



ICML

International Conference
On Machine Learning

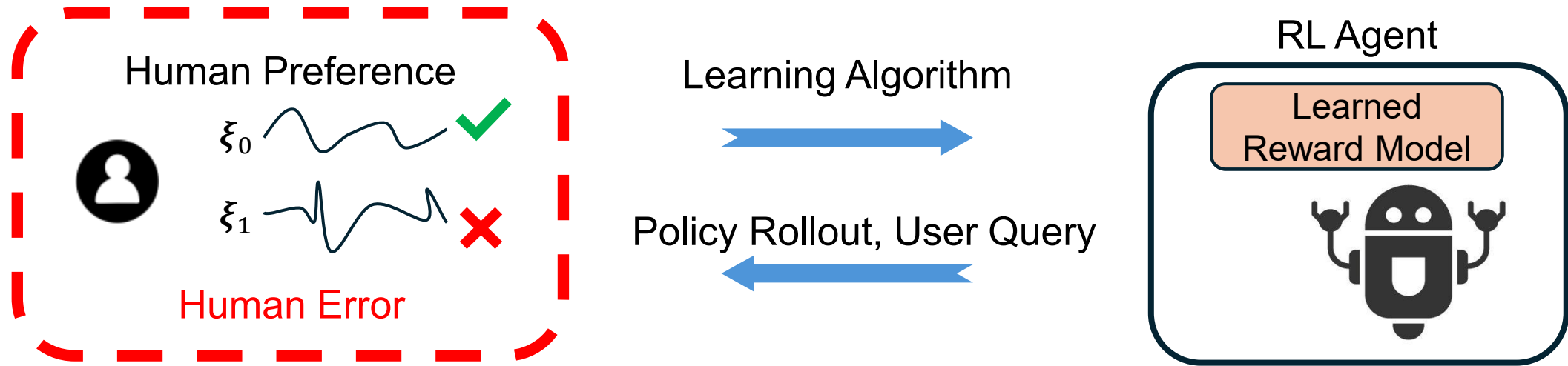
Robust Reward Alignment via Hypothesis Space Batch Cutting

Zhixian Xie, Haode Zhang, Yizhe Feng, Wanxin Jin



SHANGHAI JIAO TONG
UNIVERSITY

Motivation: Preference-Based RL Suffers from Erroneous Feedback



- B-Pref Benchmark: Up to 20% downgrade with 10% error rate¹

Recent effort: Filtering False Feedback \longrightarrow Prior Knowledge
Computation Overhead

Target:

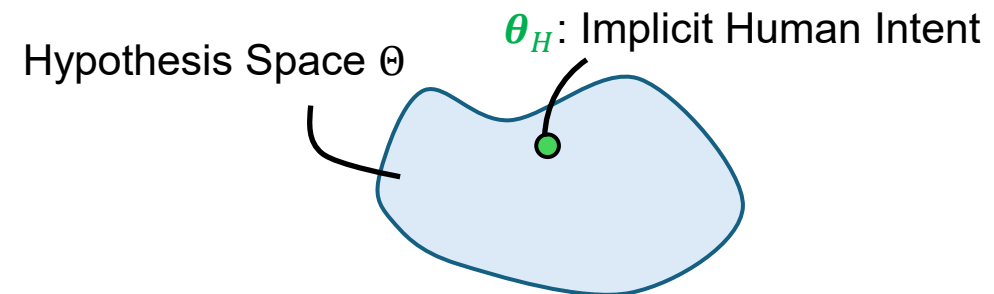
1. **Certifiable human data efficiency**
2. **Provable robustness**

Without explicitly identifying the error.

¹: Lee, Kimin, et al. "B-pref: Benchmarking preference-based reinforcement learning." arXiv preprint arXiv:2111.03026 (2021).

Preferences as Hypothesis Space Cuts

Hypothesis Space: The space of reward function (parameters), containing implicit human intent.



Searching for θ_H : Use preferences to induce cuts and remove the hypothesis space inconsistent with preference!

