

Reinforcement Learning from Demonstrations

Methods and Applications in Surgical Digital Twin Simulations

Doctoral Examination, 24.06.2024

Candidate: Ivan Ovinnikov

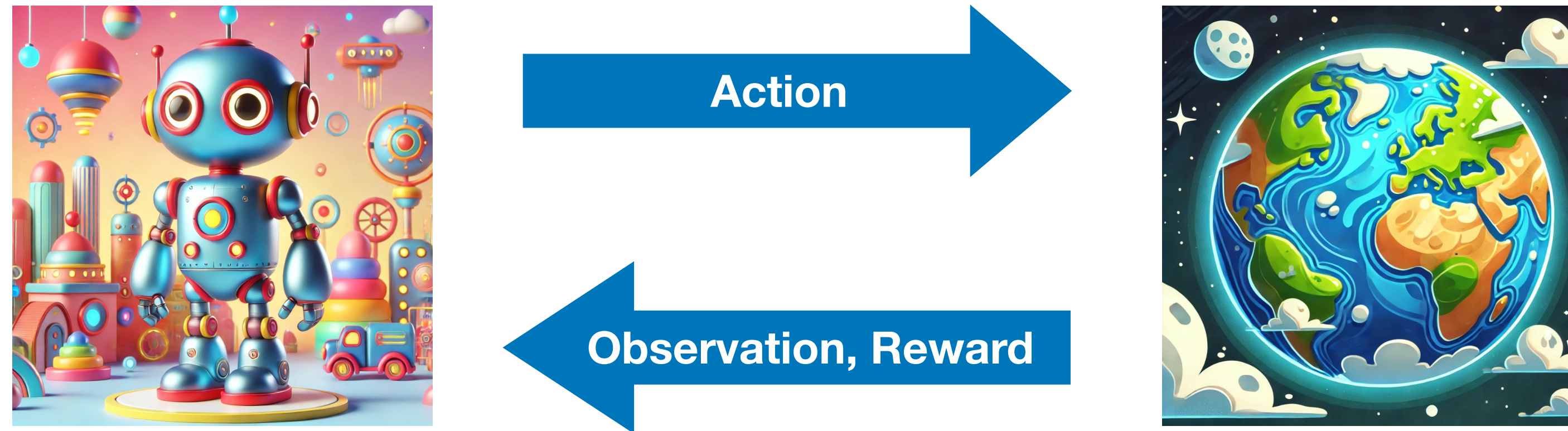
Supervisor: Prof. Dr. Joachim M. Buhmann

Co-examiner: Prof. Dr. Andreas Krause

Co-examiner: Dr. Raimundo Sierra

Chair: Prof. Dr. Markus Gross

Reinforcement Learning as Behaviour Formalism



Main Objective: maximize gathered reward



Image sources: ChatGPT, Google Deepmind, NVIDIA IsaacGym

Digital twins

Definition & Examples

A set of virtual information constructs that mimics the structure, context and behaviour of an individual / unique physical asset, or a group of physical assets, is dynamically updated with data from its physical twin throughout its life cycle and informs decisions that realize value.

AIAA Position Paper, 2020

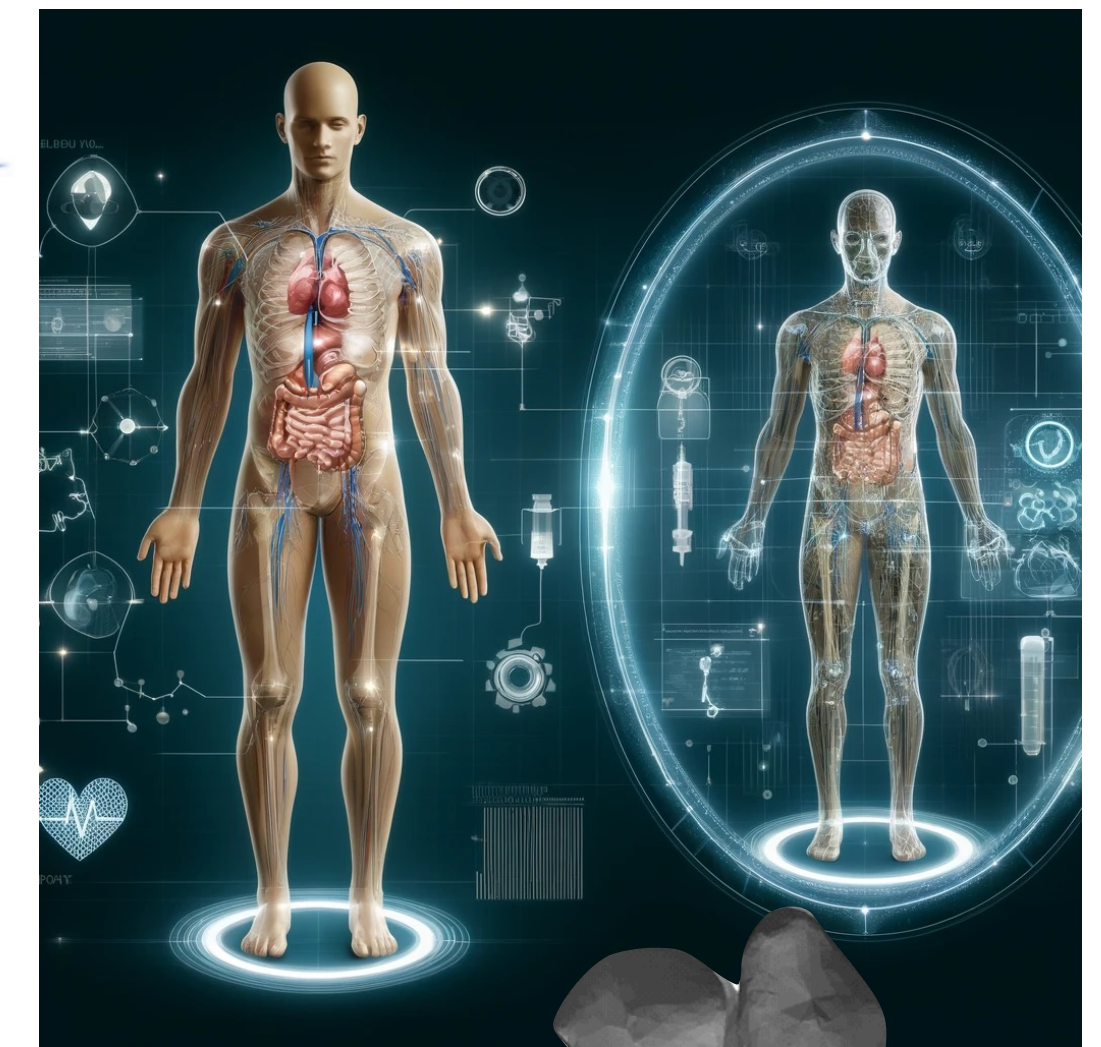
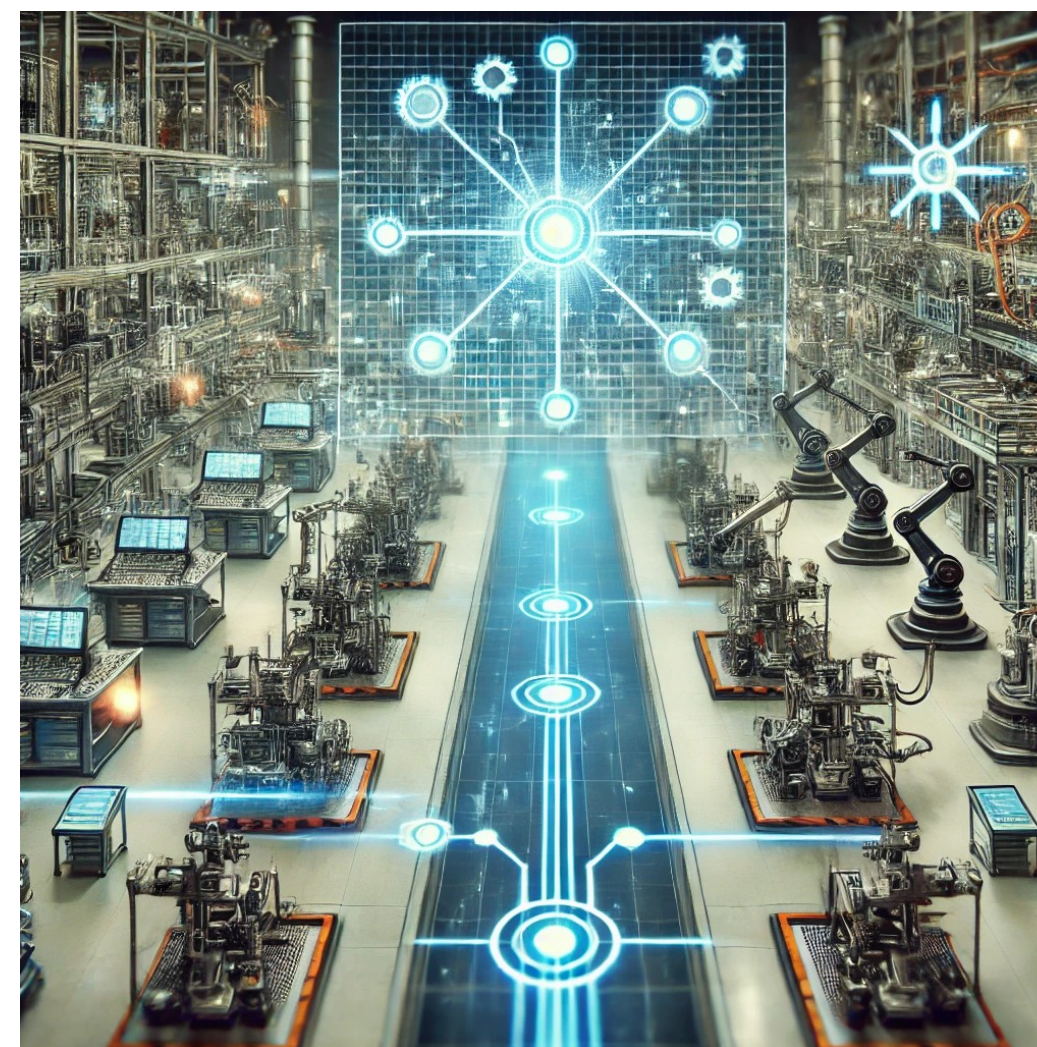
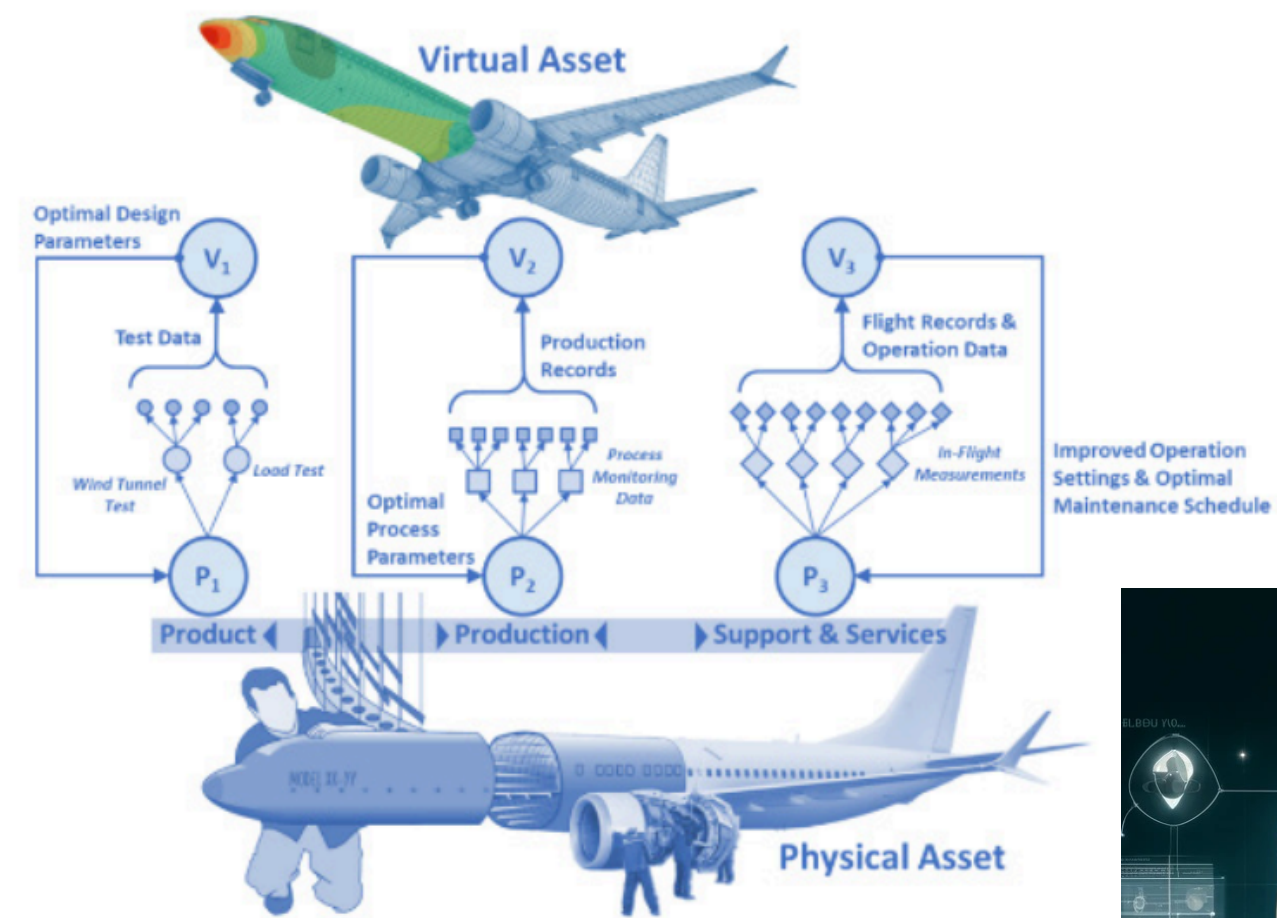
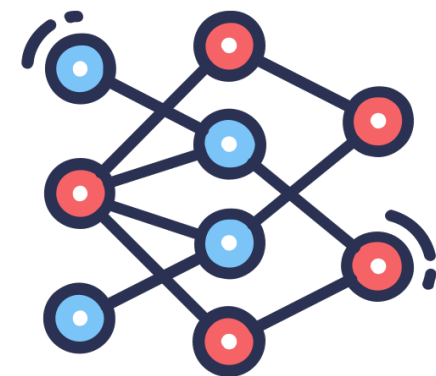


Image sources: AIAA, ChatGPT, F. Laumer

RL in digital twins

Environment models for agents

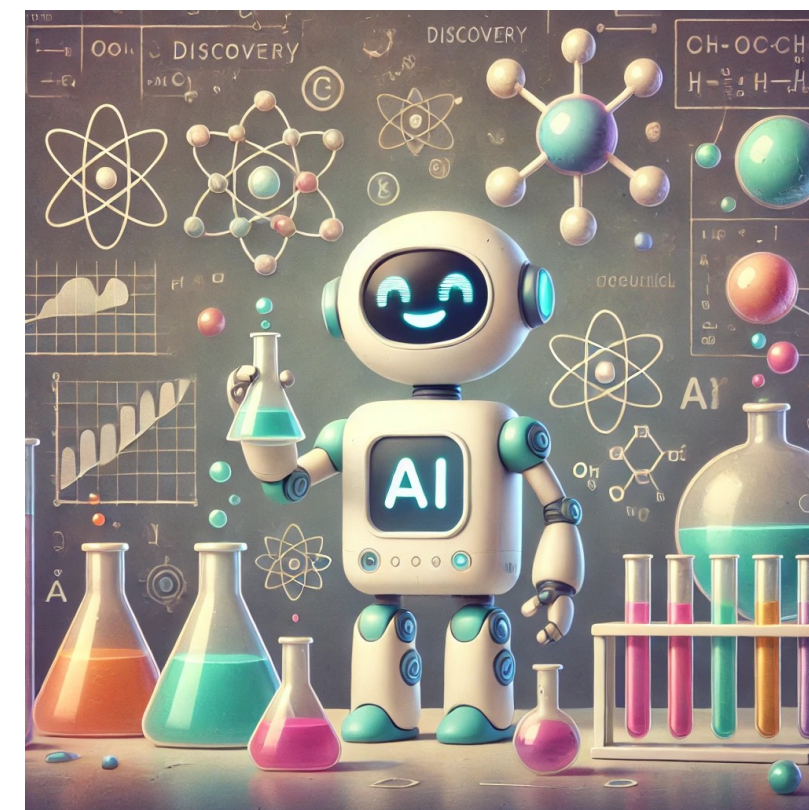
Digital twins provide an interactive data source to learning models



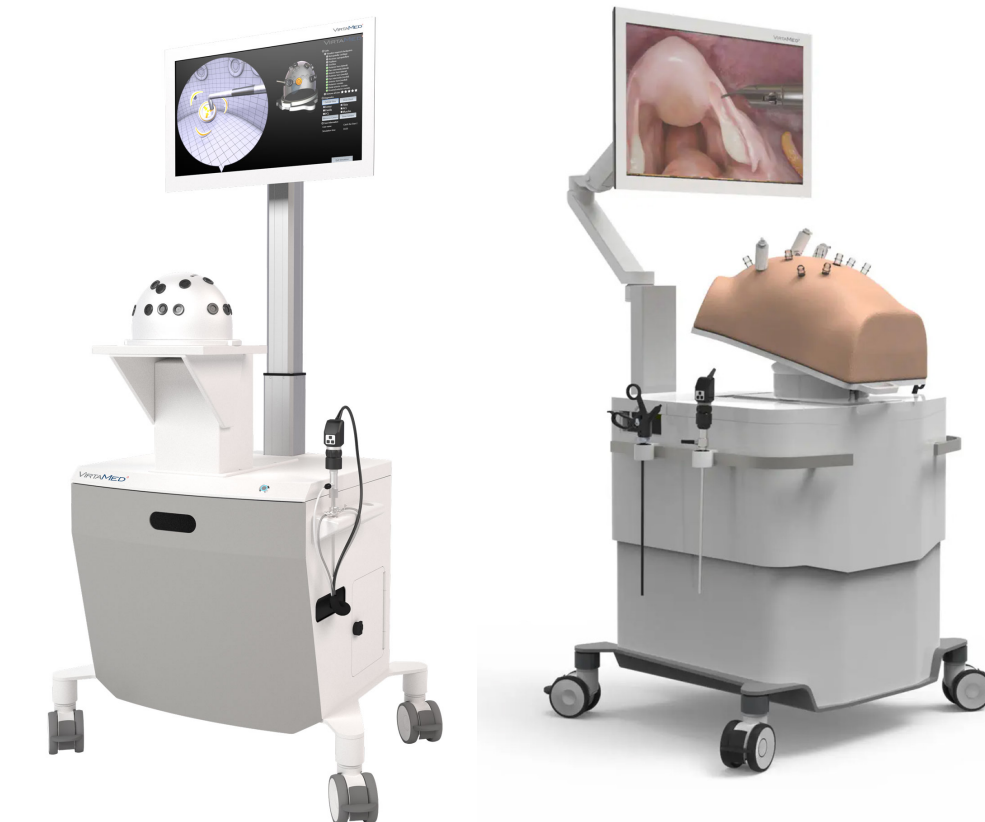
```
1 iow staff 21K 21 Jan 00:13 td3++.py
1 iow staff 18K 26 Feb 16:28 td3.py
1 iow staff 29K 9 Jan 18:49 td3_airs.py
1 iow staff 11K 29 Dec 18:04 td3_continuous_action.py
1 iow staff 13K 16 Nov 2022 td3_continuous_action_iv.py
1 iow staff 25K 26 Nov 2023 td3_gail.py
1 iow staff 24K 9 Mar 03:55 td3_gail.py
1 iow staff 38K 22 Jan 17:00 td3_swil.py
1 iow staff 47B 2 Sep 2023 test.pkl
1 iow staff 25K 6 Mar 03:26 valuedice.py
1 iow staff 13K 16 Apr 03:18 videos
1 iow staff 41B 10 Apr 03:18 videos
(dev) ~ cleanrl git:(main) x python launcher_iv.py
(dev) ~ cleanrl git:(main) x python lap/gpt_gail_lapLiver.py --verbose --capture
(dev) ~ cleanrl git:(main) x
```



