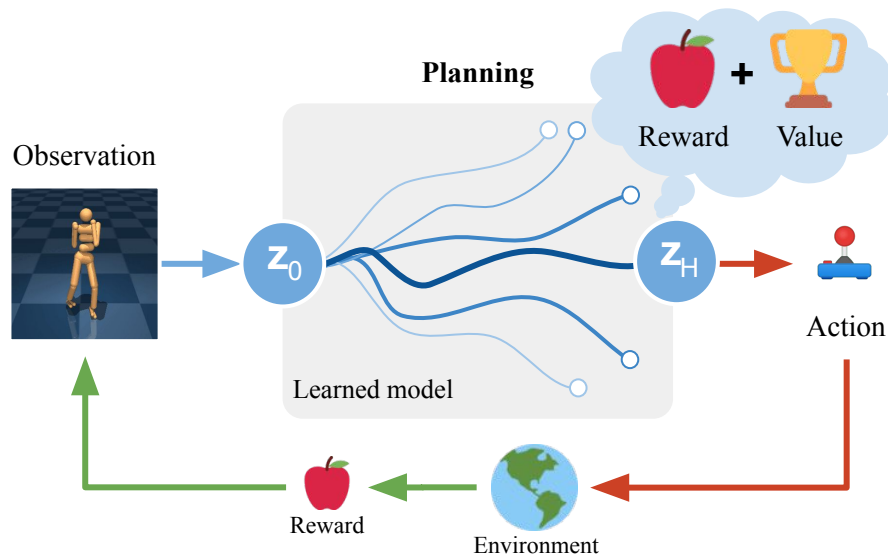  - Cost of long-horizon planning

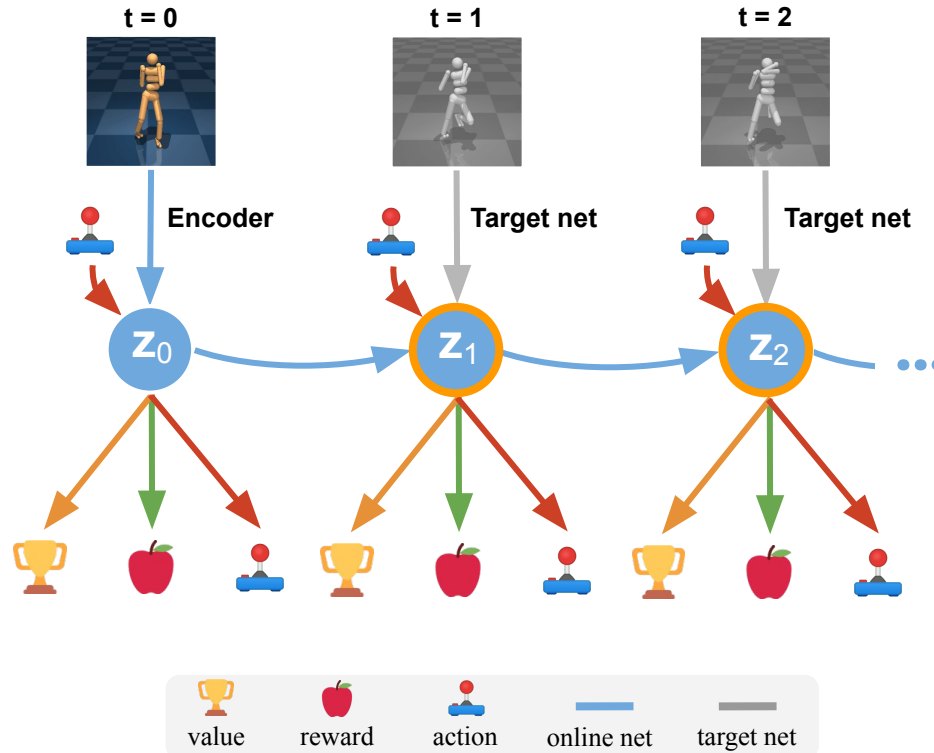$\mathbf{s_0}$

(repeat for $H$ steps)