(a) Trajectory Preference Pair (ξ_0, ξ_1)



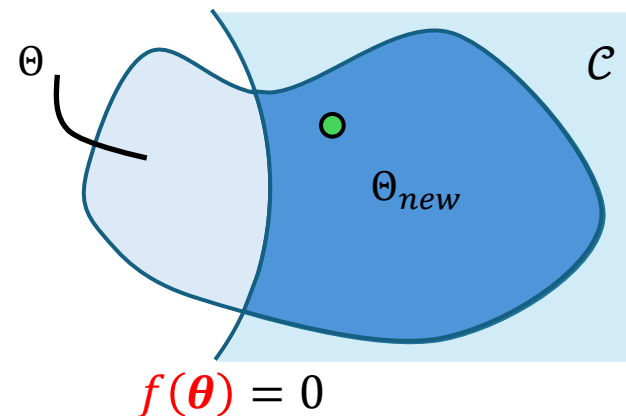
$$r_{\theta_H}(\xi_0) \geq r_{\theta_H}(\xi_1)$$

(b) Cutting Set

$$f(\theta) = r_{\theta}(\xi_0) - r_{\theta}(\xi_1)$$

$$\theta_H \in \mathcal{C} = \{\theta \mid f(\theta) \geq 0\}$$

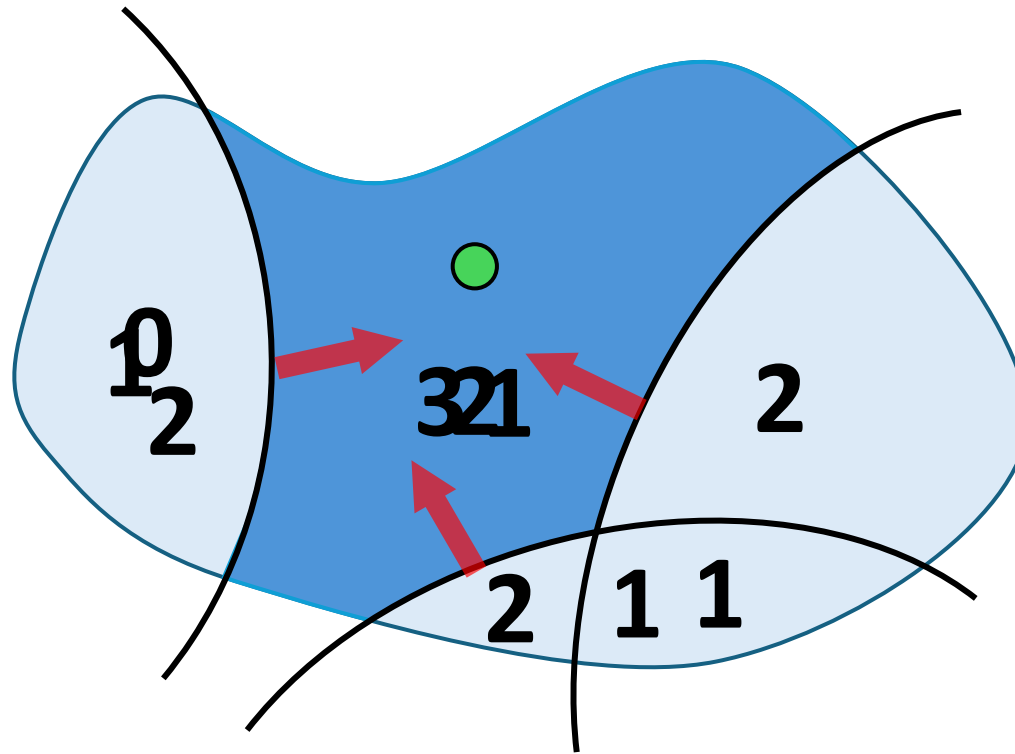
(c) Hypothesis Space Cut
 $\Theta_{new} = \Theta \cap \mathcal{C}$



Batched Cutting as Voting

Use a batch of preference to vote, update the hypothesis space base on vote function:

$$V(\theta) = \# \text{ of satisfied cuts for given } \theta$$

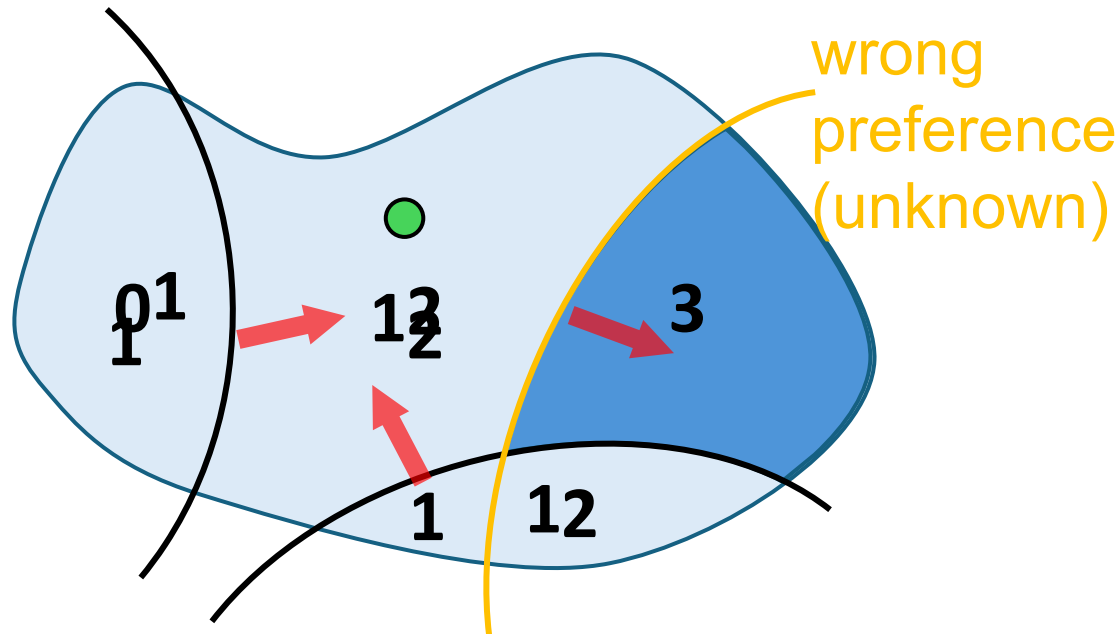


Update the hypothesis space using threshold $V(\theta) \geq 3$!

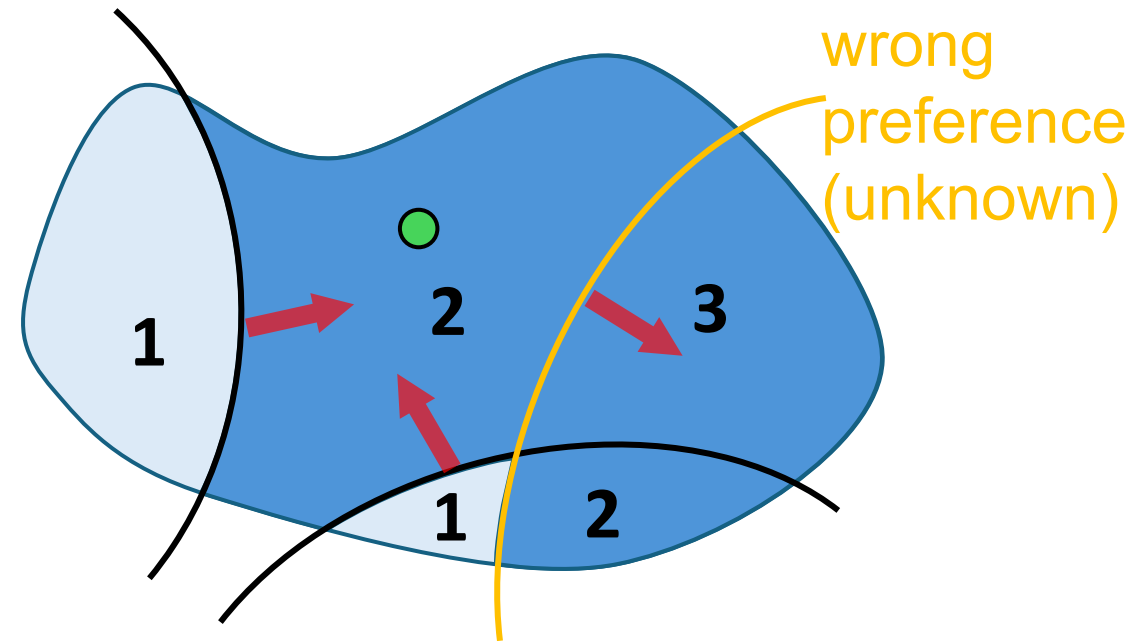
Key Idea: Batched Cutting with Conservativeness