LLM Agents



Scientific Discovery

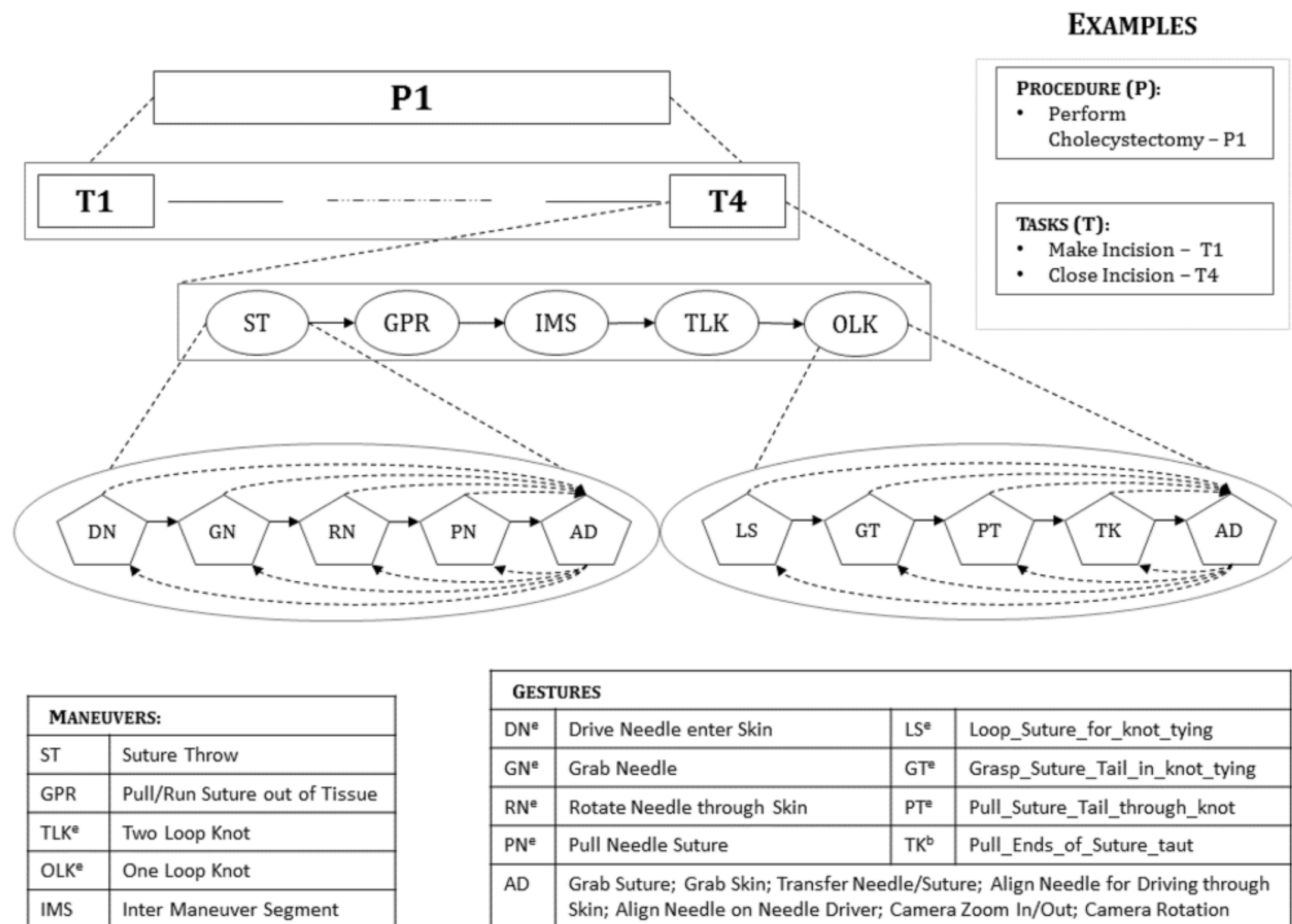


VIRTAMED⁺
WE SIMULATE REALITY

Surgical Simulation

RL in digital twins

Eliciting behaviours in digital twins



Complex environments require increasingly complex behaviours

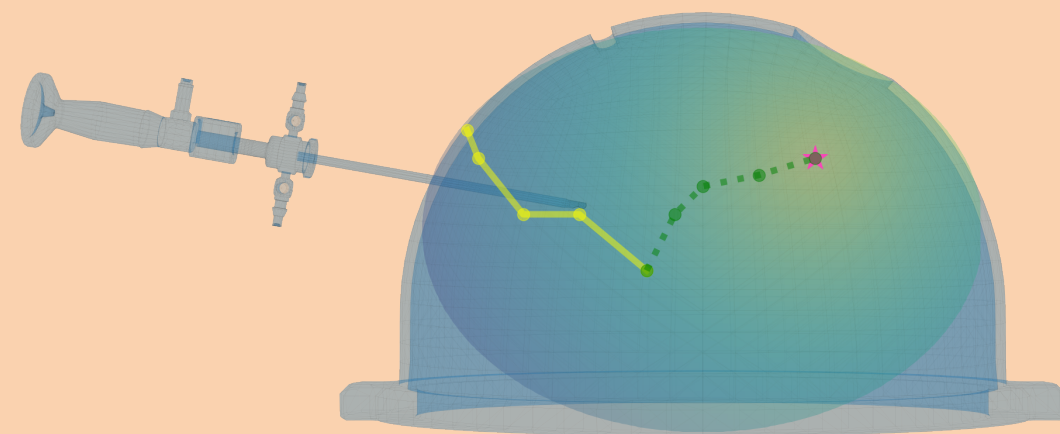


Leverage demonstration data

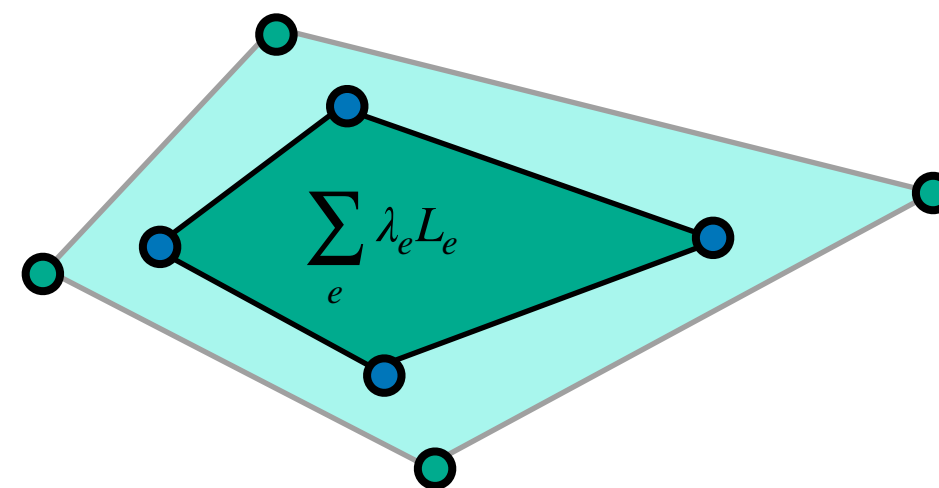
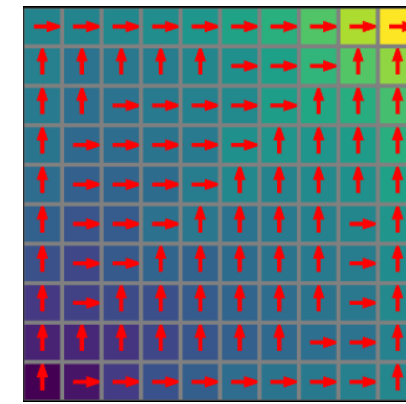
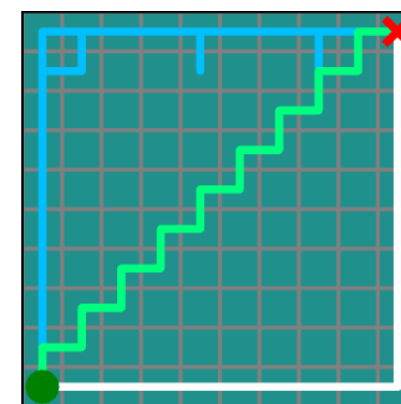
Fig 1. Hierarchical semantic decomposition of surgical activity. ^e denotes that the segment can be performed using either of the robotic arms, ^b denotes that the segment is performed using both the robotic arms.

Overview

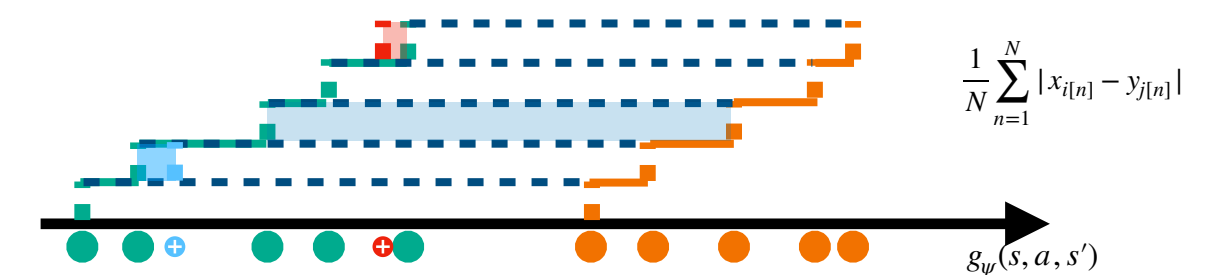
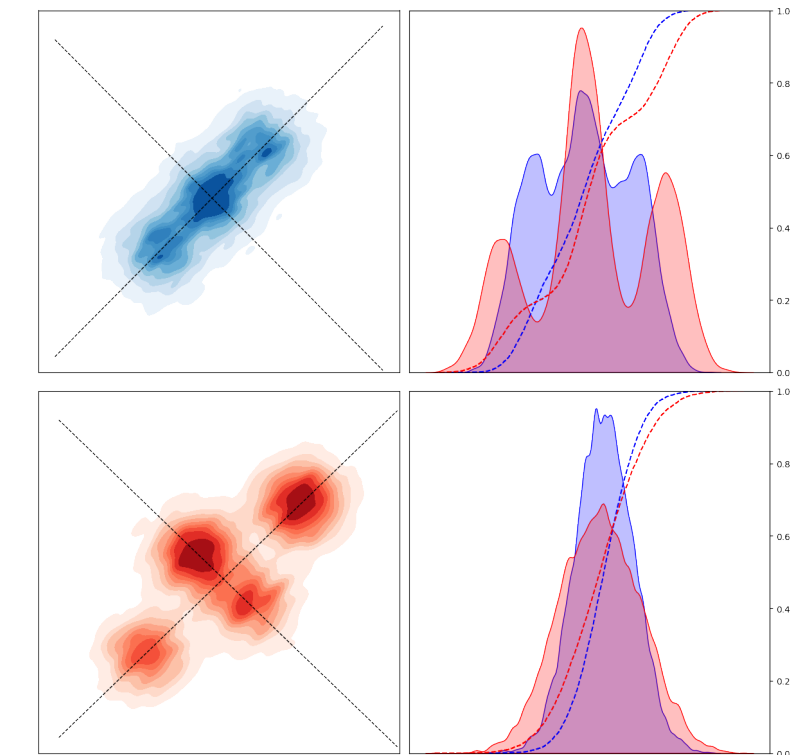
**Application:
Algorithmic RL pipeline in
Surgical Digital Twins**



**Method I:
Addressing Reward
Generalization using Causal
Invariance**

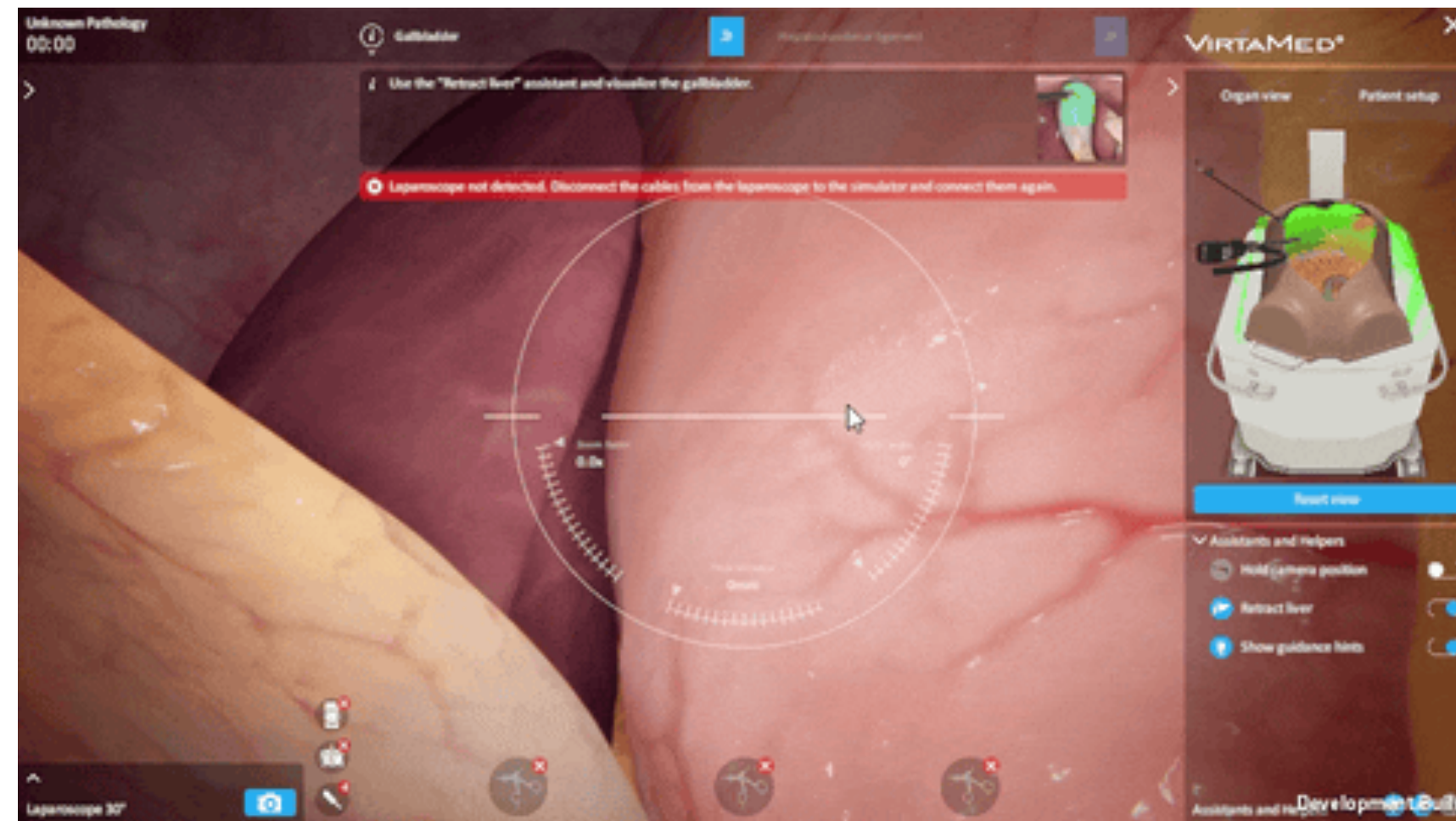


**Method II:
Addressing Data Efficiency in
Imitation Learning using Sliced
Optimal Transport**



RL in Surgical Digital Twins

Research questions

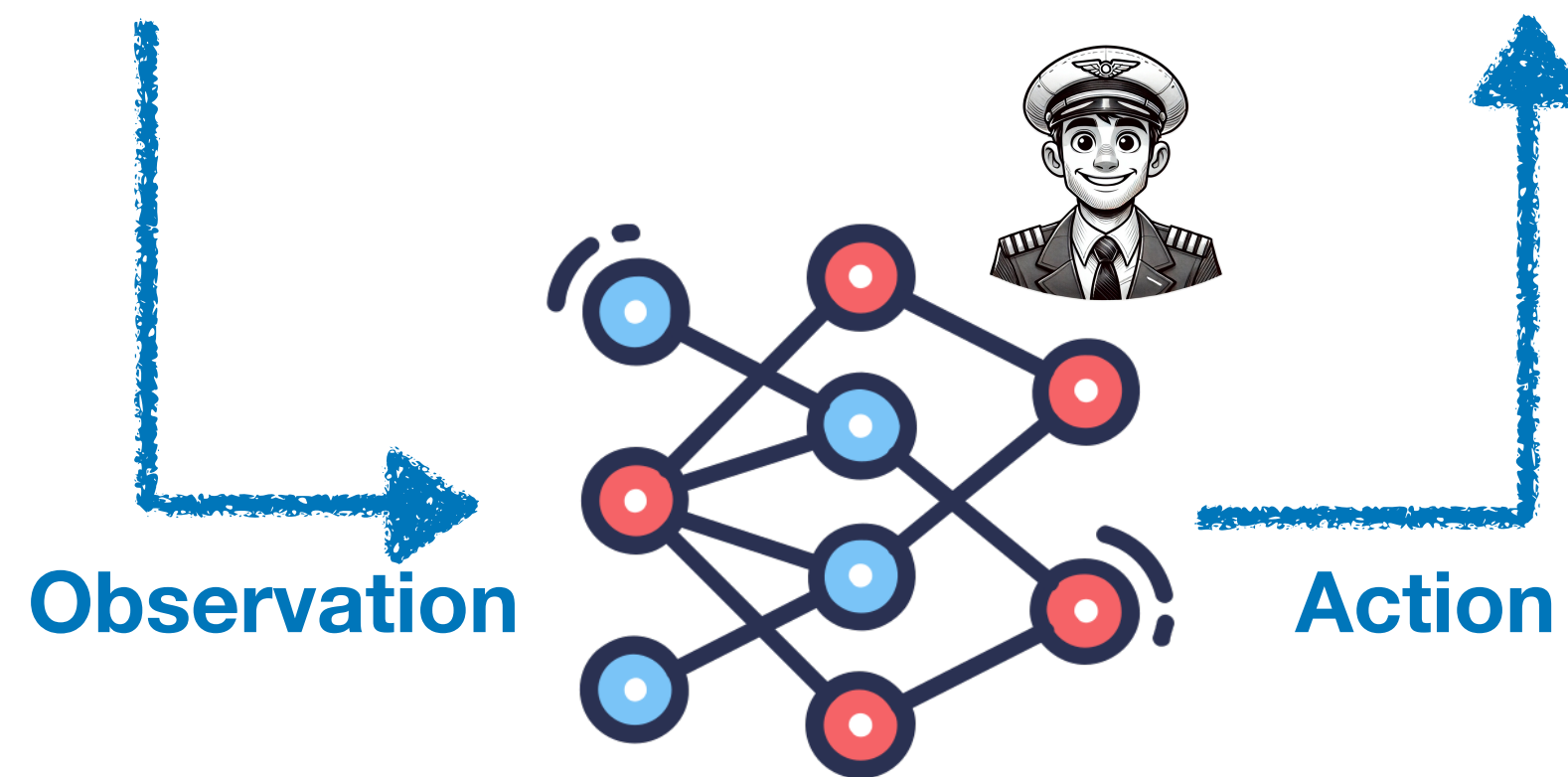


How to leverage ML technology?

