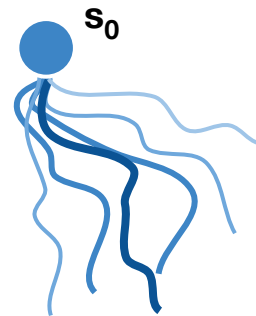# Temporal Difference Learning for Model Predictive Control

Nicklas Hansen,   Xiaolong Wang*,   Hao Su*

UC San Diego

# Data-Driven Model Predictive Control

- 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

$\mathbf{s_0}$

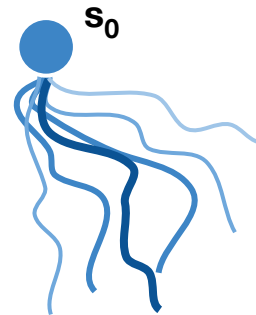(repeat for $\infty$ steps)

# Data-Driven Model Predictive Control

- 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)
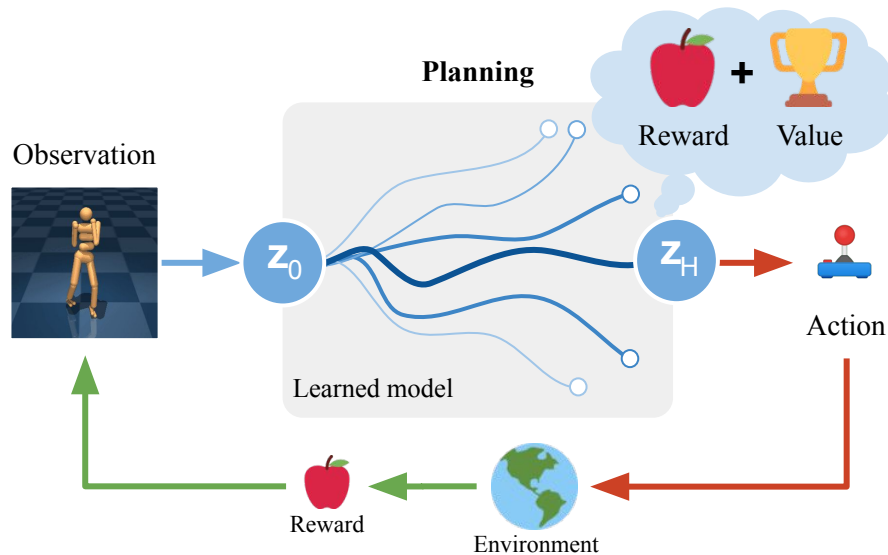
# How can TD-learning help MPC?

- Learning a ***terminal value function*** by TD-learning

    - MPC yields temporally ***local*** optimal solutions

    - A value function approximates the ***globally*** optimal solution

- Learning a ***task-oriented*** representation

    - Model-based RL typically ***models everything*** in the environment

    - Model-free RL only retains information ***predictive of reward***

# TD-MPC

**Inference** *(planning)*

- Planning in latent space

- Return estimate:

