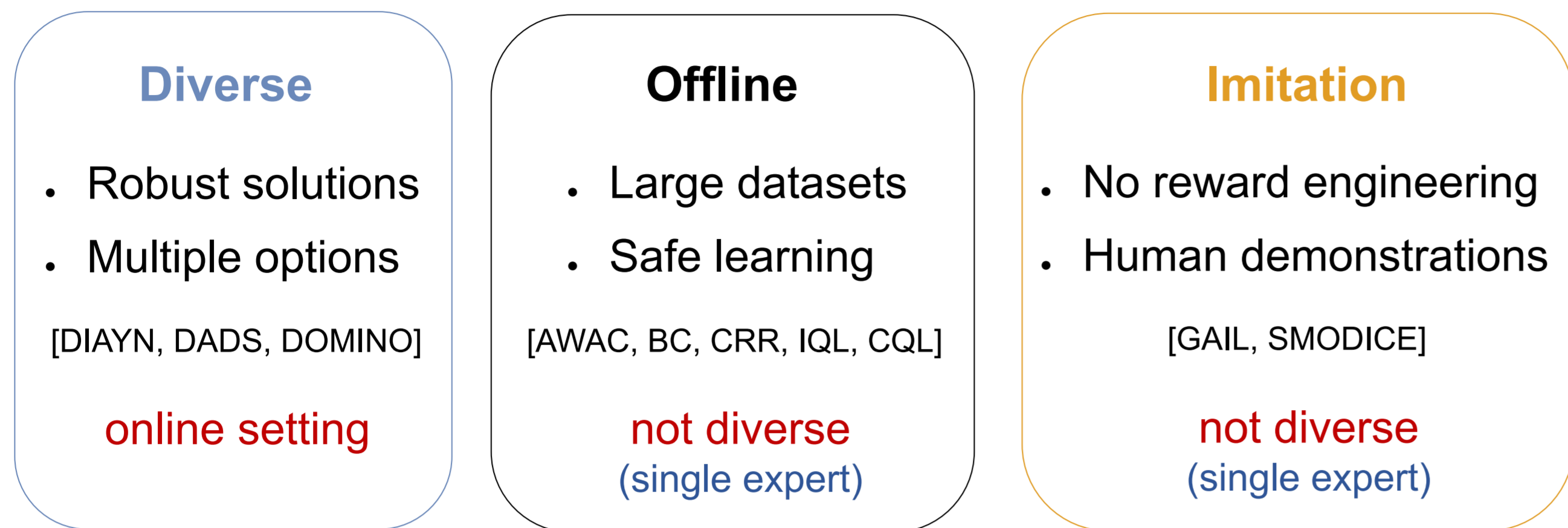


Offline Diversity Maximization Under Imitation Constraints

Marin Vlastelica, Jin Cheng, Georg Martius, Pavel Kolev



Overview



Propose: principled algorithm for **Diverse Offline Imitation** (DOI) learning

Diverse Offline Imitation

$$\max_{\{d_z(S)\}_{z \in Z}} \mathcal{I}(S; Z)$$

subject to

$$D_{\text{KL}}(d_z(S) \| d_E(S)) \leq \epsilon \quad \forall z$$

Input: $\mathcal{D}_E \sim d_E(s)$ state-only expert dataset
 $\mathcal{D}_O \sim d_O(s, a)$ state-action behavior dataset

Def. State-action occupancy

$$d^\pi(s, a) := (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \Pr[s_t = s, a_t = a \mid s_0 \sim \rho_0, a_t \sim \pi(\cdot | s_t), s_{t+1} \sim \mathcal{P}(\cdot | s_t, a_t)]$$

Relaxed Problem Formulation

$$\max_{\{d_z(S)\}_{z \in Z}} \mathcal{I}(S; Z)$$

Mutual Info. variational lower bound

$$\mathcal{I}(S; Z) \geq \sum_z \mathbb{E}_{d_z(s)} \left[\frac{\log(|Z|q(z|s))}{|Z|} \right]$$

$q(z|s)$ train a skill-discriminator

subject to

$$D_{\text{KL}}(d_z(S, A) \| d_{\tilde{E}}(S, A)) \leq \epsilon \quad \forall z$$