- How to evaluate trainees in a data-driven manner?
- How to provide meaningful feedback and assistance?



Devise algorithmic pipeline based on RL agents serving as surgical co-pilots



Behavioural Policy (and Reward)

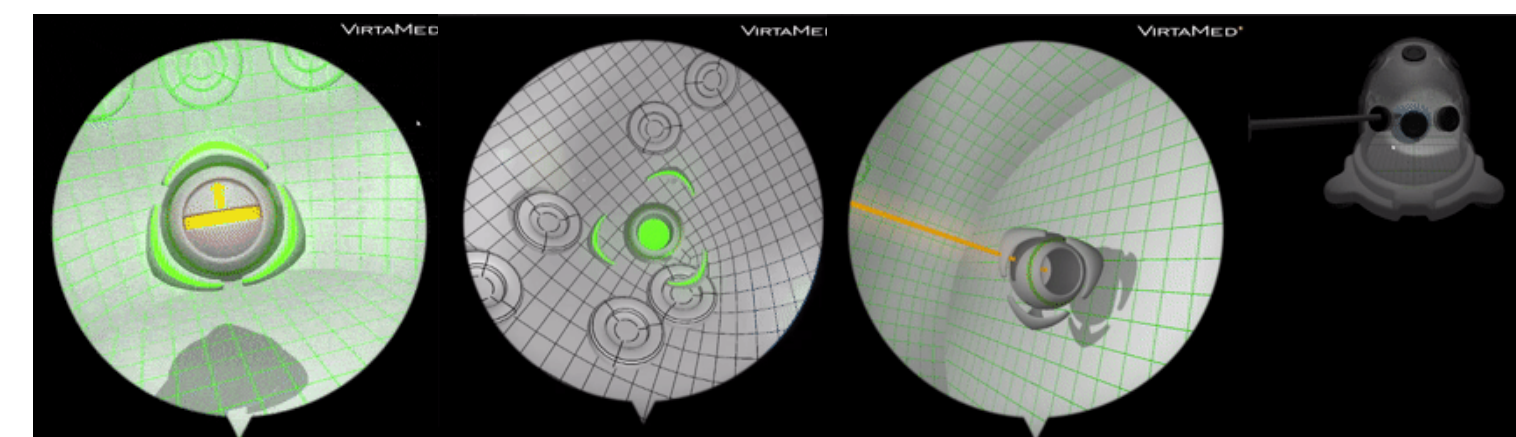
Training surgical agents

Leveraging control loop approach

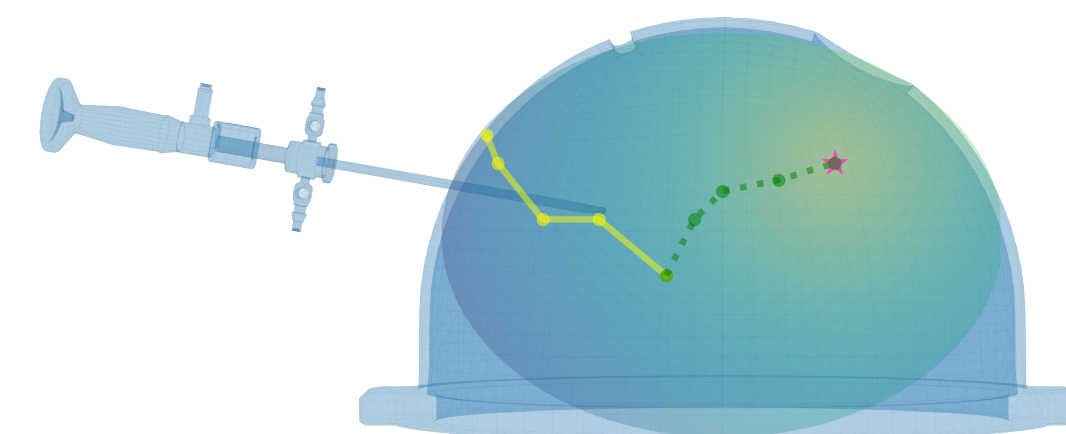


Behavioural Policy (and Reward)

Benchmark suite of tasks for surgical RL agents

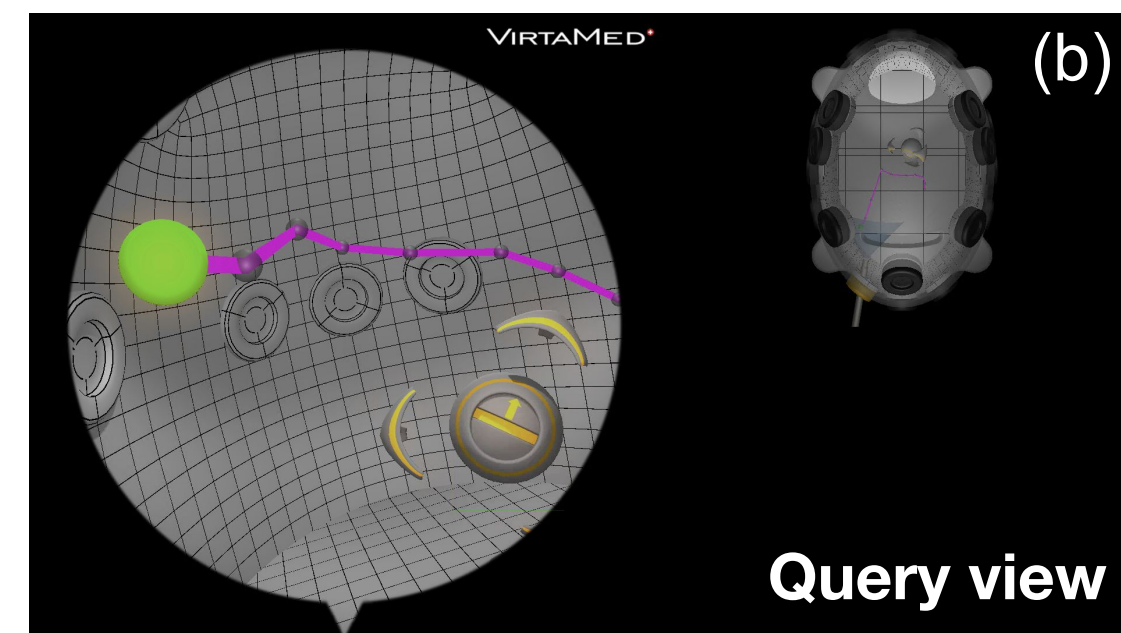
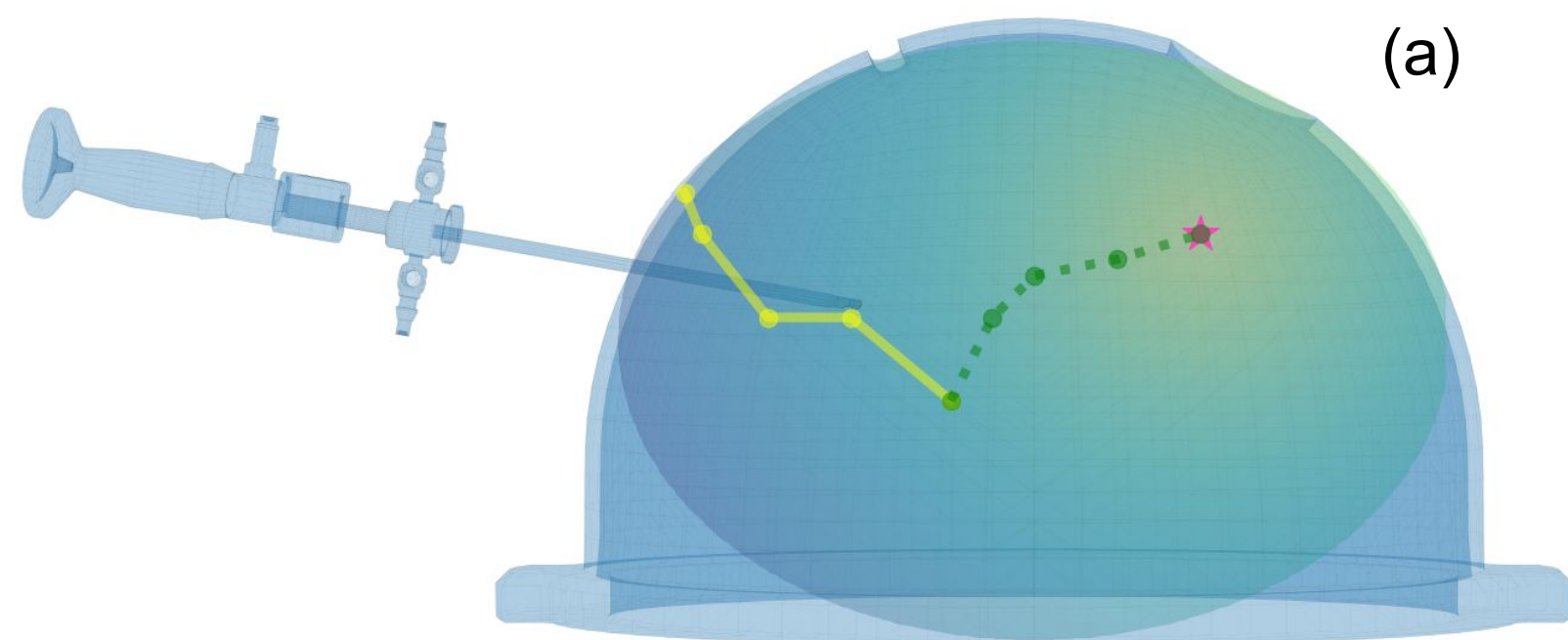
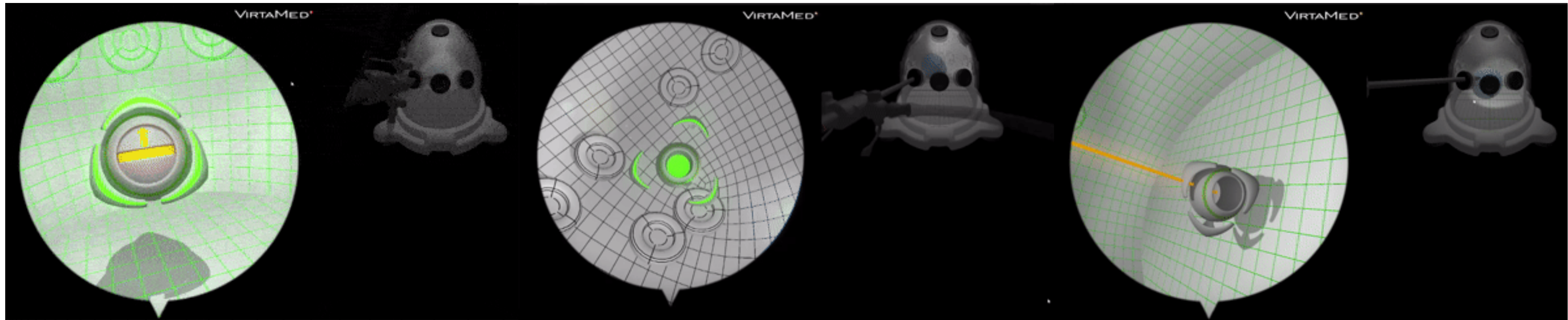


Learning from demonstrations for purposes of evaluation



FASTRL Benchmark suite

Fundamentals of Arthroscopic Surgery Training

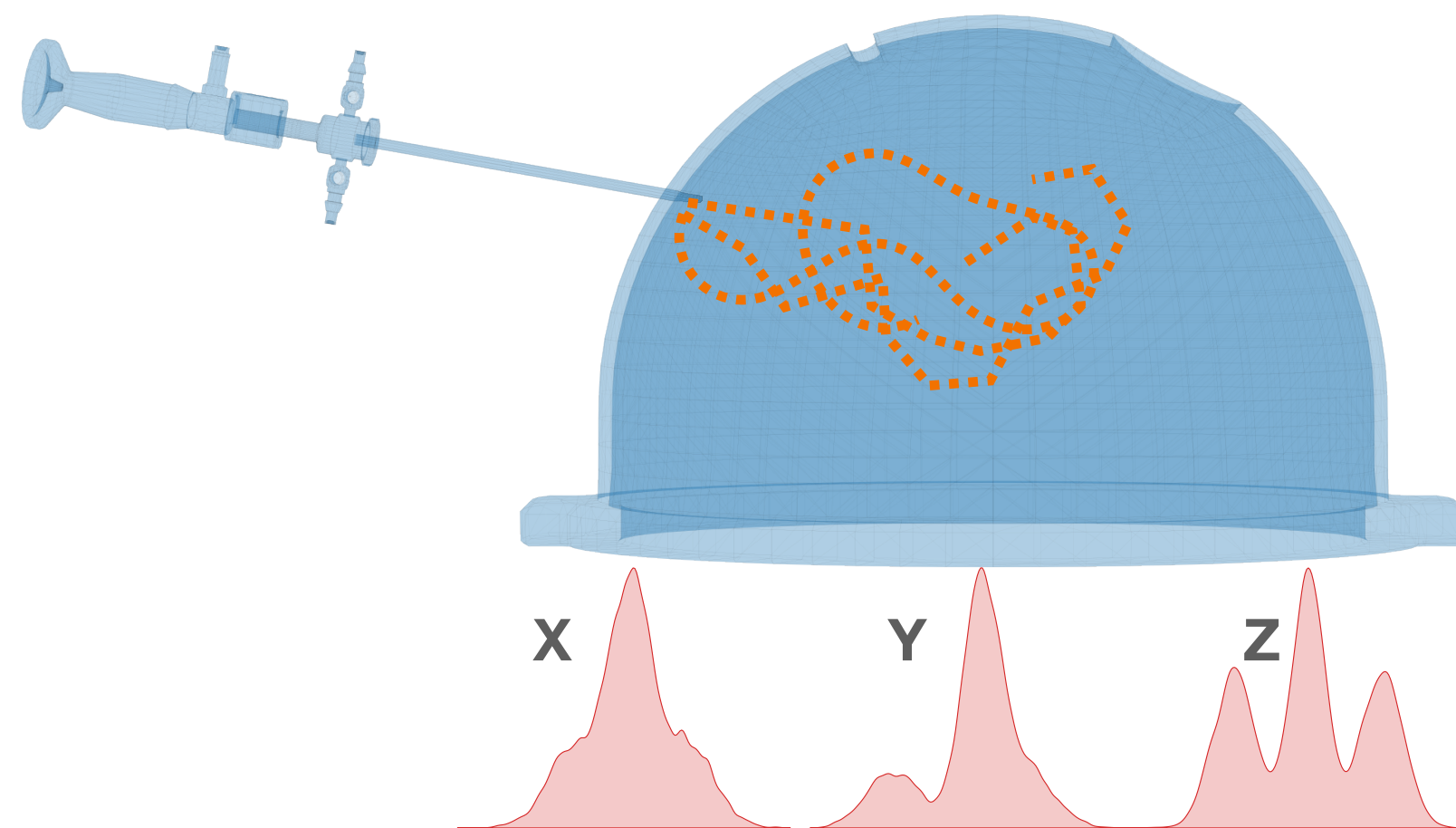


Algorithmic pipeline for surgical assistance

Learning from demonstrations

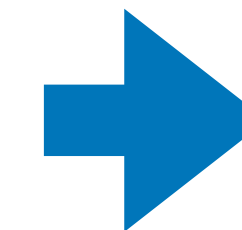
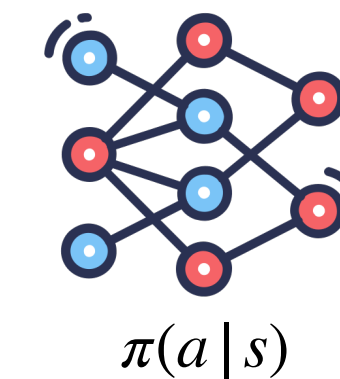
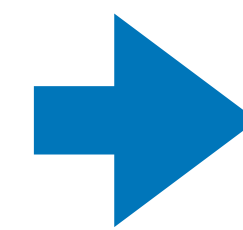
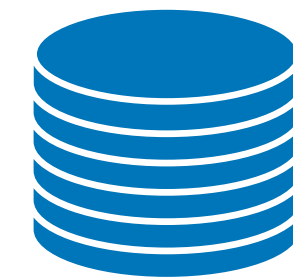
Given a dataset of trajectories

$$\mathcal{D}_E = \{\xi_i\}_{i \leq N} = \{(s_0, a_0, \dots, a_{T_i-1}, s_{T_i})\}_{i \leq N}$$



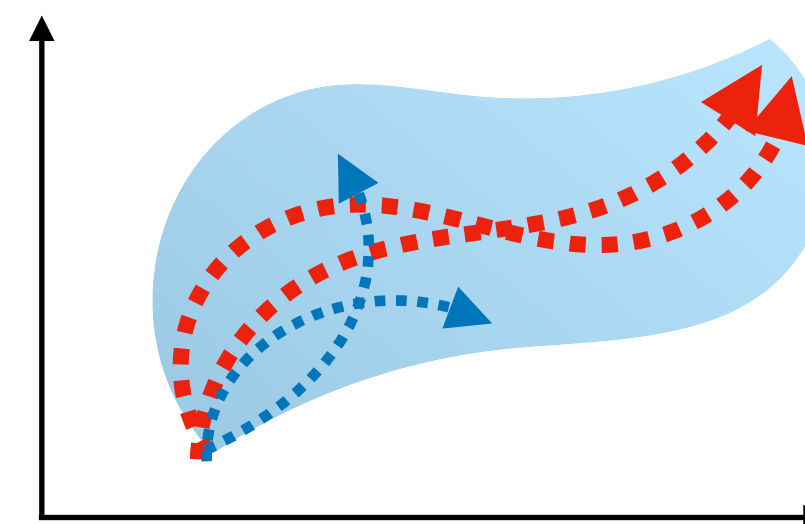
How to learn behaviours from demonstrations?