$$\mathbb{E}_{\Gamma} \left[ \underbrace{\gamma^{H} Q_{\theta}(\mathbf{z}_{H}, \mathbf{a}_{H})}_{\textbf{Value}} + \underbrace{\sum_{t=0}^{H-1} \gamma^{t} R_{\theta}(\mathbf{z}_{t}, \mathbf{a}_{t})}_{\textbf{Rewards}} \right]$$
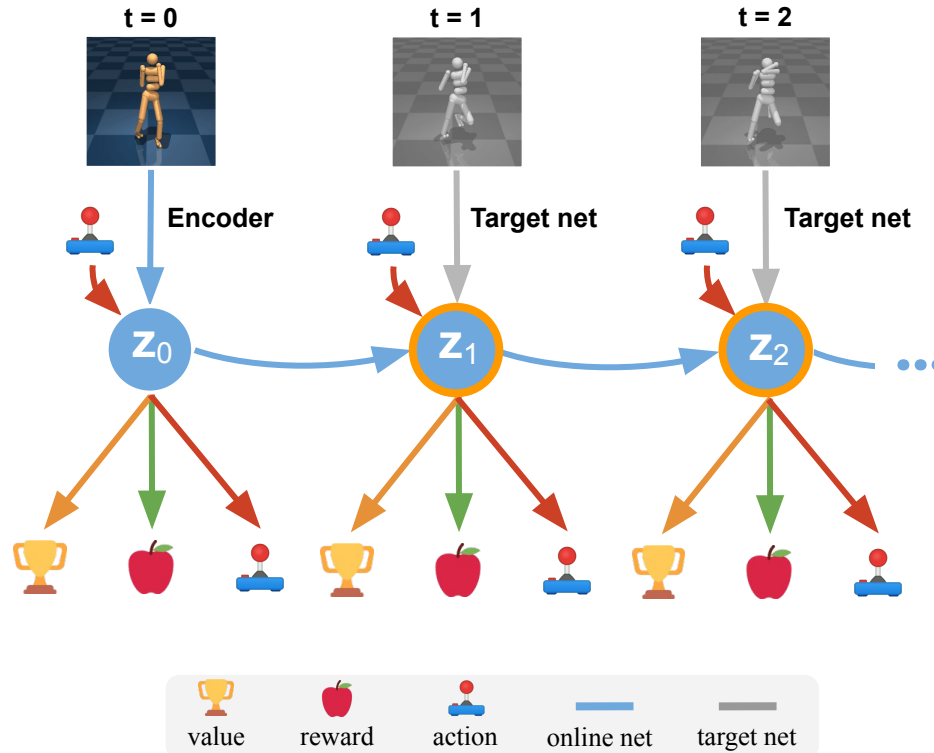
# TD-MPC

Task-Oriented Latent Dynamics (**TOLD**) model

- Model only parts of environment that are ***predictive of reward***

- Learned ***jointly*** with value function ***using TD-learning***

# TD-MPC

# TD-MPC

TOLD minimizes the objective

$$\mathcal{J}(\theta;\Gamma) = \sum_{i=t}^{t+H} \lambda^{i-t} \mathcal{L}(\theta;\Gamma_i), \qquad (7)$$
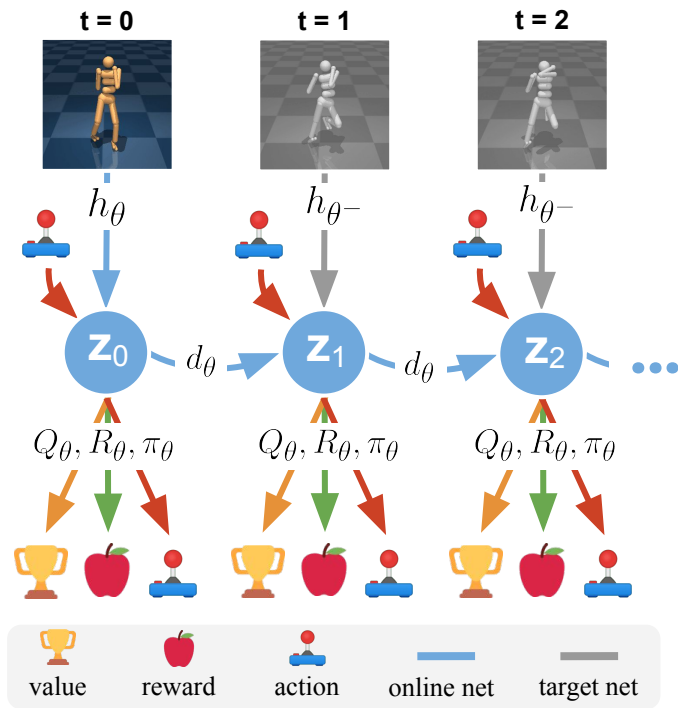
where

$$\mathcal{L}(\theta;\Gamma_i) = c_1 \underbrace{\|R_\theta(\mathbf{z}_i, \mathbf{a}_i) - r_i\|_2^2}_{\text{reward}} \qquad (8)$$

$$+ c_2 \underbrace{\|Q_\theta(\mathbf{z}_i, \mathbf{a}_i) - (r_i + \gamma Q_{\theta^-}(\mathbf{z}_{i+1}, \pi_\theta(\mathbf{z}_{i+1})))\|_2^2}_{\text{value}} \quad (9)$$

$$+ c_3 \underbrace{\|d_\theta(\mathbf{z}_i, \mathbf{a}_i) - h_{\theta^-}(\mathbf{s}_{i+1})\|_2^2}_{\text{latent state consistency}} \qquad (10)$$

and the policy minimizes

$$\mathcal{J}_\pi(\theta;\Gamma) = -\sum_{i=t}^{t+H} \lambda^{i-t} Q_\theta(\mathbf{z}_i, \pi_\theta(\text{sg}(\mathbf{z}_i))), \qquad (11)$$

# TD-MPC

Why learn a policy?