Use Lower threshold to perform conservative update:

Threshold $V(\theta) \geq 3$: Cut out θ_H



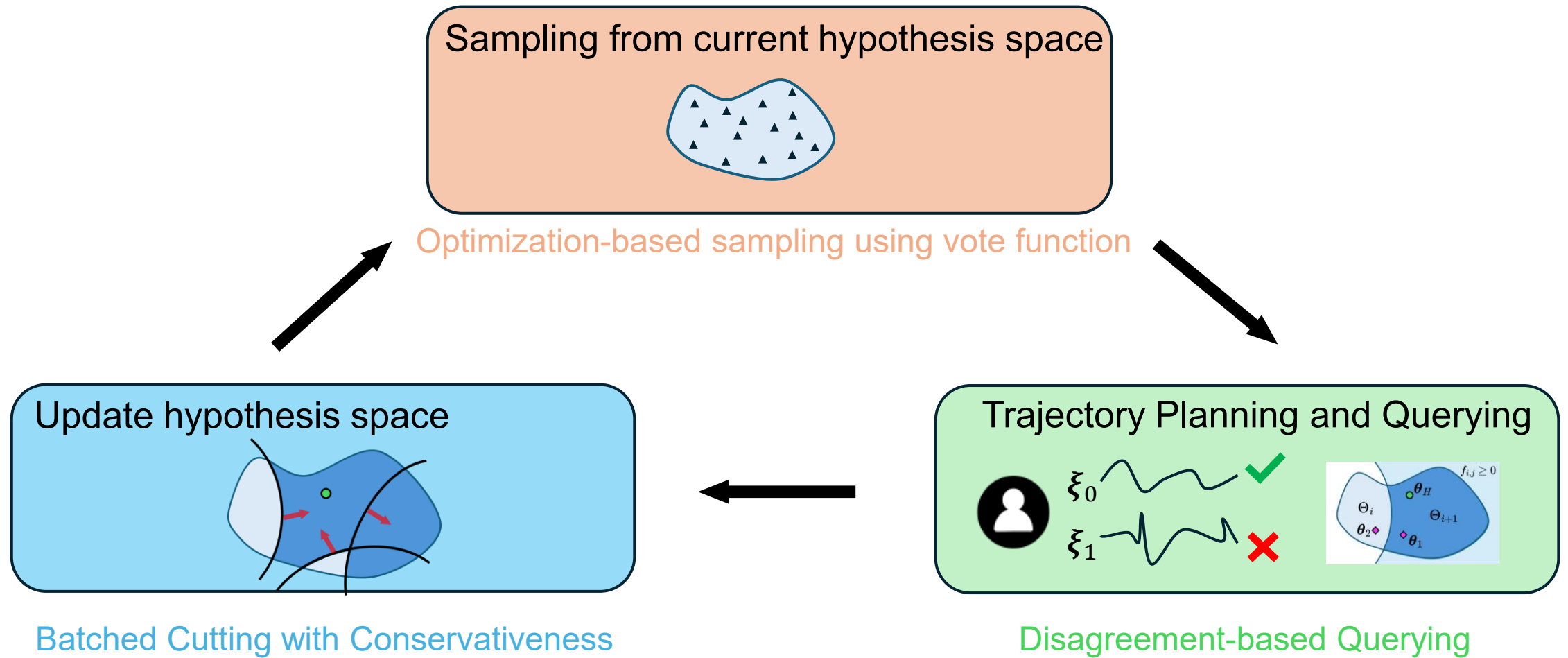
Threshold $V(\theta) \geq 2$: Preserve θ_H in update