1. Behavioural cloning



Supervised Learning

2. Distribution matching



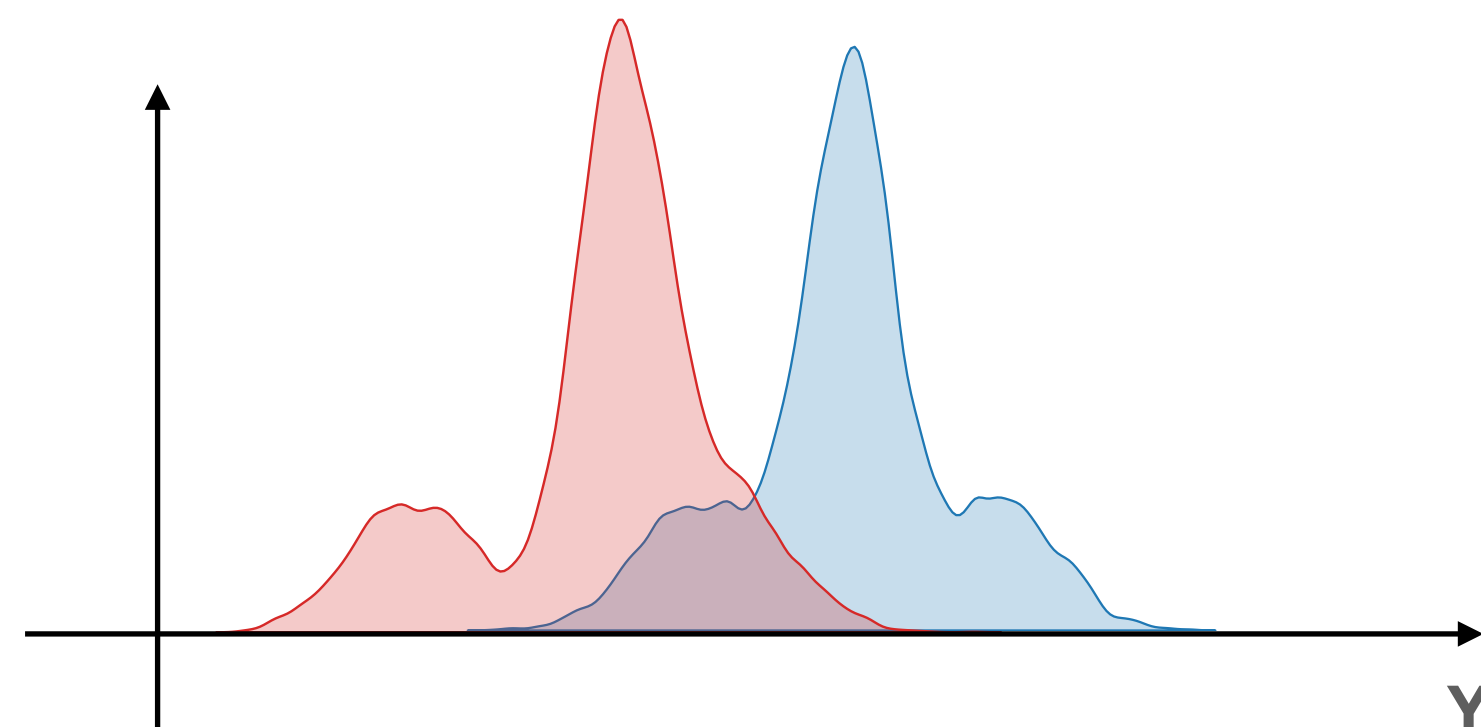
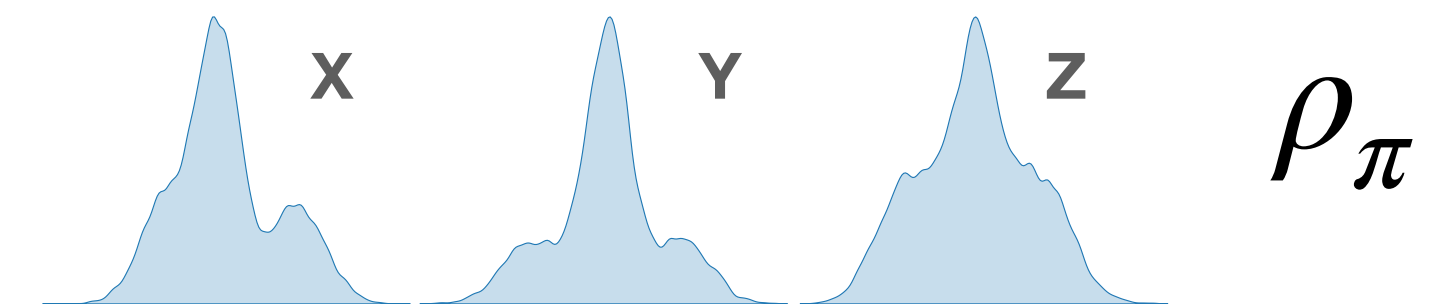
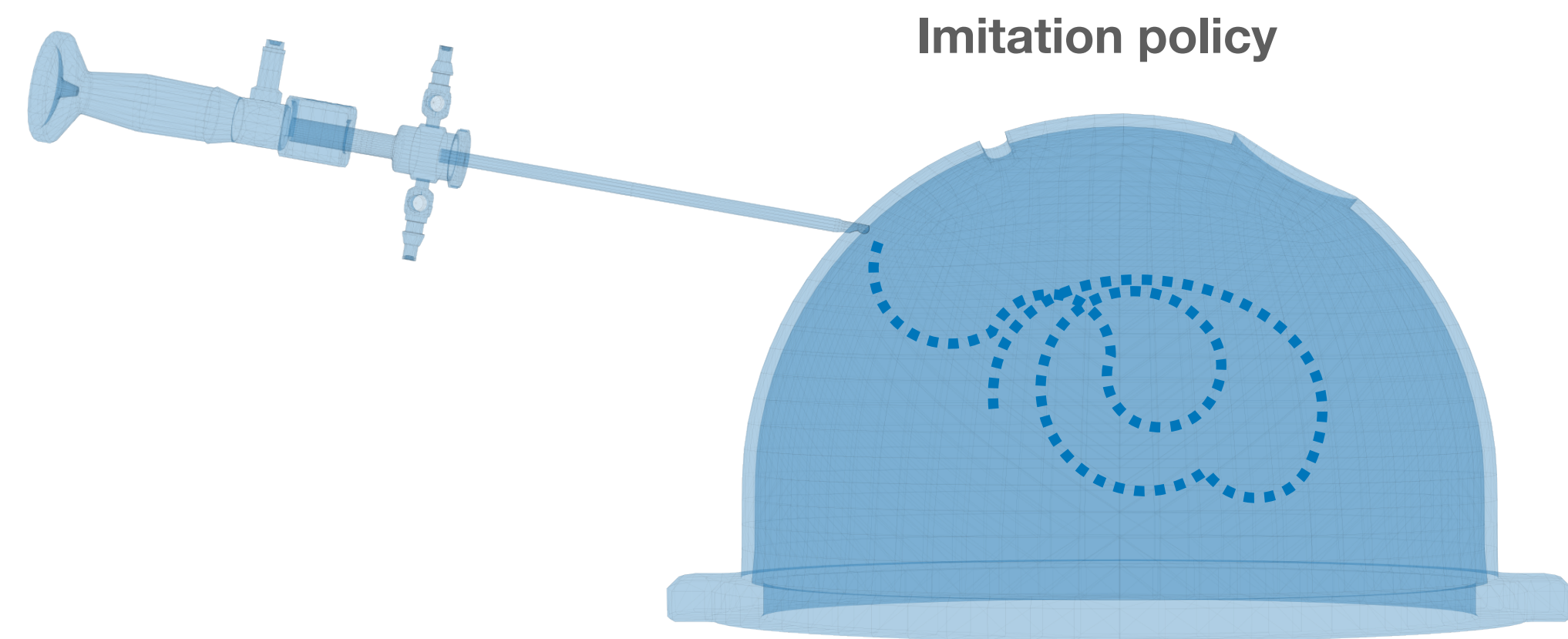
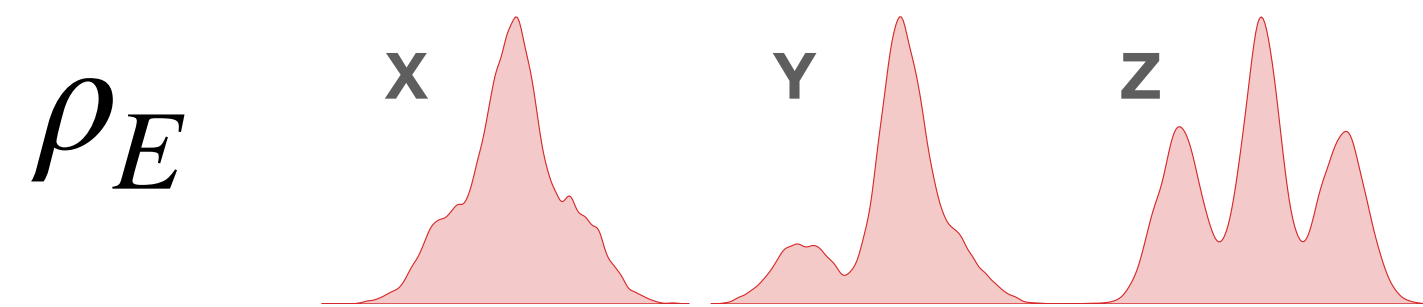
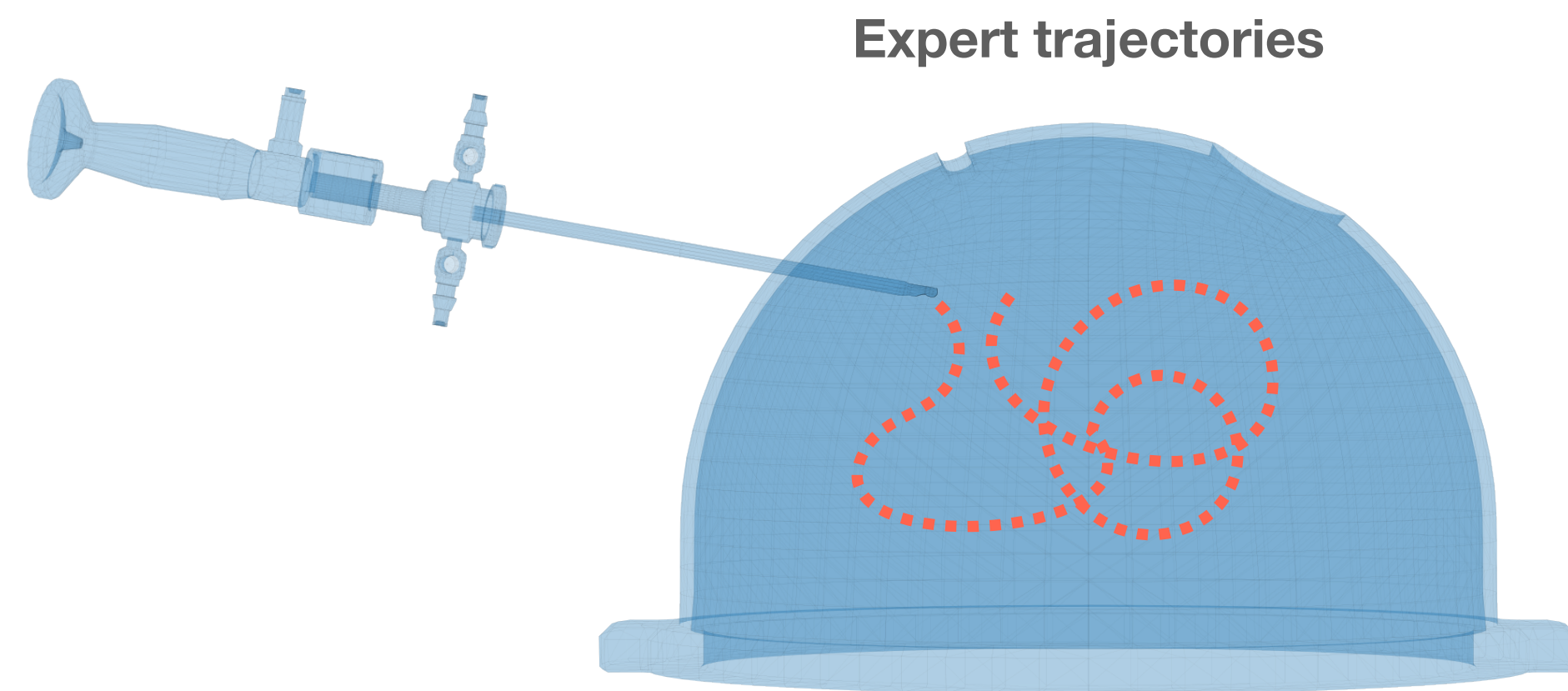
..... Expert trajectories

..... Sampling Policy



Algorithmic pipeline for surgical assistance

Distribution matching as basis for imitation learning



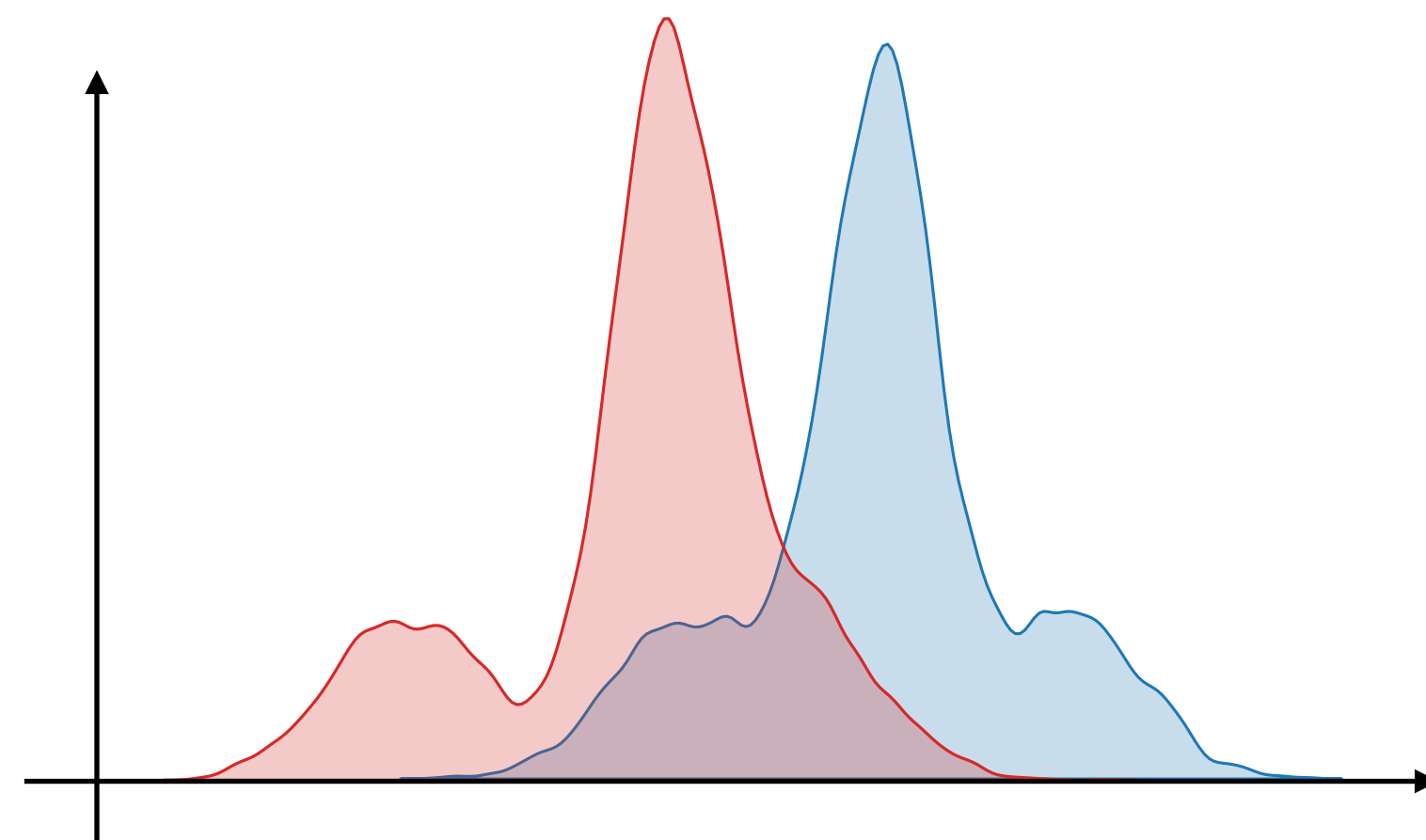
$$\min_{\pi} D(\rho_E, \rho_\pi)$$

Imitation Learning Objective

Algorithmic pipeline for surgical assistance

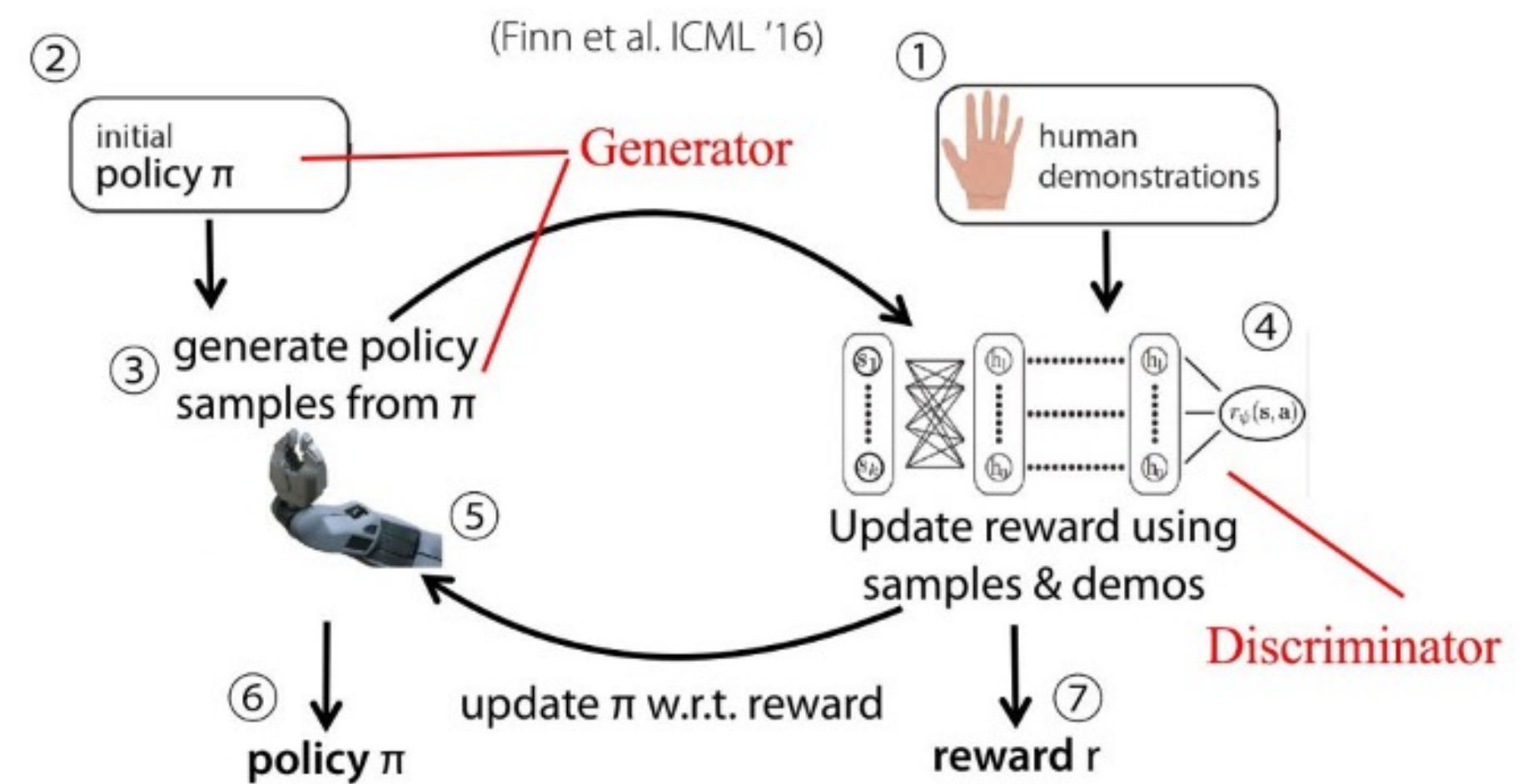
Distribution matching: algorithmic choices

$$\min_{\pi} D(\rho_E, \rho_{\pi})$$



How can we recover a policy and reward?