- **Planning:** policy ***action proposals*** speed up convergence

- **Learning:** estimating Q-targets via planning is ***very slow***; use policy instead

$$\mathcal{L}(\theta; \Gamma_i) = c_1 \underbrace{\|R_\theta(\mathbf{z}_i, \mathbf{a}_i) - r_i\|_2^2}_{\text{reward}} \qquad (8)$$
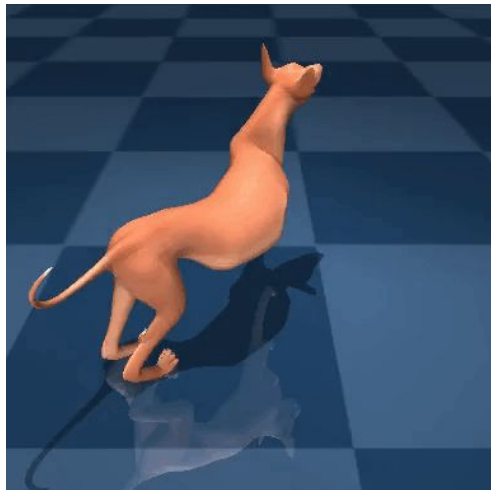
$$+ c_2 \underbrace{\|Q_\theta(\mathbf{z}_i, \mathbf{a}_i) - (r_i + \gamma Q_{\theta^-}(\mathbf{z}_{i+1}, \pi_\theta(\mathbf{z}_{i+1})))\|_2^2}_{\text{value}} \quad (9)$$

$$+ c_3 \underbrace{\|d_\theta(\mathbf{z}_i, \mathbf{a}_i) - h_{\theta^-}(\mathbf{s}_{i+1})\|_2^2}_{\text{latent state consistency}} \qquad (10)$$
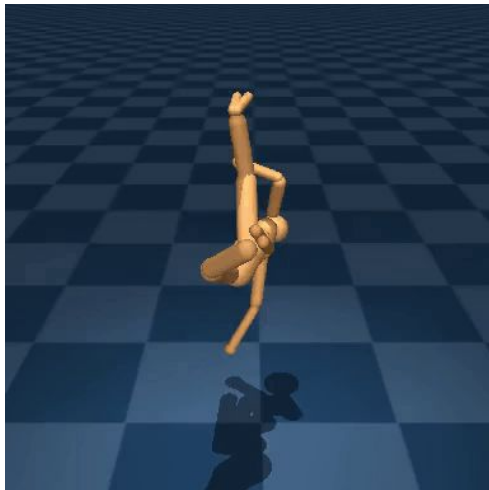
# Results

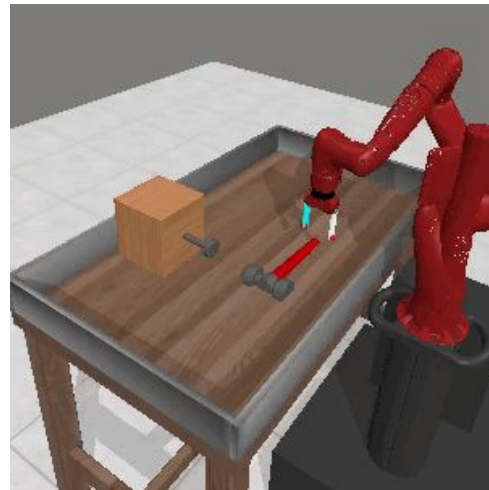**TD-MPC** solves ***challenging*** continuous control problems