SMODICE expert (offline)

$$d_{\tilde{E}}(S, A) \approx \arg \min_{d(S, A)} D_{\text{KL}}(d(S, A) \| d_E(S, A))$$

subject to

$$D_{\text{KL}}(d_z(S, A) \| d_{\tilde{E}}(S, A)) \leq \epsilon \quad \forall z$$

SMODICE expert (offline)

$$\eta_{\tilde{E}}(s, a) = \frac{d_{\tilde{E}}(s, a)}{d_O(s, a)}$$

Algorithmic Approach

Diversity

$$\max_{d_z(s, a)} \min_{\lambda \geq 0} \sum_{z \in Z} \mathbb{E}_{d_z(s, a)} \left[\frac{\log(|Z|q(z|s))}{|Z|} \right] + \sum_{z \in Z} \lambda_z [\epsilon - D_{\text{KL}}(d_z(S, A) \| d_{\tilde{E}}(S, A))]$$

Imitation

$$\max_{d_z(s, a)} \min_{\lambda > 0} \sum_z \lambda_z \left\{ \epsilon + \mathbb{E}_{d_z(s, a)} [R_z^\lambda(s, a)] - D_{\text{KL}}(d_z(S, A) \| d_O(S, A)) \right\}$$

Regularized RL Problem

DICE (offline)

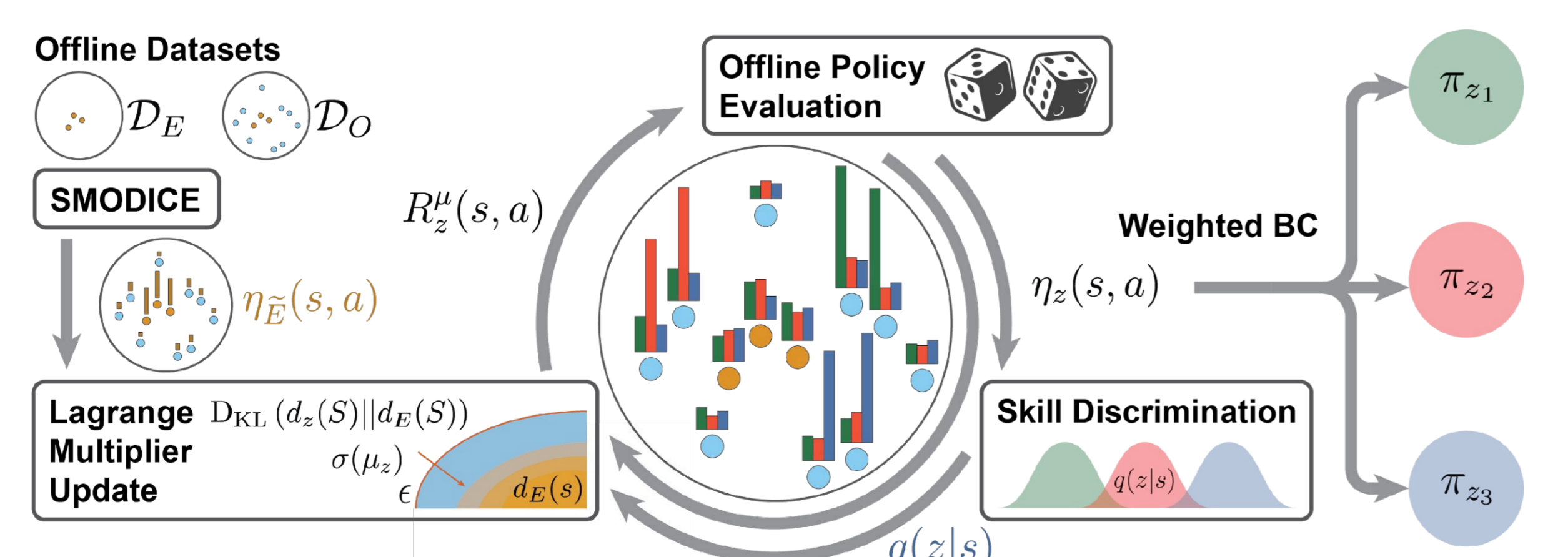
$$\eta_z(s, a) = \frac{d_z(s, a)}{d_O(s, a)}$$