How to optimize $\min_{\pi} D(\rho_E, \rho_{\pi})$? Min-max problem

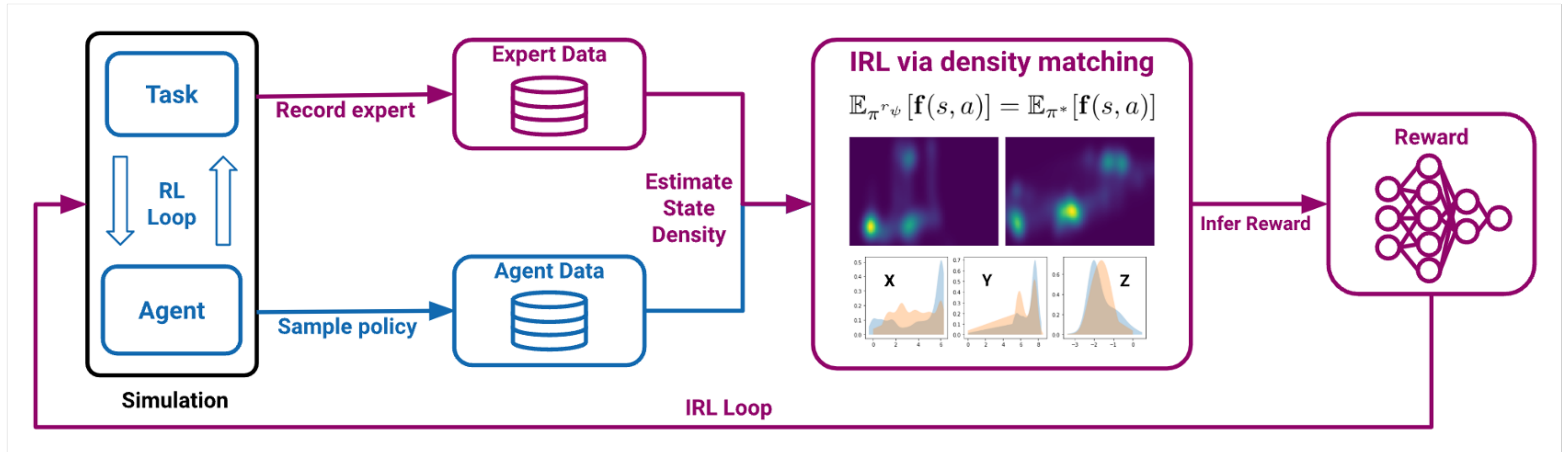


Variational dual of f-divergences

$$D_f(\rho_E, \rho_{\pi}) = \sup_g \mathbb{E}_{\rho_E}[g(x)] - \mathbb{E}_{\rho_{\pi}}[f^*(g(x))]$$

Algorithmic pipeline for surgical assistance

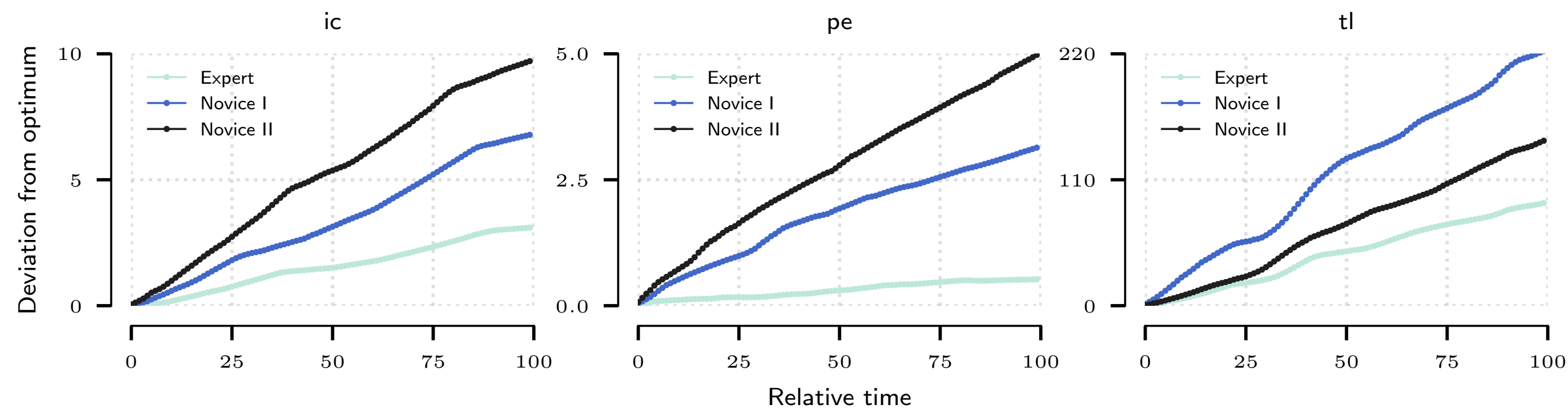
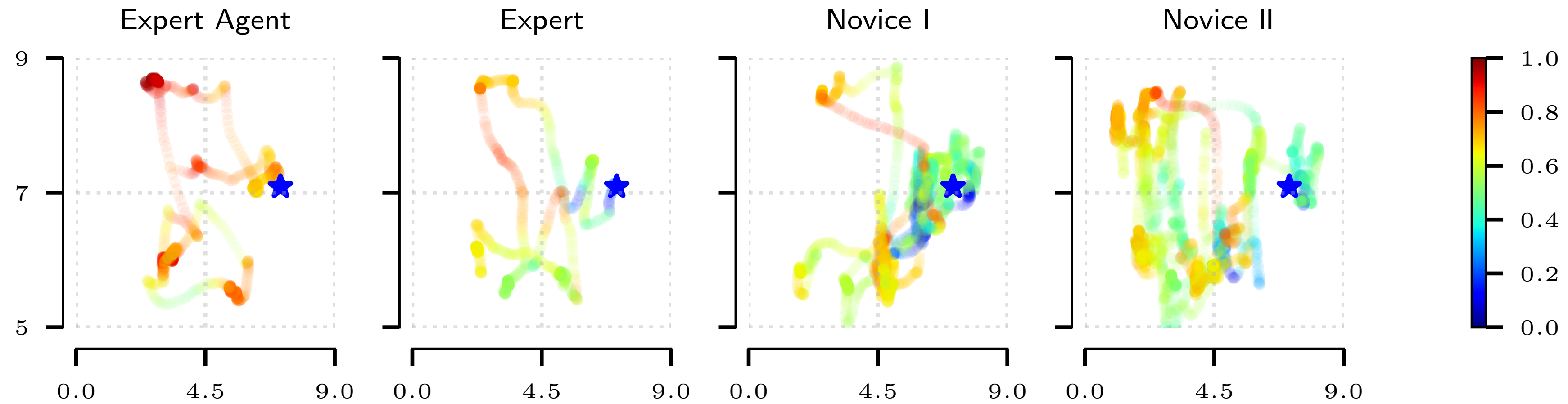
Putting it all together



Can use policy and reward for assistance and evaluation

Algorithmic pipeline for surgical assistance

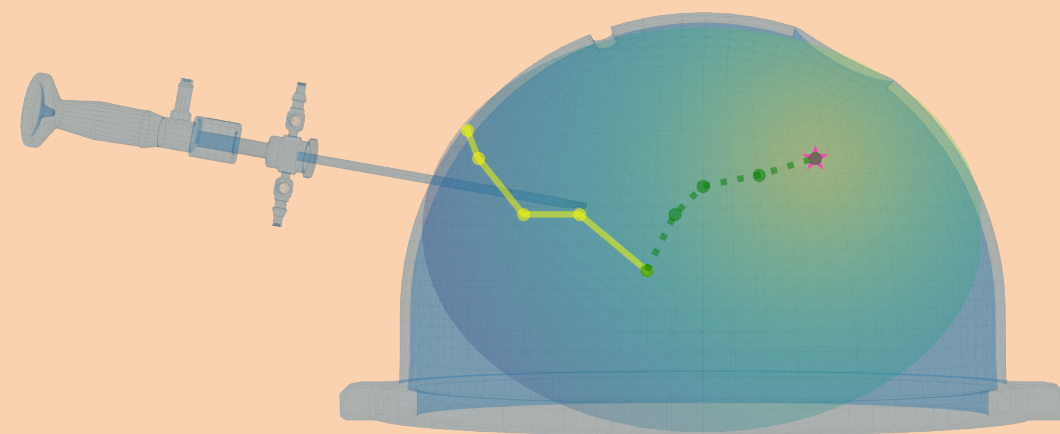
Results: human evaluation example



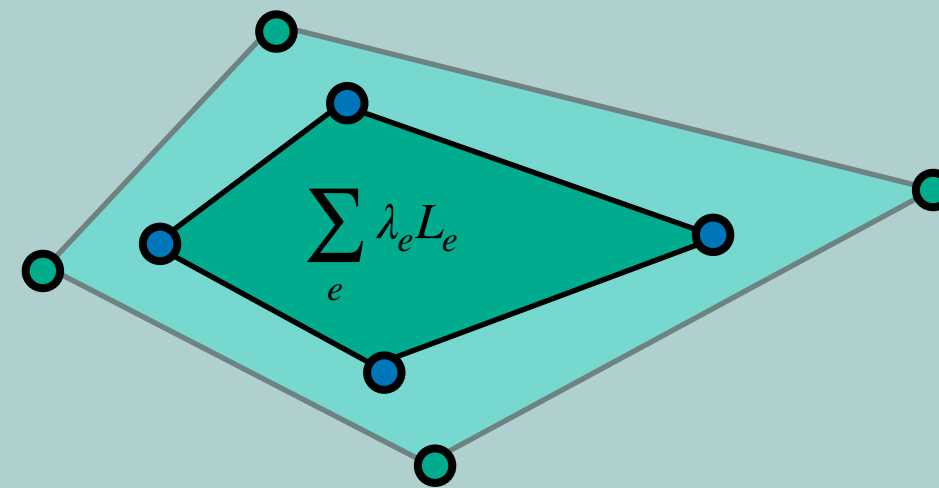
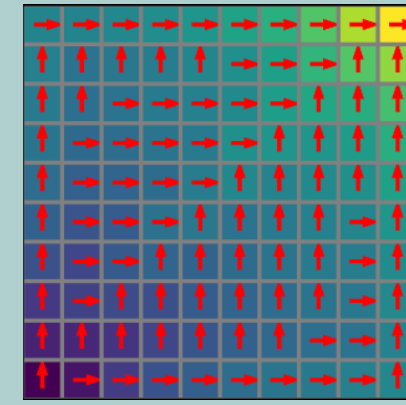
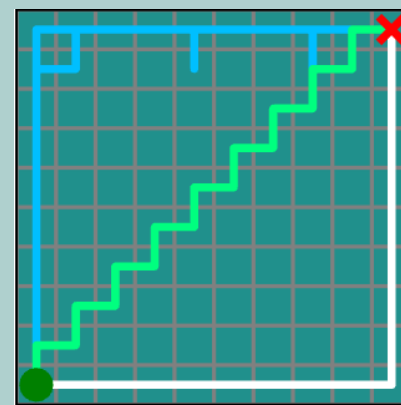
Agent ID	r_ψ	V_ϕ	r_{heur}	Trajectory length
fw	0.993 ± 0.004	0.876 ± 0.076	0.999 ± 0.0002	1600.6 ± 971.0
fm	0.732 ± 0.019	0.728 ± 0.016	0.742 ± 0.016	1821.5 ± 104.8
an	0.720 ± 0.042	0.716 ± 0.040	0.729 ± 0.039	1929.2 ± 277.4
mk	0.617 ± 0.117	0.627 ± 0.121	0.634 ± 0.106	2578.5 ± 745.8
mv	0.518 ± 0.094	0.521 ± 0.094	0.529 ± 0.095	3313.5 ± 665.5
io	0.374 ± 0.271	0.339 ± 0.249	0.384 ± 0.278	4290.8 ± 1981

Overview

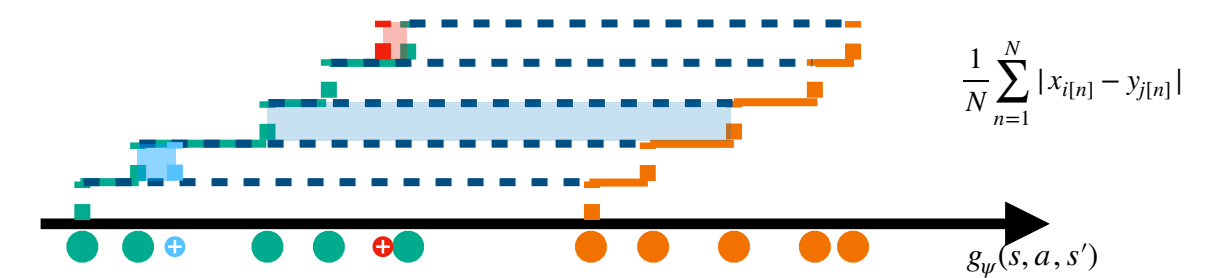
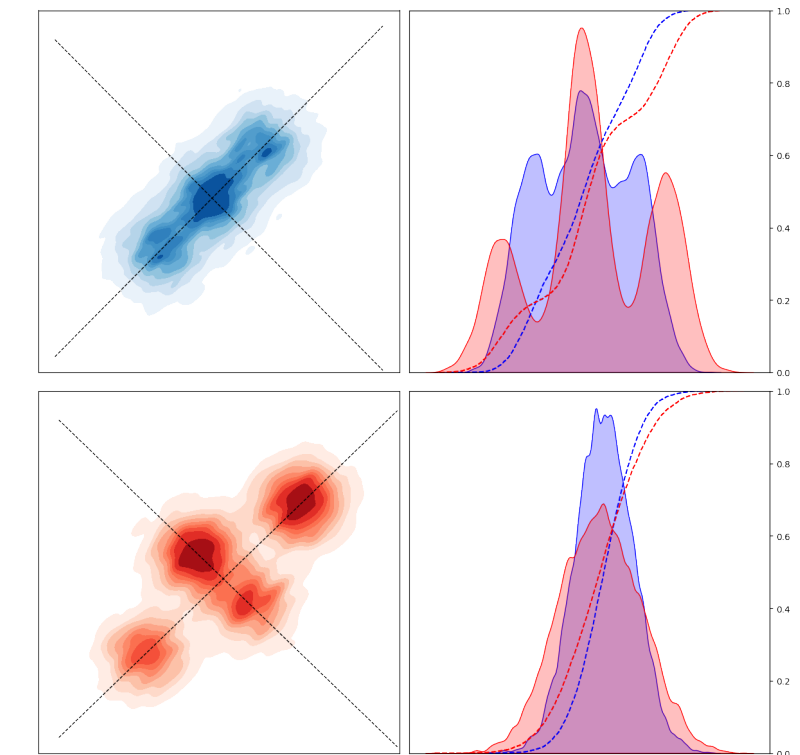
**Application:
Algorithmic RL pipeline in
Surgical Digital Twins**



**Method I:
Addressing Reward
Generalization using Causal
Invariance**



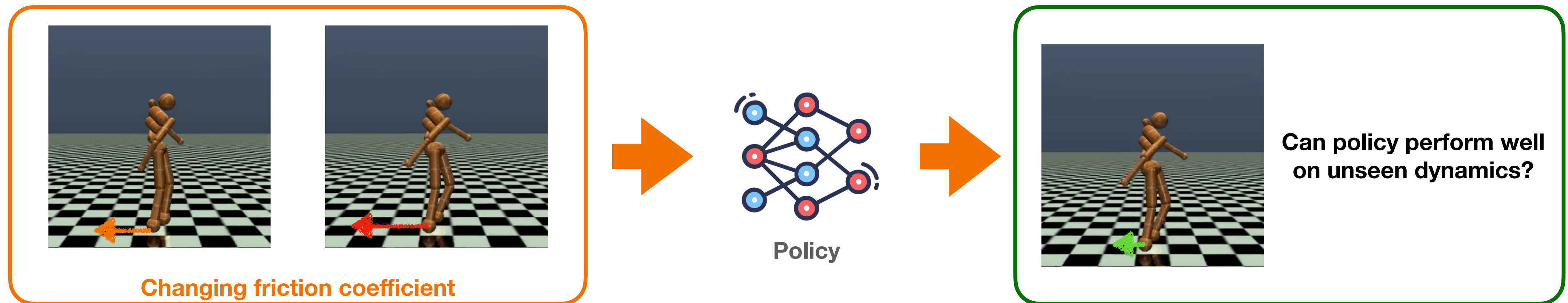
**Method II:
Addressing Data Efficiency in
Imitation Learning using Sliced
Optimal Transport**



Method: addressing generalization

Defining generalization

Generalization in reinforcement learning



Generalization in inverse reinforcement learning

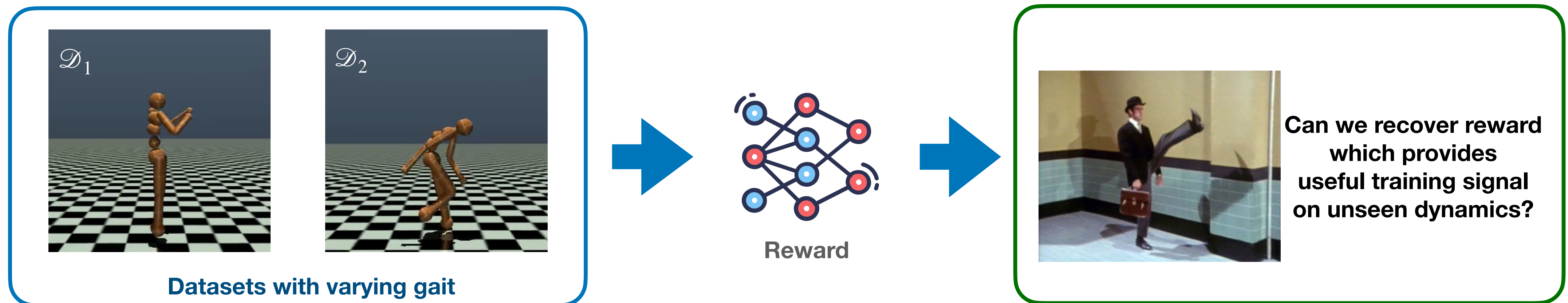
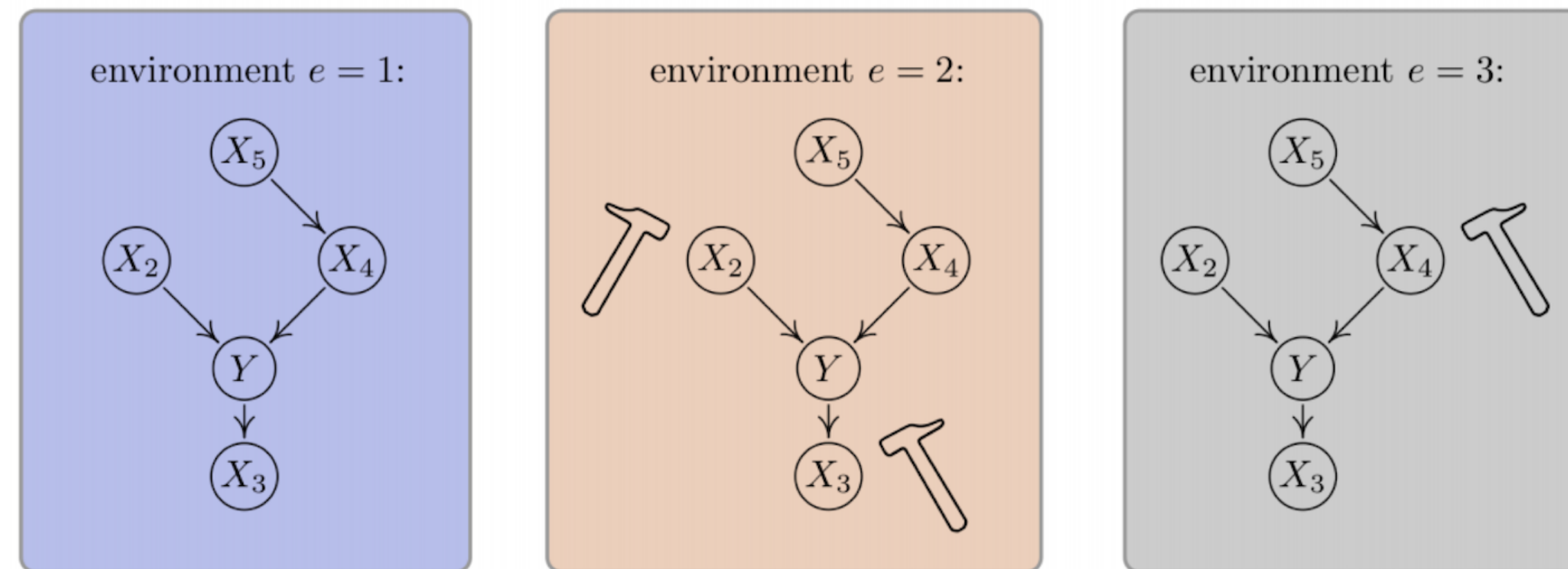