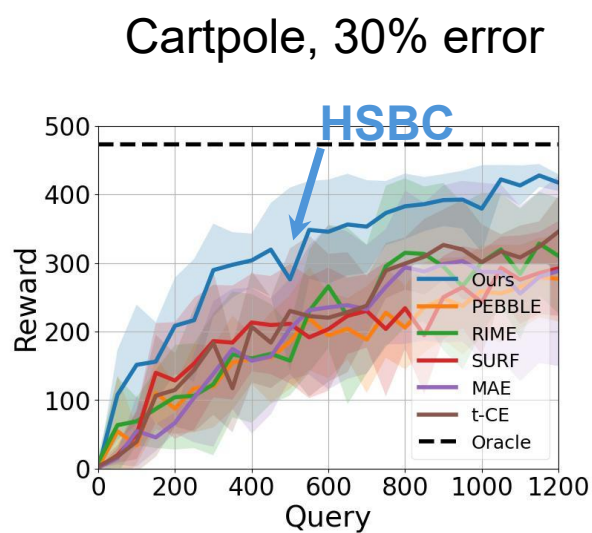
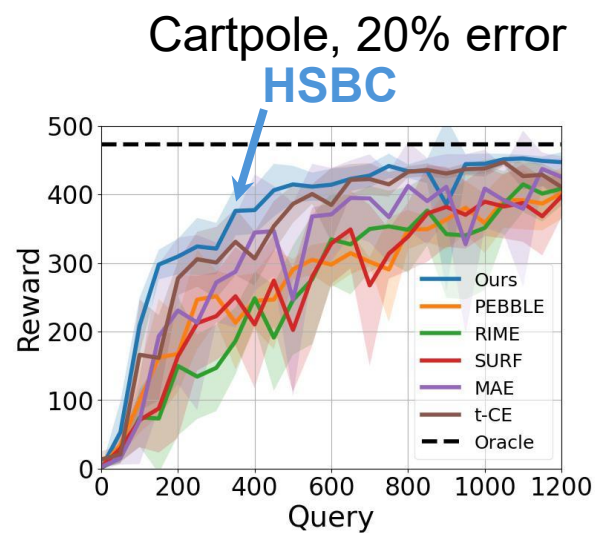
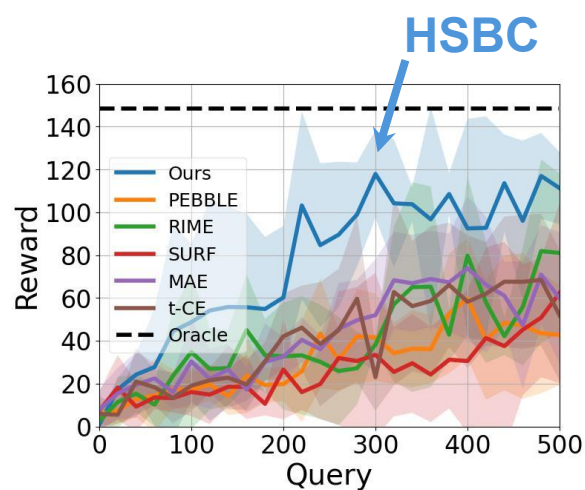
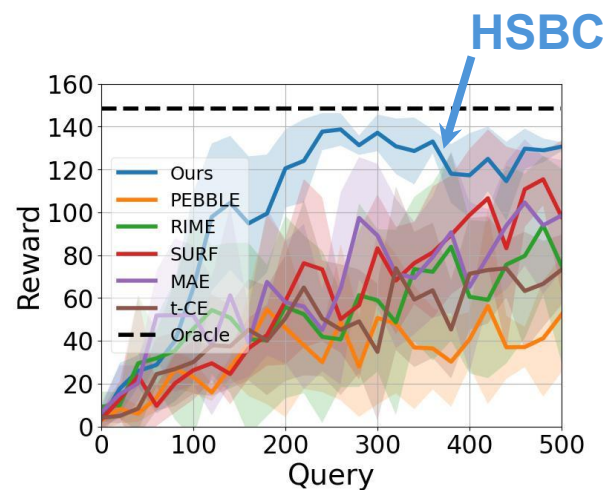
With batch size N : $V(\theta) \geq (1 - \gamma)N$, conservativeness $0 < \gamma < 1$

Provably Robustness, Error Agnostic!