# How can TD-learning help MPC?

**Inference** *(planning)*

- Planning in latent space

- Return estimate:

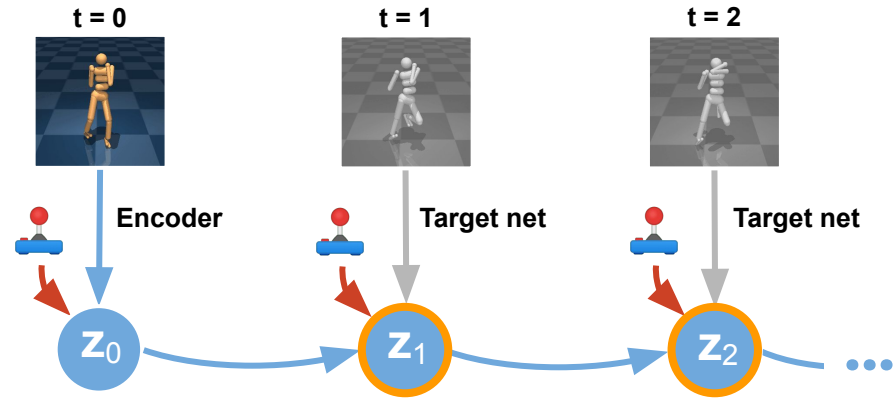$$\mathbb{E}_{\Gamma} \left[ \gamma^H Q_\theta(\mathbf{z}_H, \mathbf{a}_H) + \sum_{t=0}^{H-1} \gamma^t R_\theta(\mathbf{z}_t, \mathbf{a}_t) \right]$$

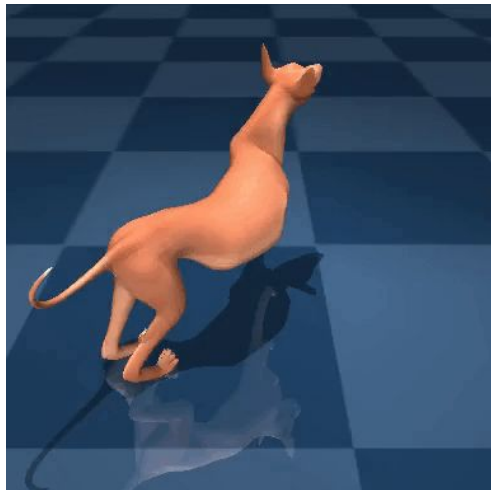**Value**     **Rewards**

# TD-MPC

# TD-MPC



Minimize diff. between **recurrent prediction** and **target encoding**

# Results

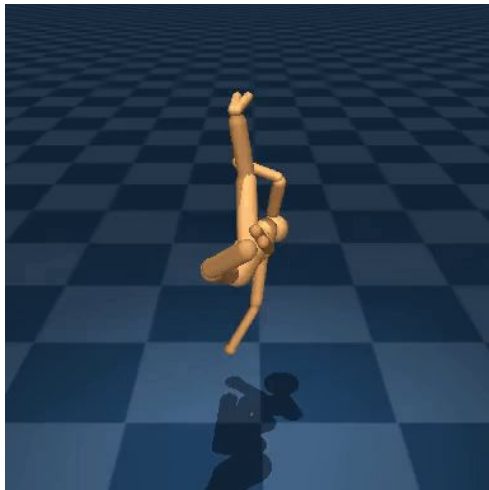**TD-MPC** solves ***challenging*** continuous control problems