Image: Ministry of Silly Walks, BBC

A little excursion

Invariant causal prediction



Peters et al. 2015

$$\mathcal{D}_e = (\mathbf{X}^e, Y^e) \quad e = \{1 \dots K\}$$

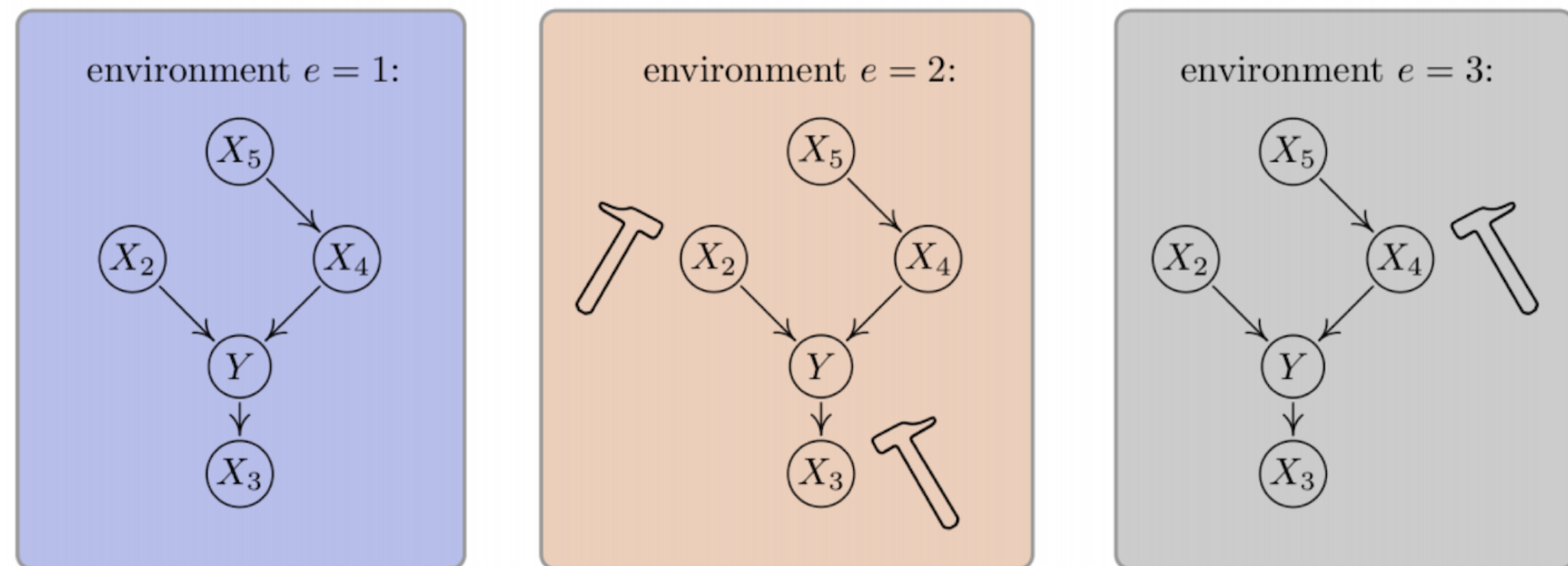
Prediction Task $g : \mathcal{X} \rightarrow \mathcal{Y}$

$$P(Y | X_1, \dots, X_n)$$

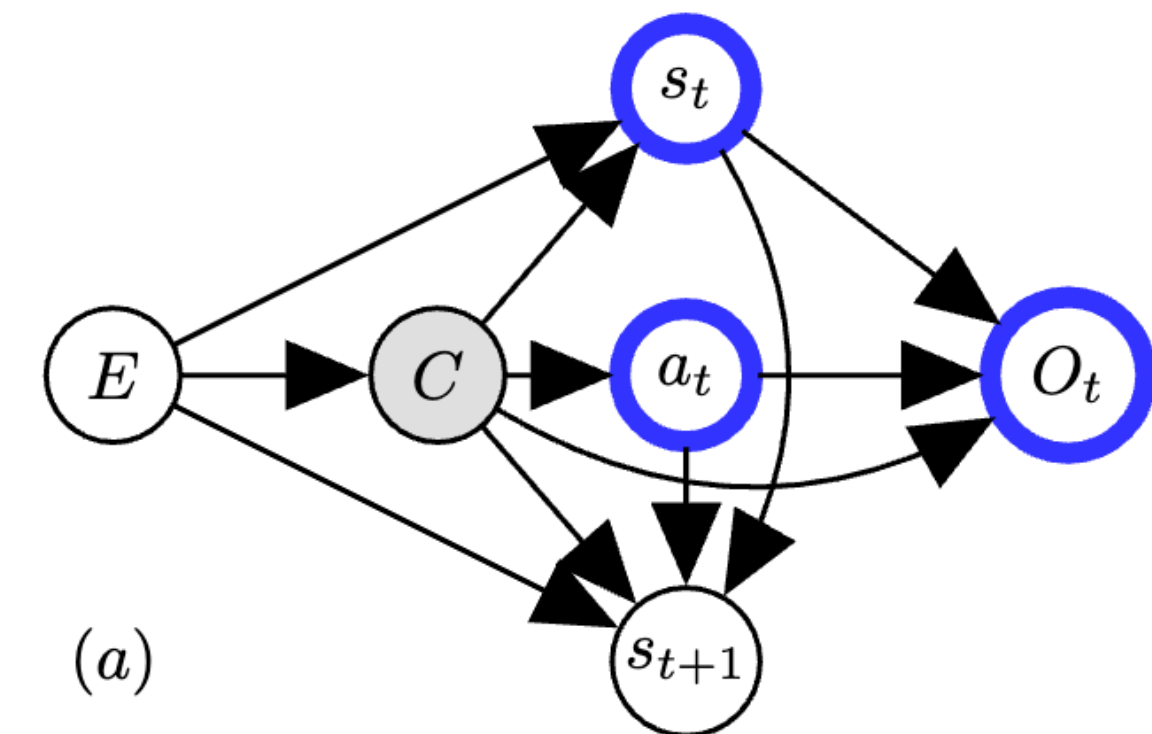
Should be stable across settings!

Reward generalization

Mapping ICP intuition to IRL



Peters et al. 2015



**Inverse reinforcement learning corresponds to learning
the optimality conditional $P(O_t | s_t, a_t)$**

Main intuition: optimality label distribution should be stable across expert demonstrations

Reward generalization

How to adapt principle to IRL setting?

Primal case: Gibbs MLE Problem over trajectories

$$\max_{\varphi, \psi} \mathbb{E}_{\xi \sim \mathcal{D}_E} \left[\log \frac{1}{Z_{\psi, \varphi}} p(\xi | \psi, \varphi) \right]$$

Dual case: total variation distance over transitions

$$\max_{g \in \mathcal{G}} \min_{q \in \mathcal{P}} \mathbb{E}_{(s, a, s') \sim \mathcal{D}_E} [g(s, a, s')] - \mathbb{E}_{(s, a, s') \sim q} [g(s, a, s')]$$

Invariant Risk Minimization

IRM bi-level optimization problem

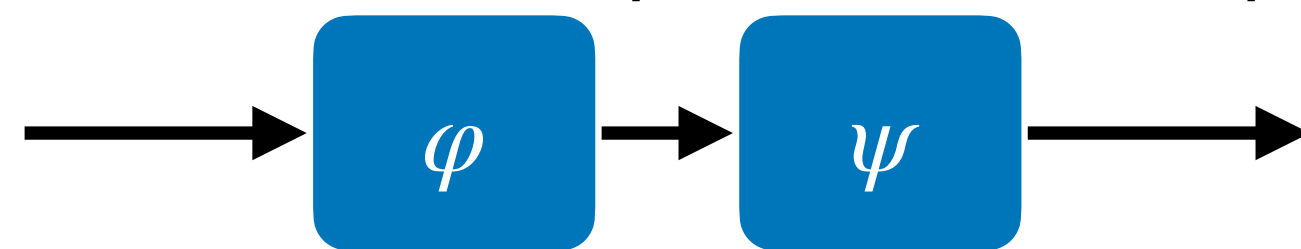
$$\max_{\varphi, \psi} \sum_{e \in \mathcal{E}_{tr}} \sum_{\xi \in \mathcal{D}_e} L_e(\xi, \psi, \varphi)$$

$$\text{s.t. } \psi \in \operatorname{argmax}_{\bar{\psi}} \sum_{\xi \in \mathcal{D}_e} L_e(\xi, \bar{\psi}, \varphi)$$

IRMv1 penalty

$$\mathbb{D}(\psi, \varphi; e) = \|\nabla_{\psi} \mathcal{L}^e(\psi, \varphi)|_{\psi=1.0}\|^2$$

Decompose reward / critic into representation φ and predictor ψ

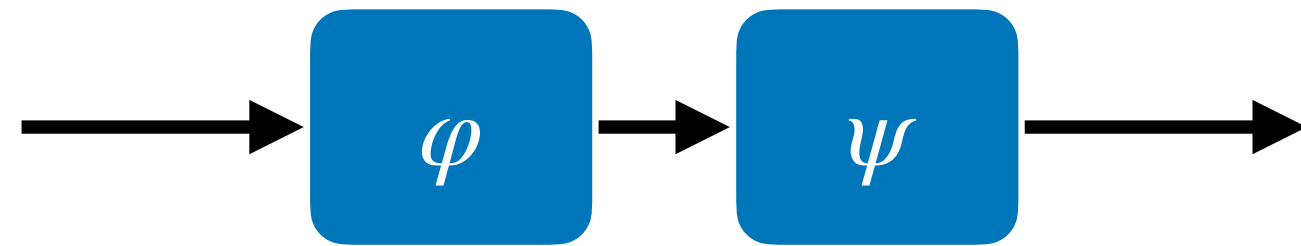


Reward generalization

How to adapt principle to IRL setting?

Invariant Risk Minimization

Decompose reward / critic into representation φ and predictor ψ



IRM bi-level optimization problem

$$\begin{aligned} \max_{\varphi, \psi} \quad & \sum_{e \in \mathcal{E}_{tr}} \sum_{\xi \in \mathcal{D}_e} L_e(\xi, \psi, \varphi) \\ \text{s.t. } \quad & \psi \in \operatorname{argmax}_{\bar{\psi}} \sum_{\xi \in \mathcal{D}_e} L_e(\xi, \bar{\psi}, \varphi) \end{aligned}$$

IRMv1 penalty

$$\mathbb{D}(\psi, \varphi; e) = \|\nabla_{\psi} L^e(\psi, \varphi)|_{\psi=1.0}\|^2$$

IRMv1 penalty applied to IRL as regularizer

$$\text{Primal: } \max_{\varphi, \psi} \sum_{e \in \mathcal{E}_{tr}} \left(\mathbb{E}_{\xi \in \mathcal{D}_e} \left[\log \left(\frac{1}{Z_{\psi, \varphi}} \exp(\psi^T \varphi(\xi)) \right) + \lambda \mathbb{D}(\psi, \varphi, e) \right] \right)$$

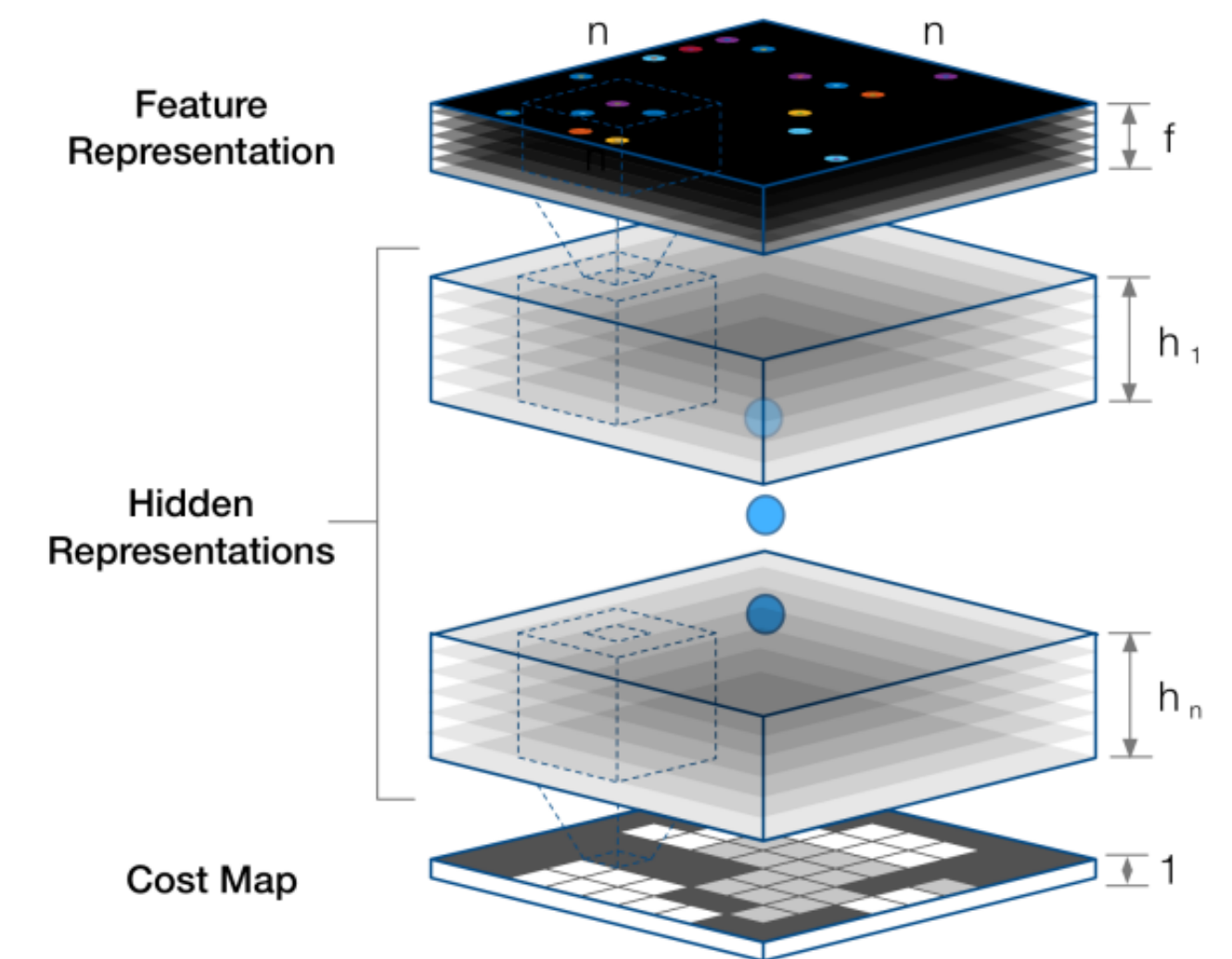
$$\text{Dual: } \max_{\psi, \varphi} \sum_{e \in \mathcal{E}_{tr}} \min_q \left(\mathbb{E}_{(s, a, s') \sim \mathcal{D}_E^e} [g_{\psi, \varphi}(s, a, s')] - \mathbb{E}_{(s, a, s') \sim q} [g_{\psi, \varphi}(s, a, s')] + \lambda \mathbb{D}(\psi, \varphi, e) \right)$$

Reward generalization

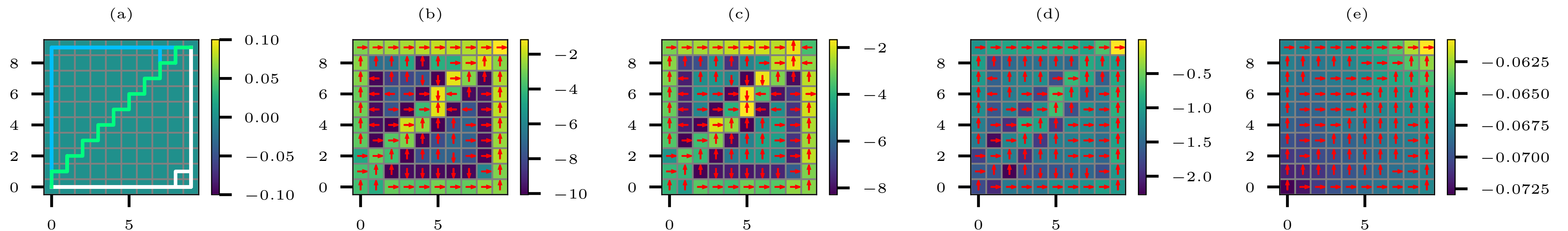
Evaluation: gridworld navigation

Solve primal problem using diverse demonstration data

$$\max_{\varphi, \psi} \sum_{e \in \mathcal{E}_{tr}} \left(\mathbb{E}_{\xi \in \mathcal{D}_e} \left[\log \left(\frac{1}{Z_{\psi, \varphi}} \exp(\psi^T \varphi(\xi)) \right) + \lambda \mathbb{D}(\psi, \varphi, e) \right] \right)$$