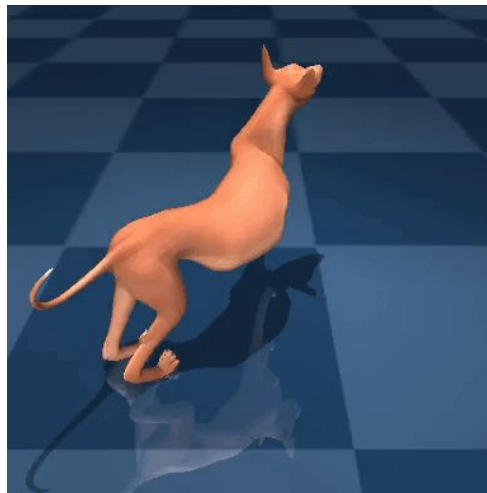| Dog Run | Humanoid Run | Hammer |
|---------|--------------|--------|

# Results

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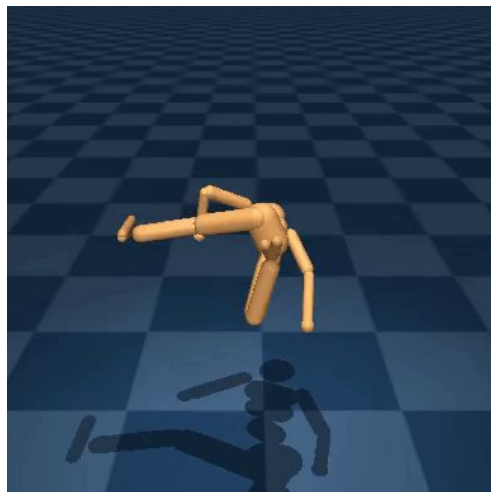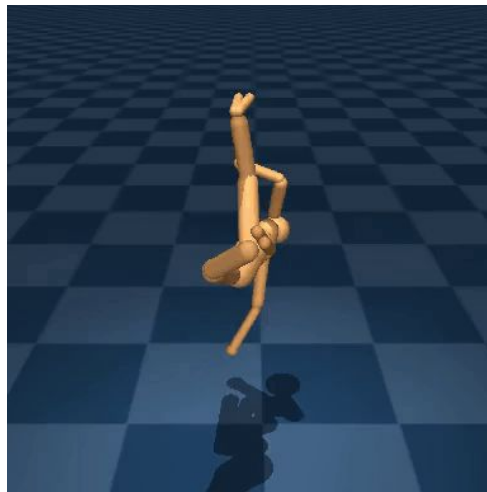
# Results

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

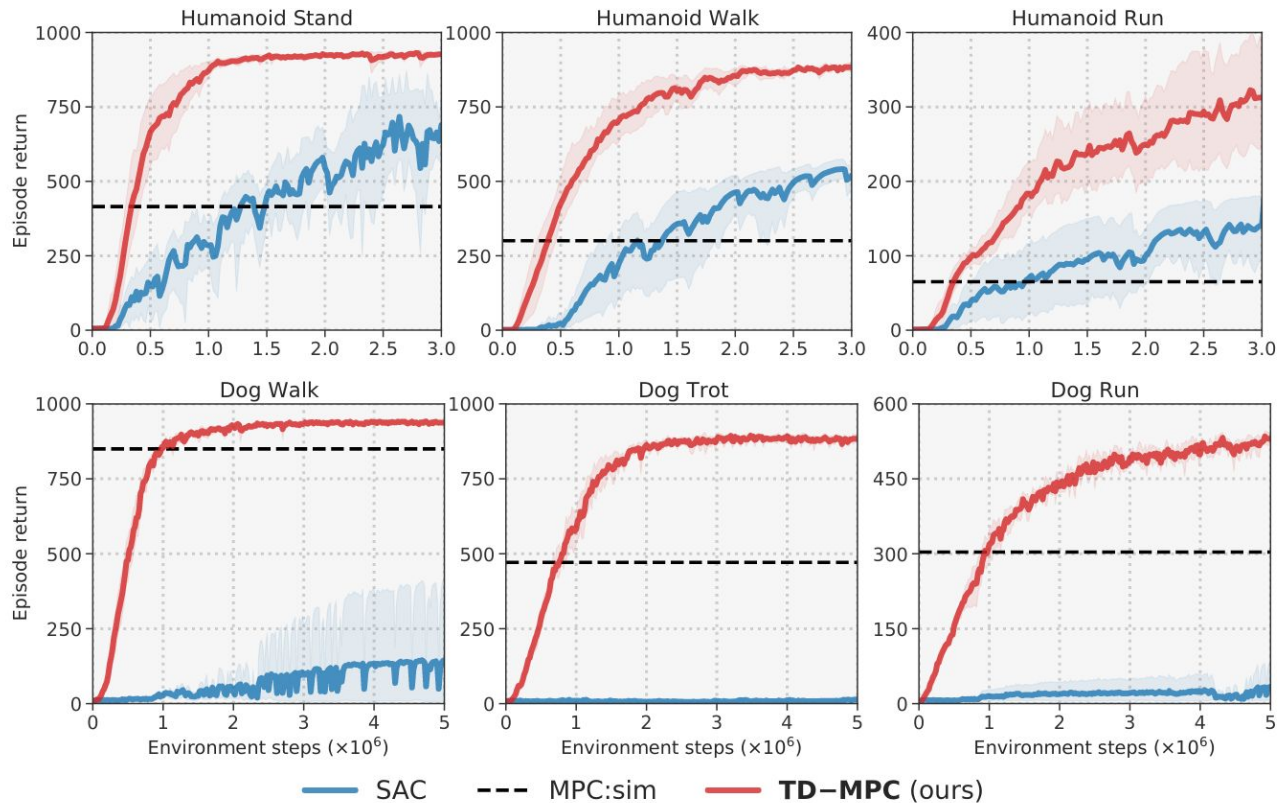**SAC**                                        **TD-MPC**