SMODICE expert (offline)

$$\eta_{\tilde{E}}(s, a) = \frac{d_{\tilde{E}}(s, a)}{d_O(s, a)}$$

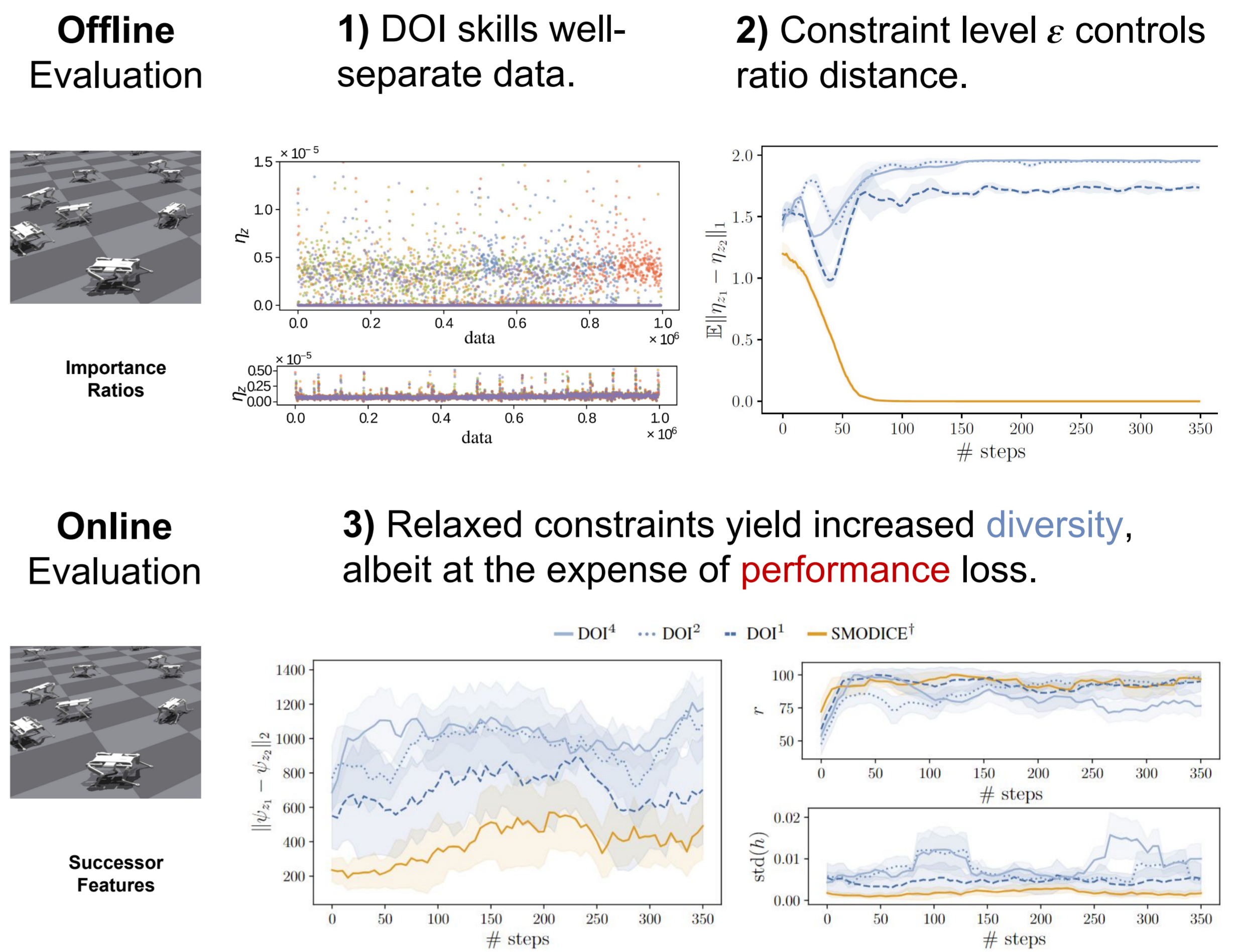
Constraint Violation Skill Diversity Expert Imitation