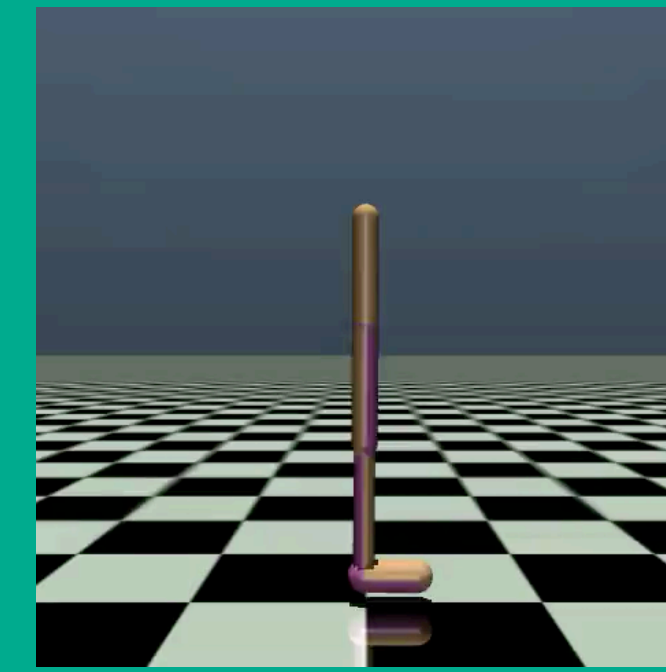
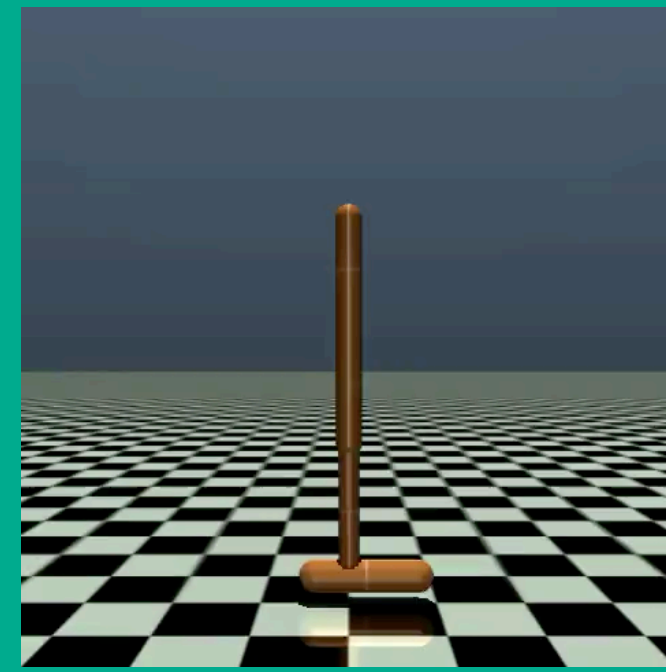
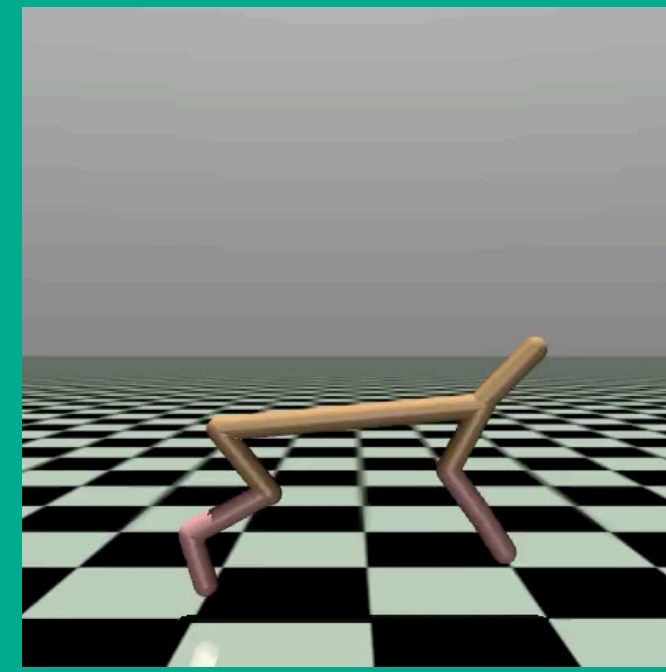
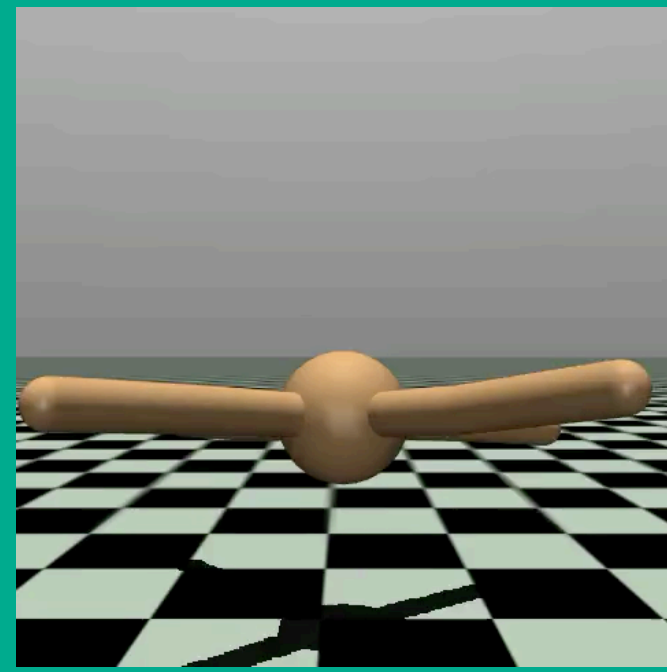
Wulfmeier et al., 2016



Reward generalization

Evaluation: robotic locomotion transfer setting

Mujoco Benchmark

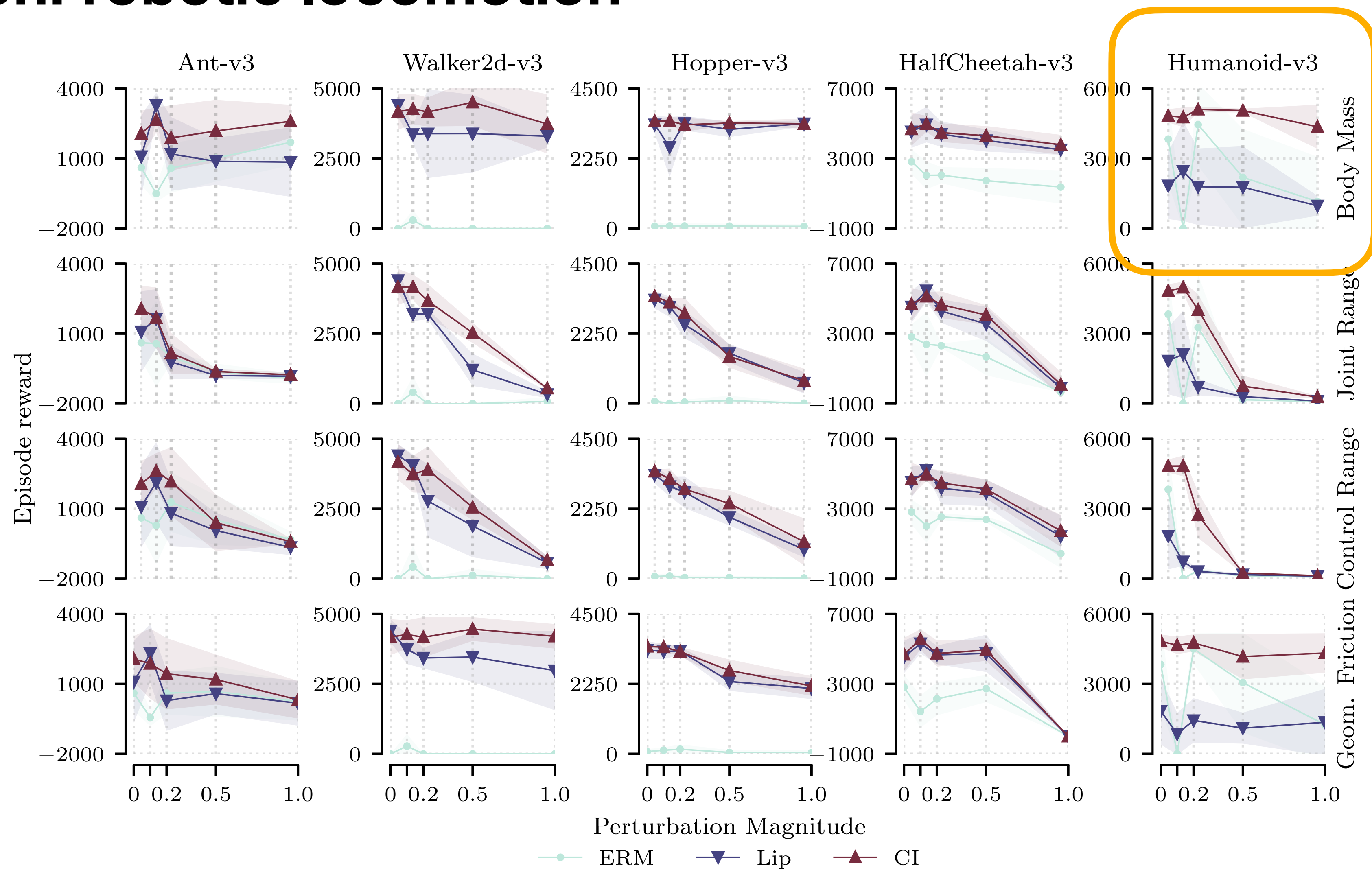


Evaluation protocol:

1. Generate diverse trajectories by applying structured noise
2. Solve dual problem by training adversarial IRL algorithms to recover rewards
3. Use recovered reward to train policies on perturbed dynamics

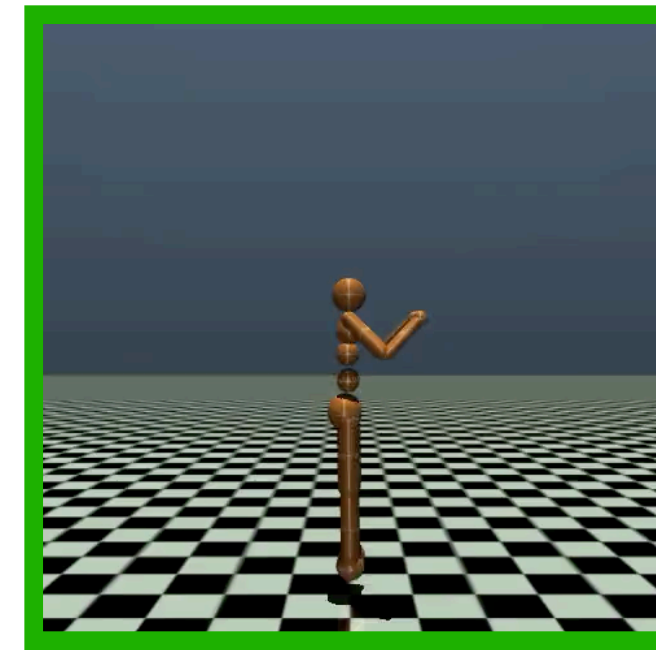
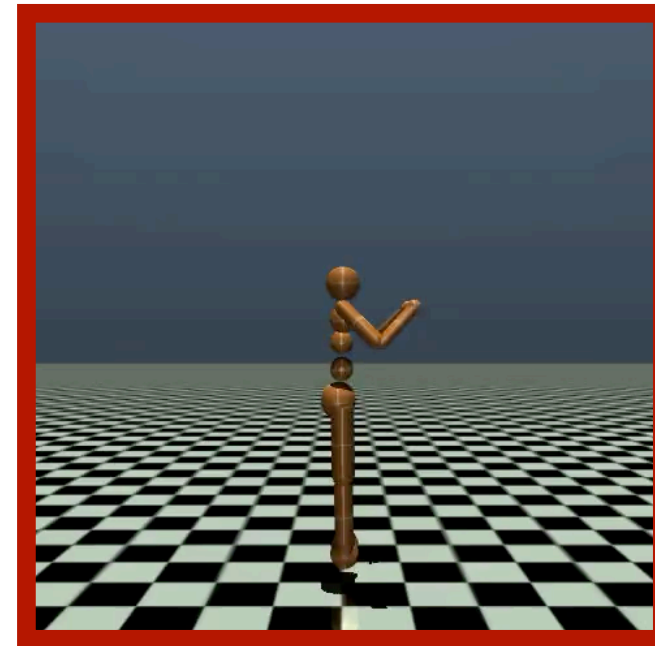
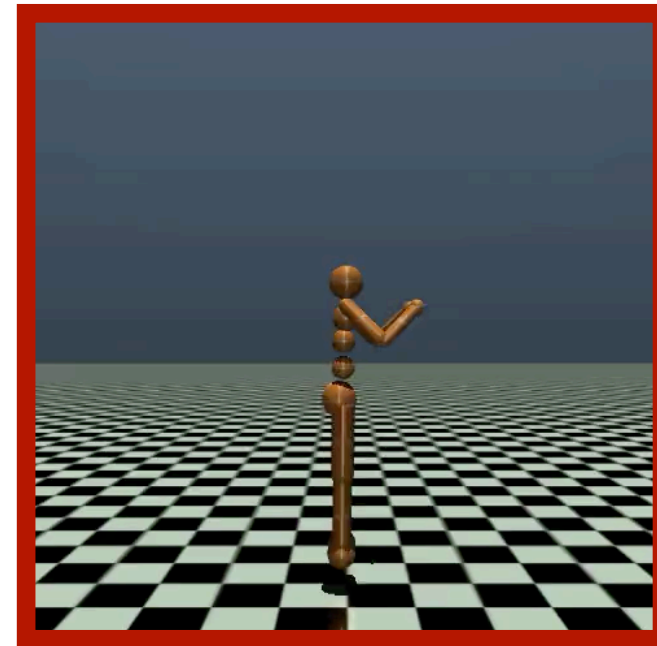
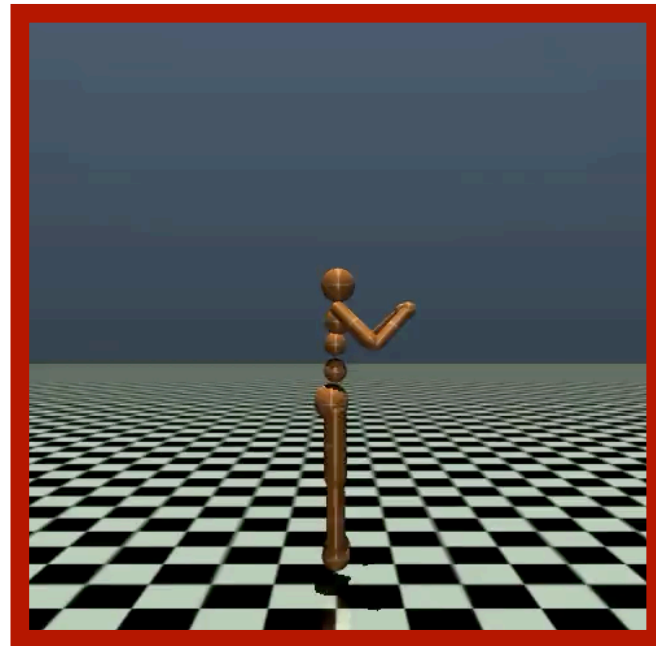
Reward generalization

Evaluation: robotic locomotion

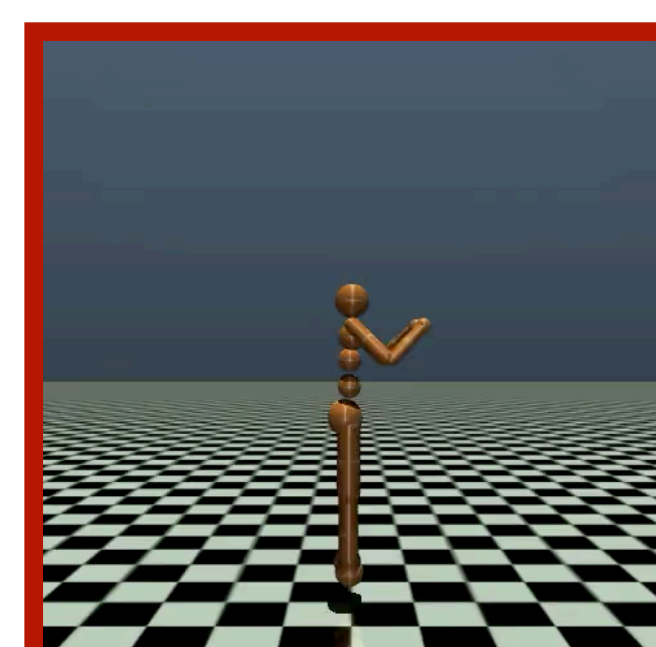
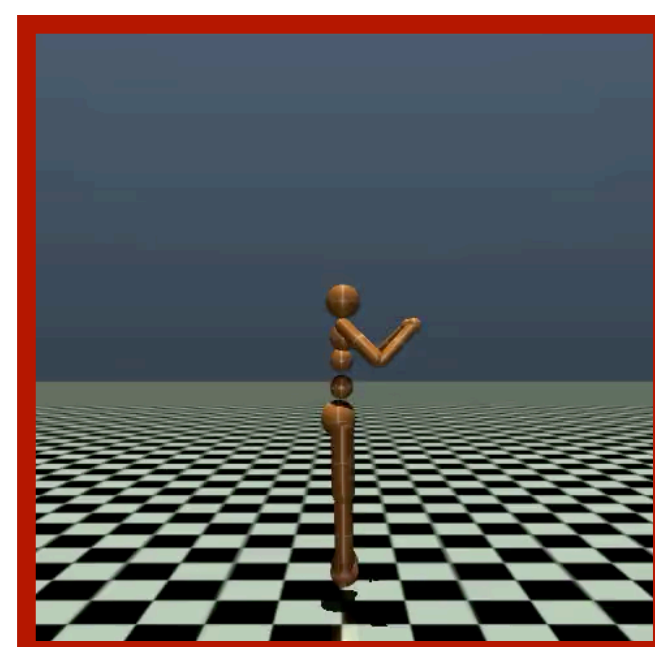
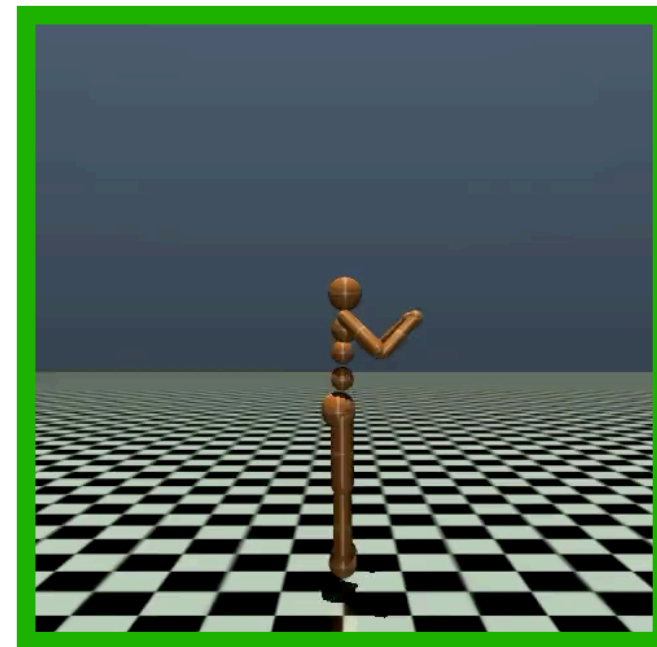
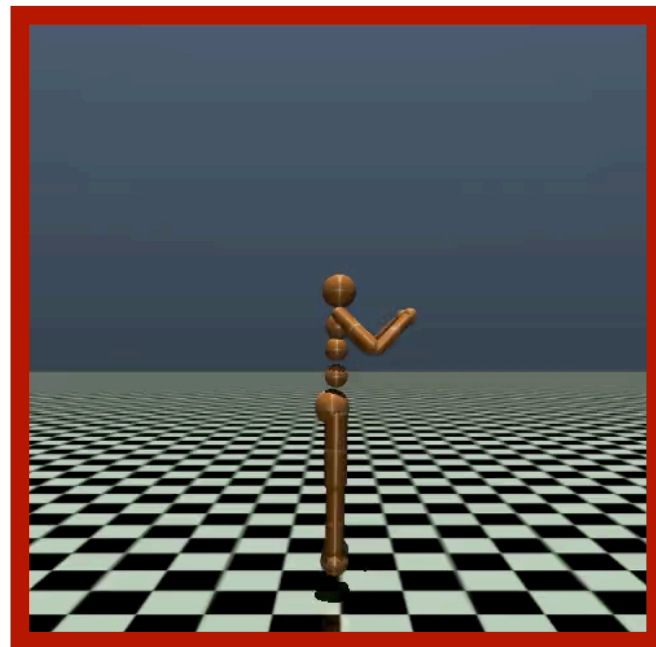