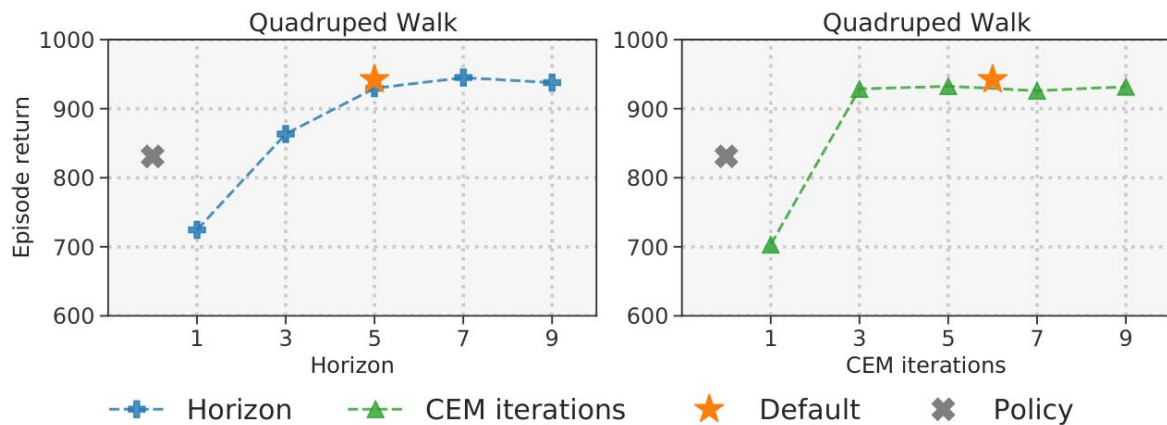# Results

# Results

More **planning** → better **performance**
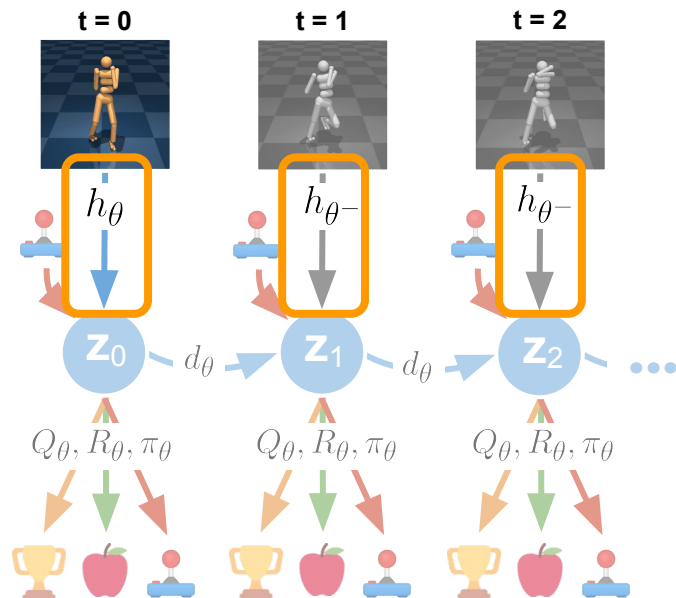