Alternating Optimization Scheme



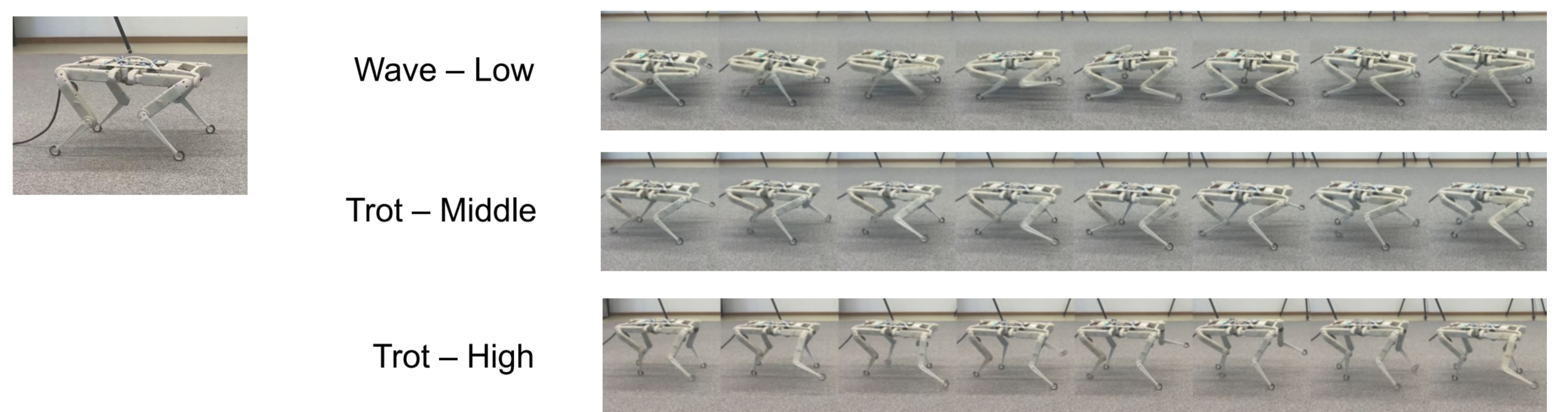
Experiments (Solo12)

I. Locomotion Task (SIM)



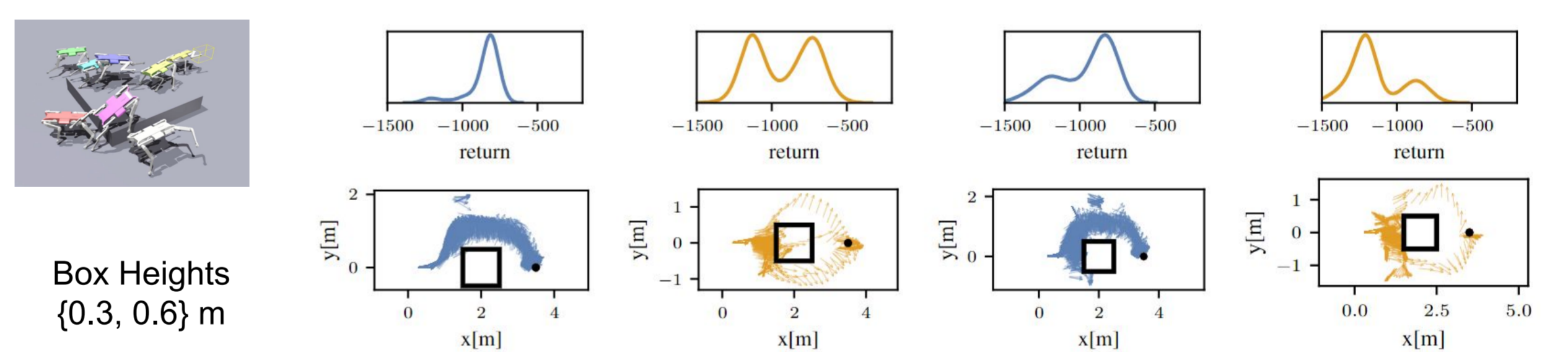
I. Locomotion Task (REAL)

4) DOI skills trained in **SIM** (with domain randomization) are successfully deployed in the **Real System**.