Hypothesis Space Batch Cutting (HSBC) Overview

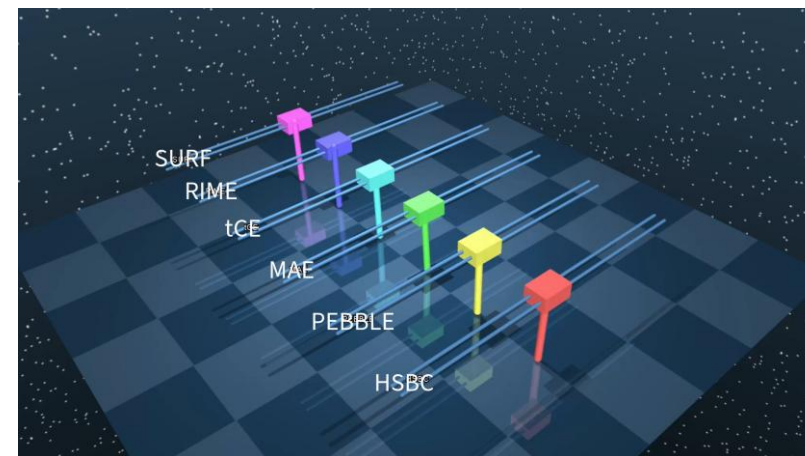


Result

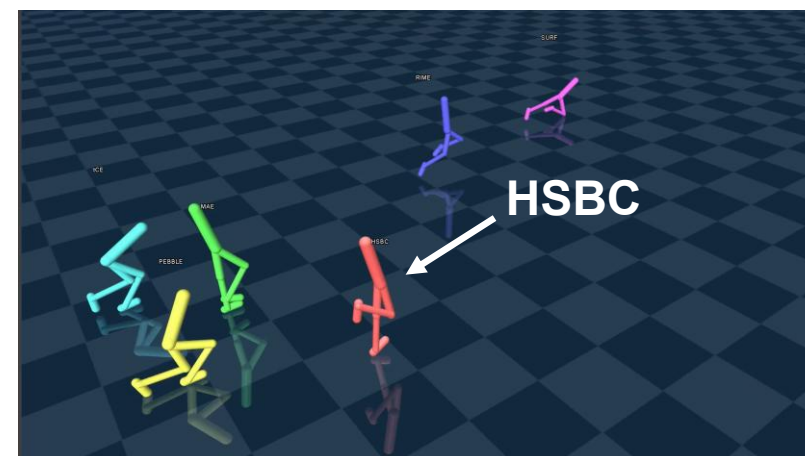


Walker, 20% error

Walker, 30% error



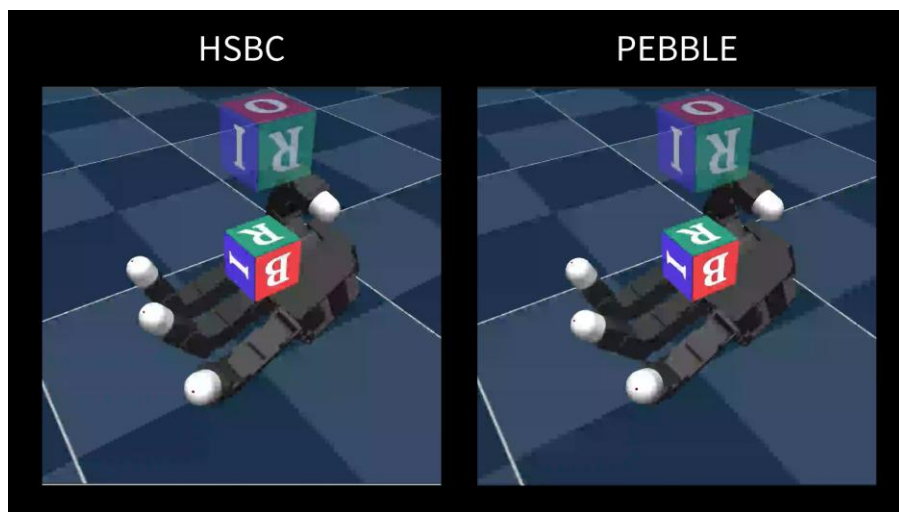
Cartpole, 30% error