Variable budget *at test-time*

# Results

Replace MLP encoder with CNN  →  competitive performance on *image-based RL*

| *100k env. steps* | Model-free | | | | Model-based | | | | Ours |
|---|---|---|---|---|---|---|---|---|---|
| | SAC State | SAC Pixels | CURL | DrQ | PlaNet | Dreamer | MuZero* | Eff.Zero* | **TD-MPC** |
| Cartpole Swingup | $812\pm45$ | $419\pm40$ | $597\pm170$ | **$759\pm92$** | $563\pm73$ | $326\pm27$ | $219\pm122$ | **$813\pm19$** | **$770\pm70$** |
| Reacher Easy | $919\pm123$ | $145\pm30$ | $517\pm113$ | $601\pm213$ | $82\pm174$ | $314\pm155$ | $493\pm145$ | **$952\pm34$** | $628\pm105$ |
| Cup Catch | $957\pm26$ | $312\pm63$ | $772\pm241$ | **$913\pm53$** | $710\pm217$ | $246\pm174$ | $542\pm270$ | **$942\pm17$** | **$933\pm24$** |
| Finger Spin | $672\pm76$ | $166\pm128$ | $779\pm108$ | **$901\pm104$** | $560\pm77$ | $341\pm70$ | — | — | **$943\pm59$** |
| Walker Walk | $604\pm317$ | $42\pm12$ | $344\pm132$ | **$612\pm164$** | $221\pm43$ | $277\pm12$ | — | — | **$577\pm208$** |
| Cheetah Run | $228\pm95$ | $103\pm38$ | **$307\pm48$** | **$344\pm67$** | $165\pm123$ | $235\pm137$ | — | — | $222\pm88$ |

# Results

**TD-MPC** is *input-agnostic*; just change *h*

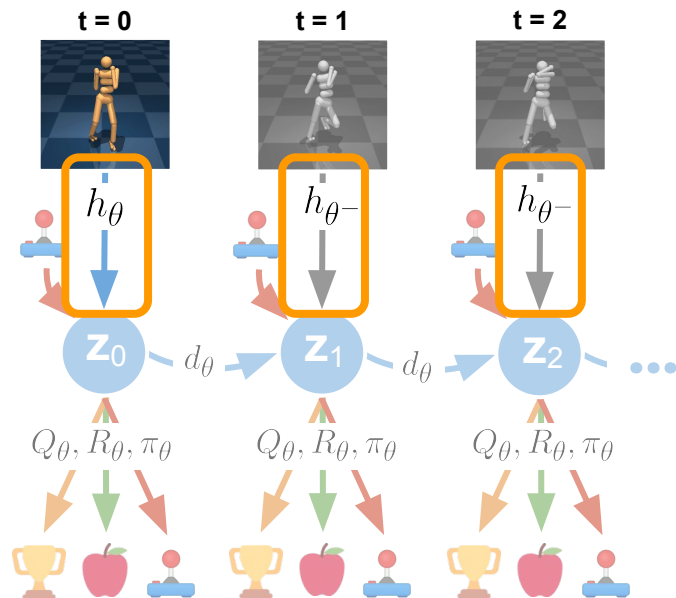- Trivially extended to multi-modal RL

# Results

**TD-MPC** is ***input-agnostic***; just change ***h***

- Trivially extended to multi-modal RL



Proprioceptive data + egocentric camera

# Results

**TD-MPC** matches the ***time to solve*** of SAC but uses far less data

| *Wall-time* (h) | Walker Walk | | | | Humanoid Stand | |
| --- | --- | --- | --- | --- | --- | --- |
| | SAC | LOOP | MPC:sim | **TD-MPC** | SAC | **TD-MPC** |
| time to solve ↓ | 0.41 | 7.72 | 0.91 | 0.47 | 9.31 | 9.39 |
| h/500k steps ↓ | 1.41 | 18.5 | — | 5.60 | 1.82 | 12.94 |

**nicklashansen.github.io/td-mpc**