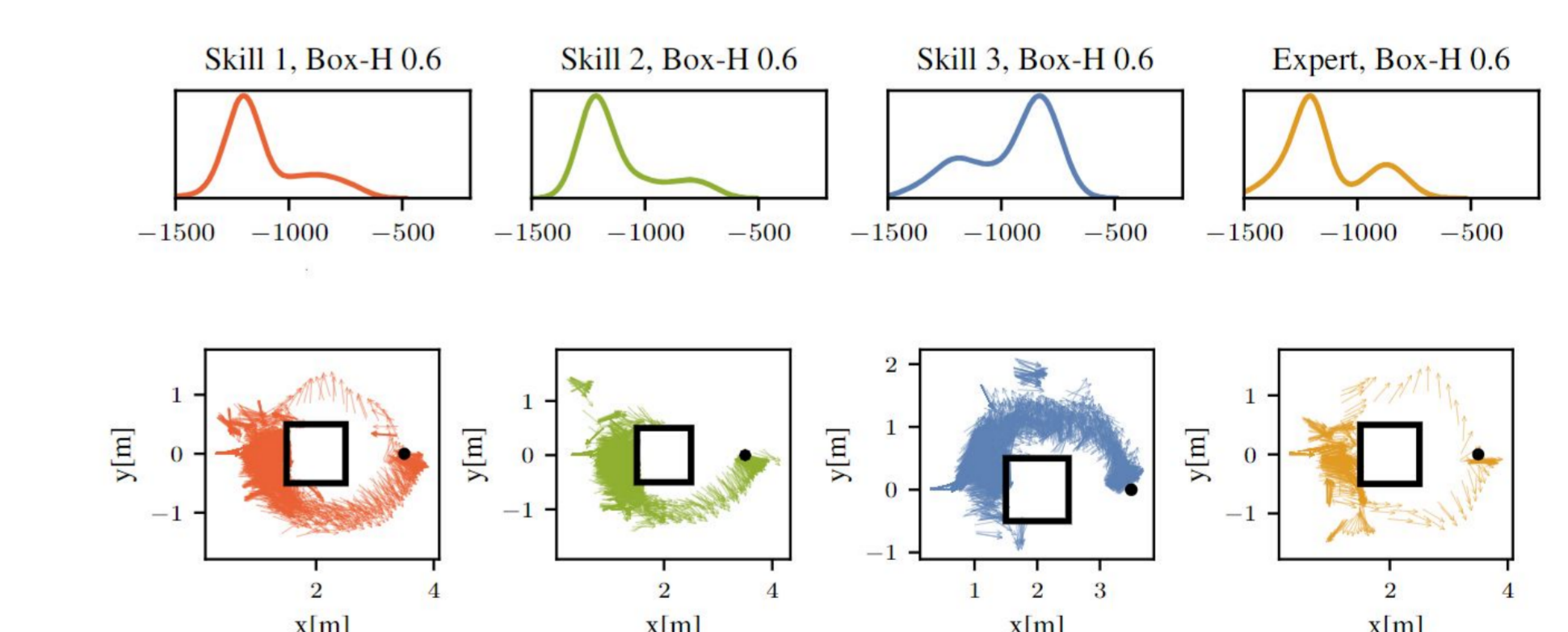


II. OBSTACLE NAVIGATION TASK (SIM)

5) SMODICE expert struggles with out-of-distribution (higher) box heights, while a robust DOI skill successfully navigates by detouring (to the left side).



Limitation: 6) Not all learned DOI skills are robust. Selection is required.



Contact: {mvlastelica, pkolev}@uni-tuebingen.de



Paper



Project Website



Autonomous Learning

References

- [SMODICE] Y. Ma, A. Shen, D. Jayaraman, O. Bastani "Versatile Offline Imitation from Observations and Examples via Regularized State-Occupancy Matching", ICML 2022
- [DICE] O. Nachum, B. Dai "Reinforcement learning via fenchel-rockafellar duality", arXiv 2020
- [DIAYN] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, Sergey Levine "Diversity is All You Need: Learning Skills without a Reward Function", ICLR 2019
- [DOMINO] T. Zahavy, Y. Schroecker, F. Behbahani, K. Baumli, S. Flennerhag, S. Hou, S. Singh "Discovering Policies with DOMINO: Diversity Optimization Maintaining Near Optimality", ICLR 2023
- [OFFLINE] R. F. Prudencio, M. R. O. A. Maximo, and E. L. Colombini. "A survey on offline reinforcement learning: Taxonomy, review, and open problems", arXiv 2022