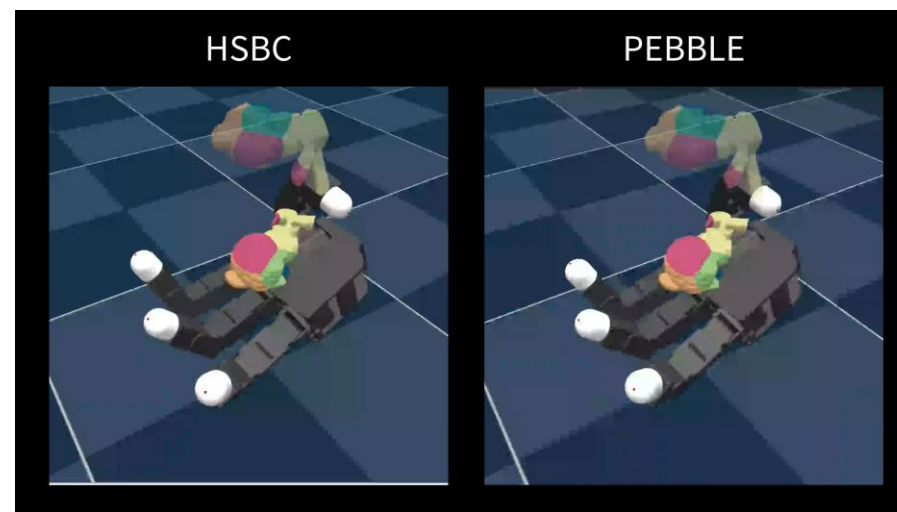
Walker, 30% error

Result

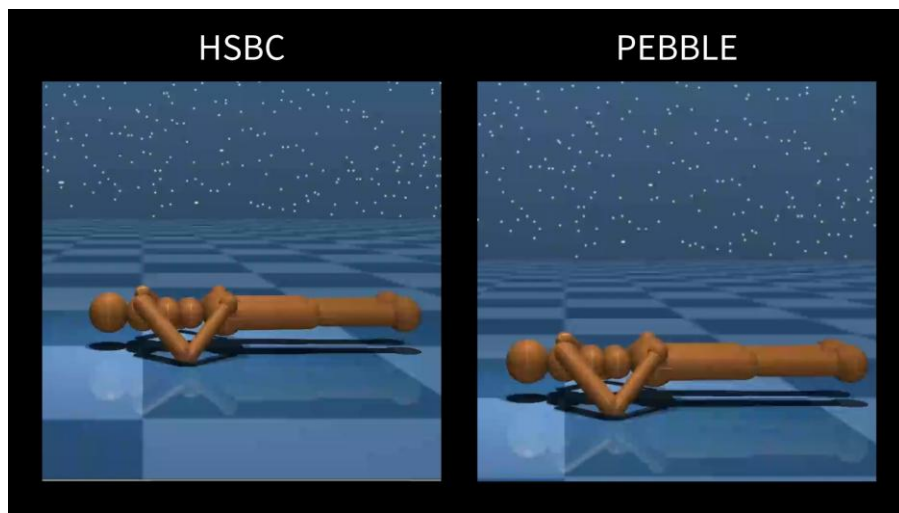
Manipulation: Cube



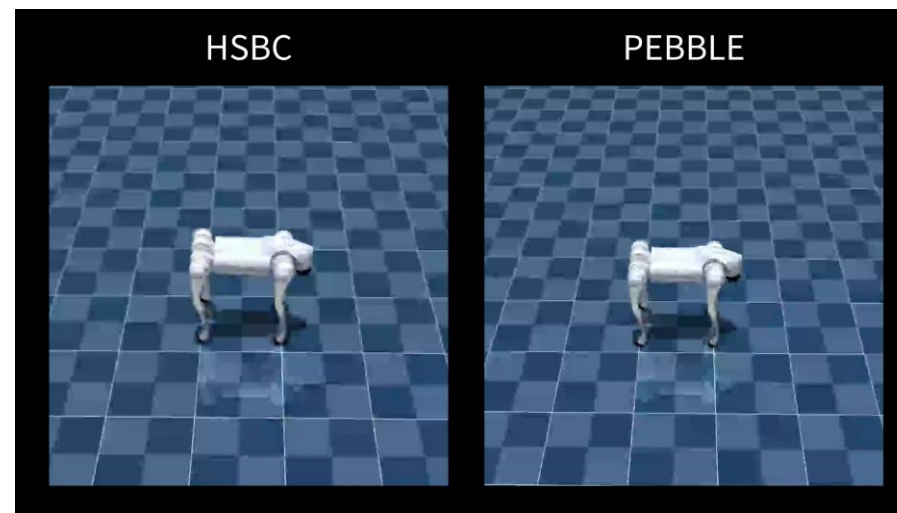
Manipulation: Bunny



Humanoid



Go2



Thanks for
Listening!

Project Page