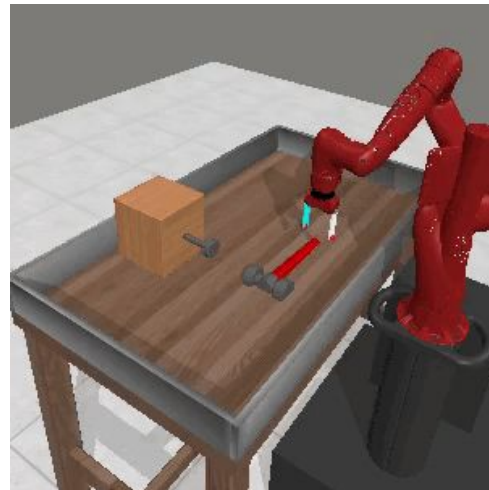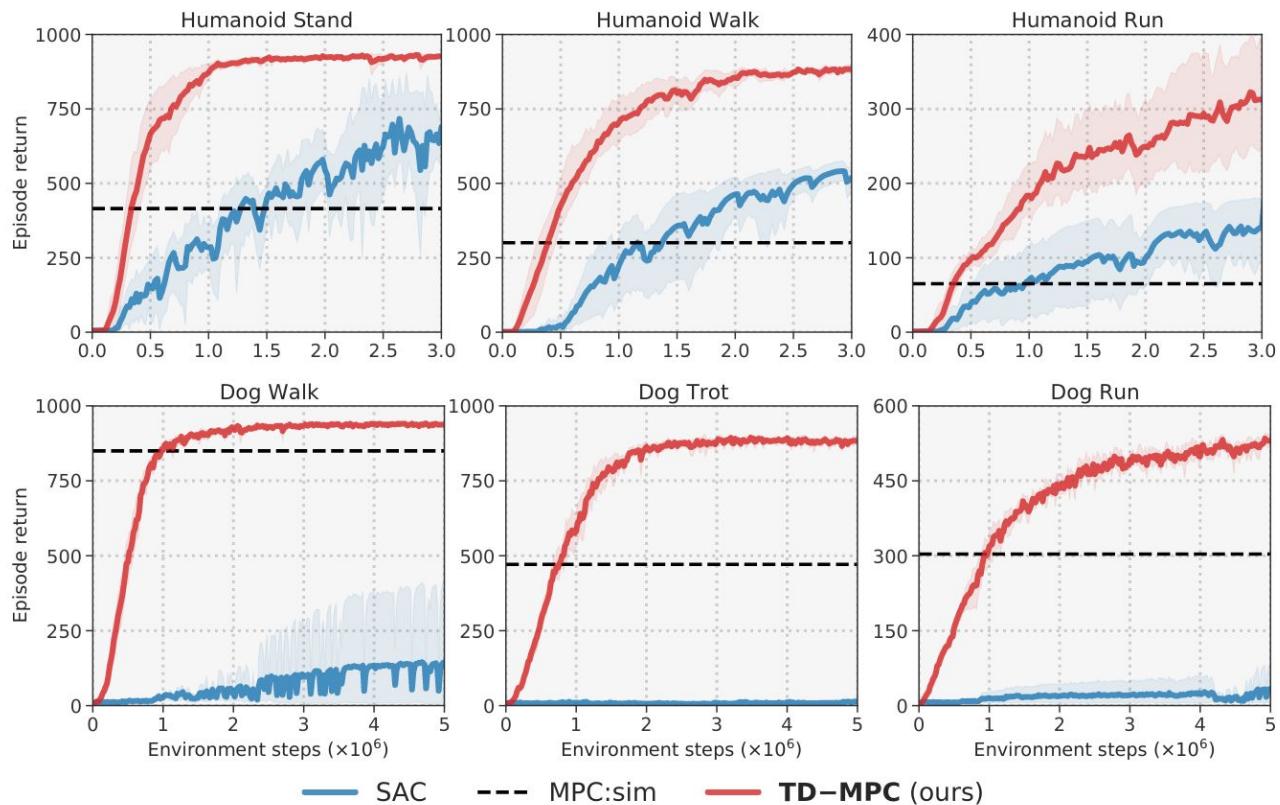**Dog Run**          **Humanoid Run**          **Hammer**

# Results

# Poster:  6-8pm today



# nicklashansen.github.io/td-mpc