Ministry of silly walks

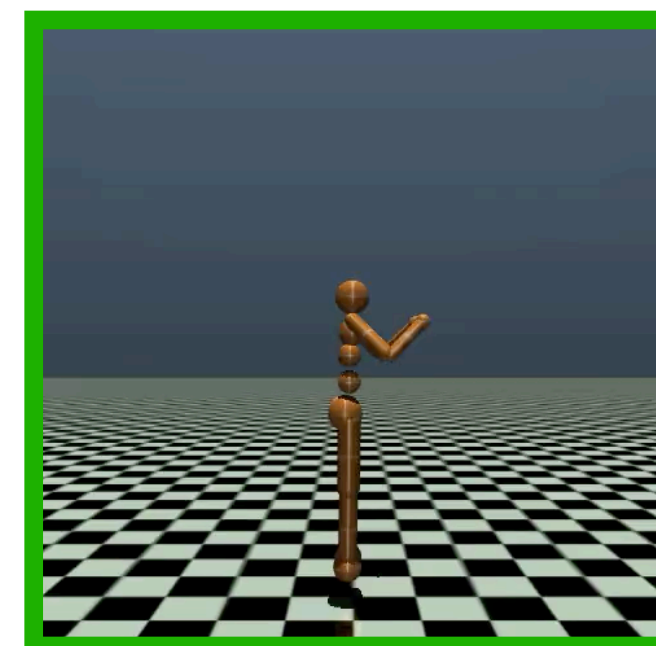
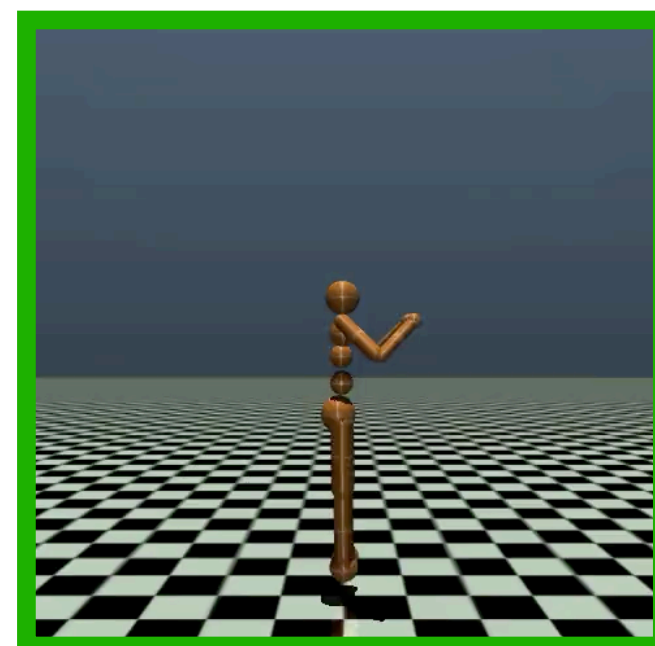
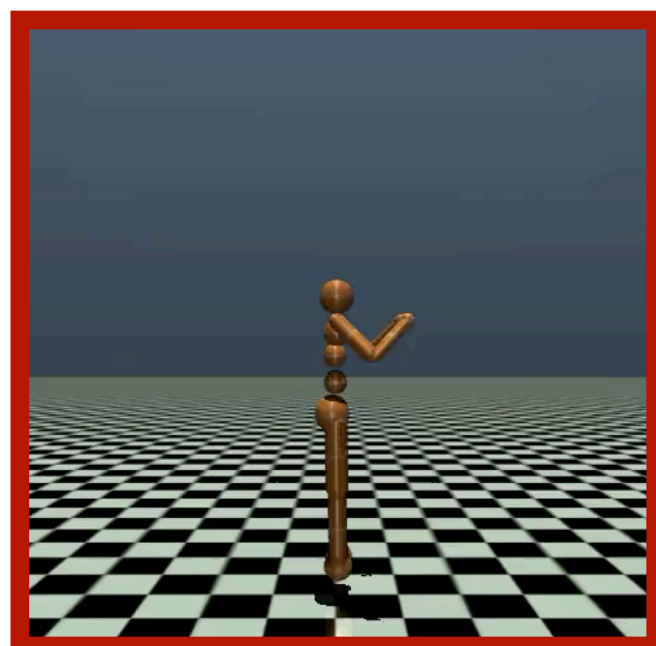
ERM



Lip

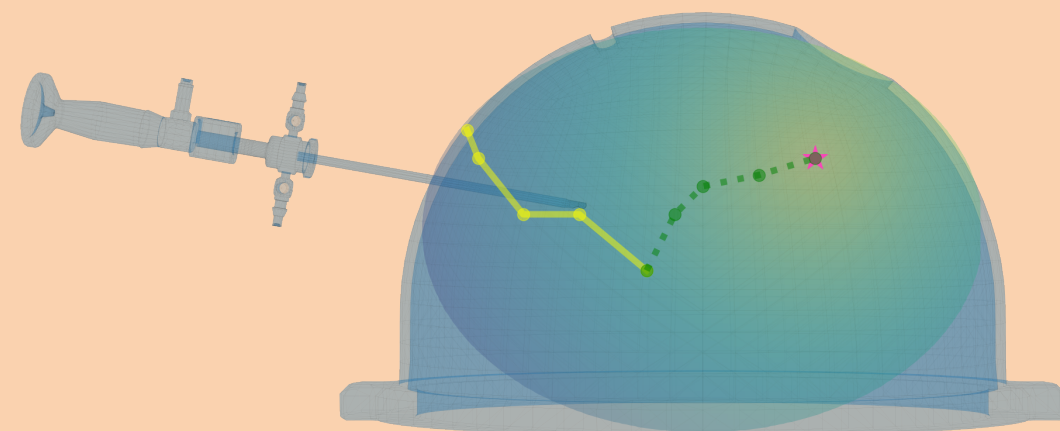


CI

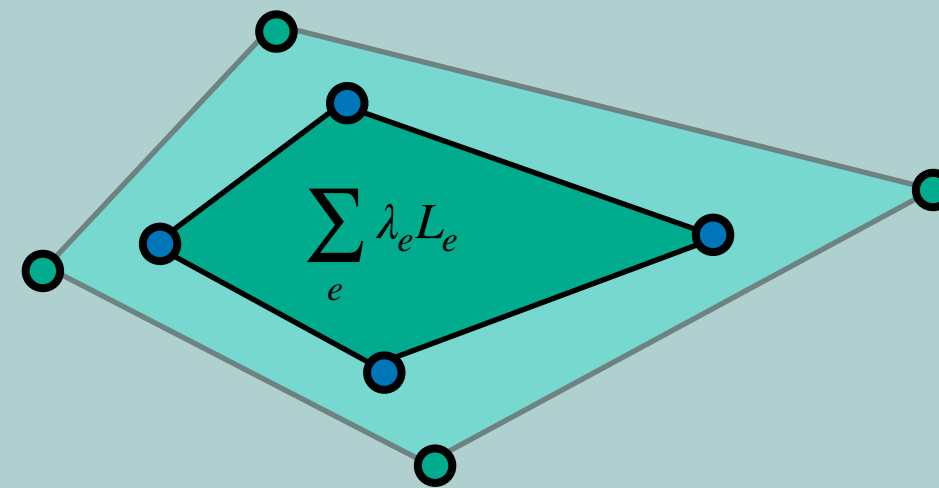
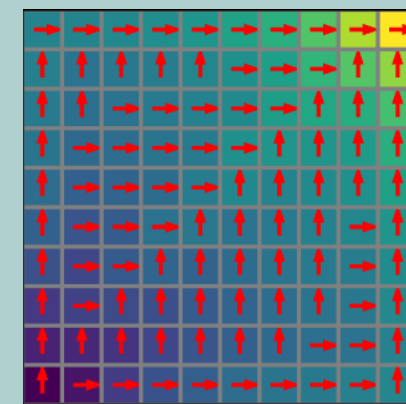
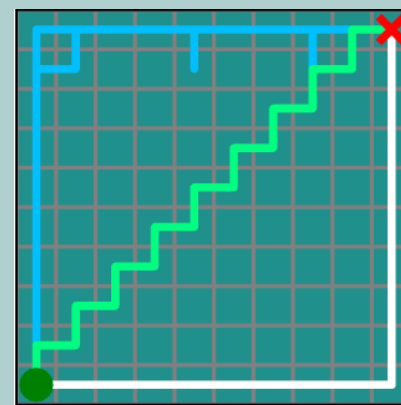


Overview

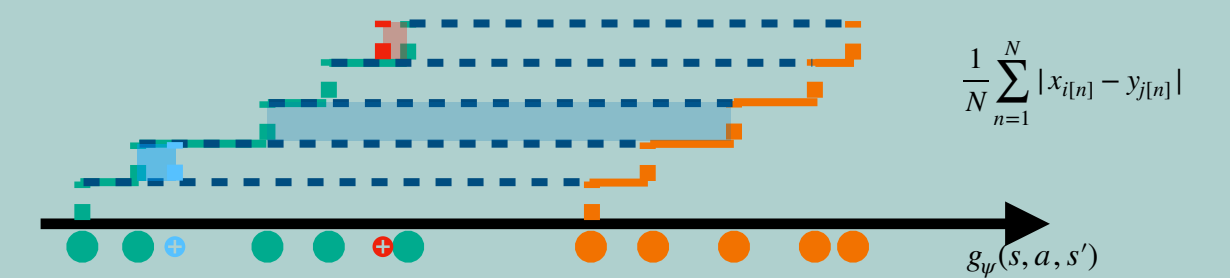
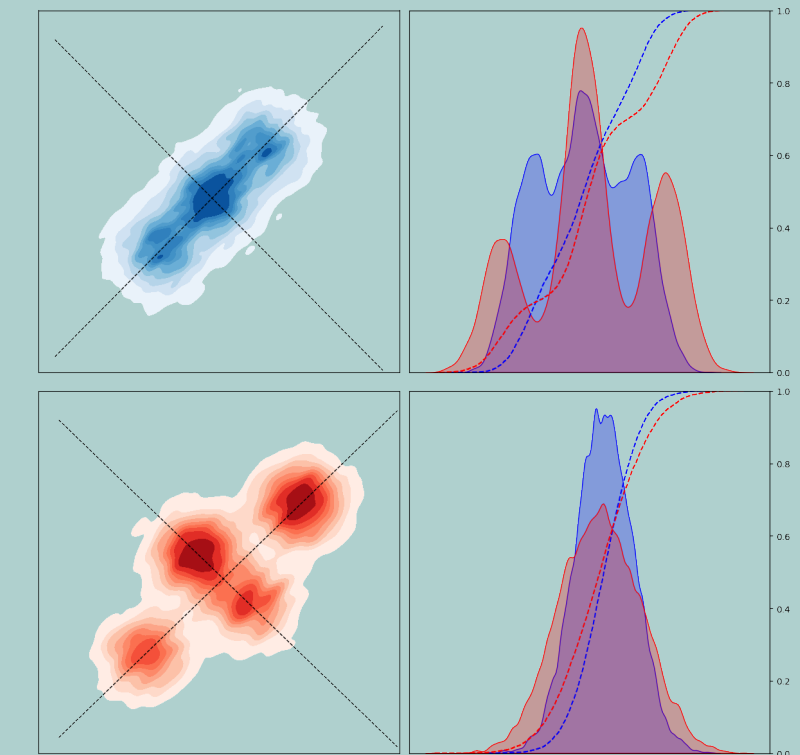
**Application:
Algorithmic RL pipeline in
Surgical Digital Twins**



**Method I:
Addressing Reward
Generalization using Causal
Invariance**



**Method II:
Addressing Data Efficiency in
Imitation Learning using Sliced
Optimal Transport**



Addressing Data Efficiency

Back to distribution matching

So far, focus on distribution matching methods based on f-divergences

$$D_f(P, Q) := \begin{cases} \mathbb{E}_Q \left[f \left(\frac{dP}{dQ} \right) \right] & P \ll Q \\ +\infty & \text{otherwise} \end{cases}$$

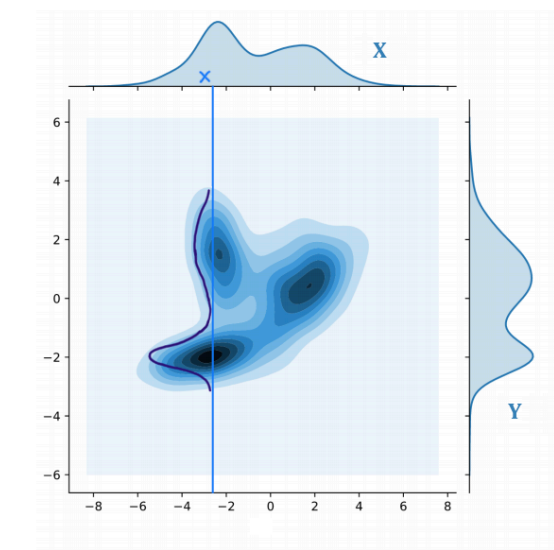
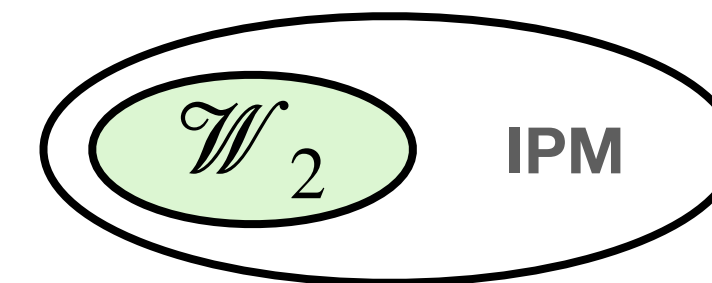
Two issues

- Ill-defined for disjoint support
- Requires brittle optimization procedures

Alternative: Integral Probability Metrics

$$v_{\mathcal{F}}(P, Q) := \sup_{f \in \mathcal{F}} |\mathbb{E}_P[f] - \mathbb{E}_Q[f]|$$

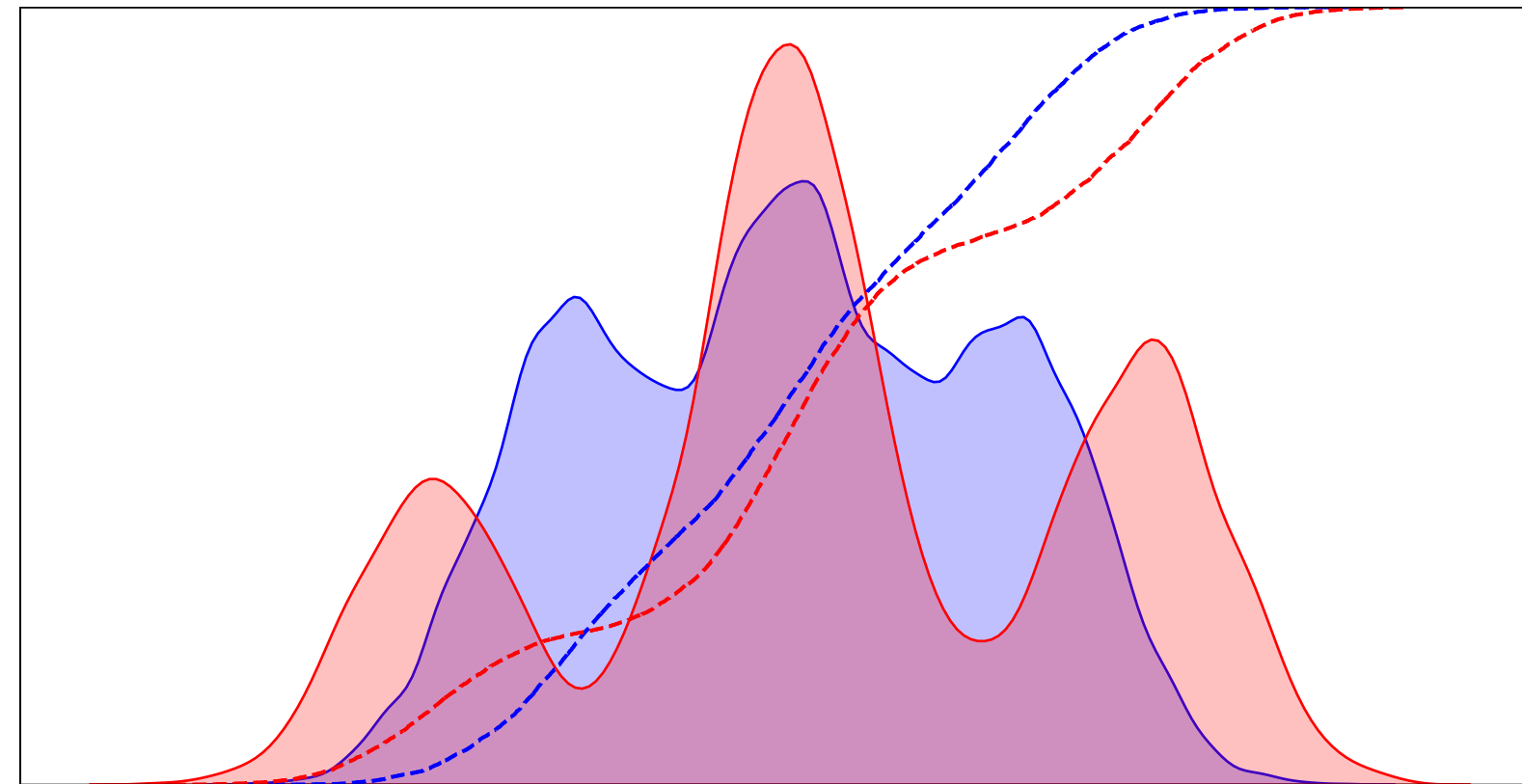
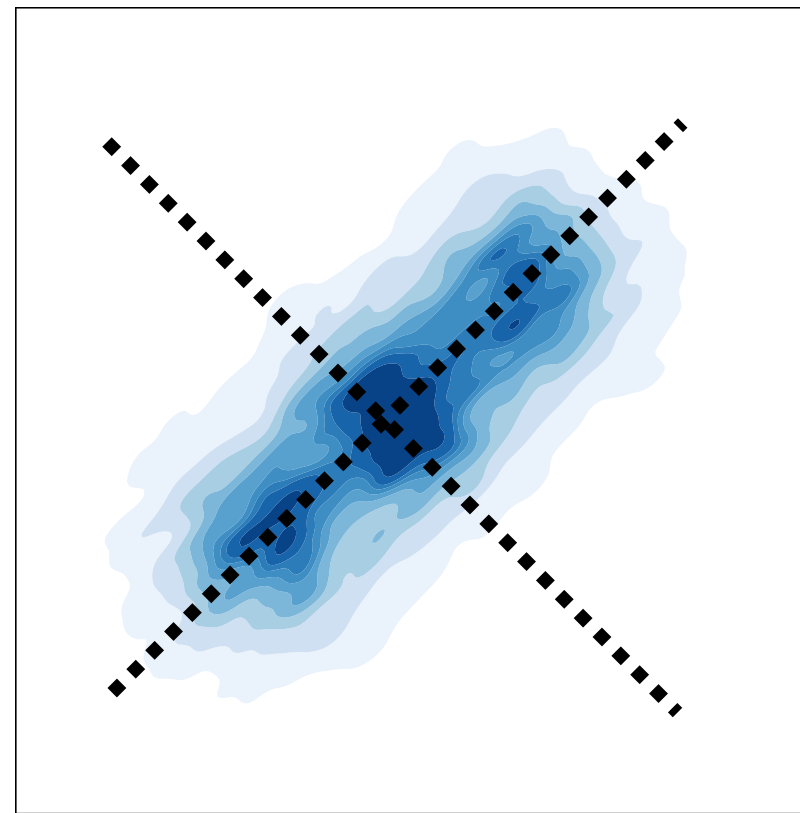
Depending on \mathcal{F} , different metrics:
MMD, Kantorovich, Dudley



Leverage advances in computational optimal transport (OT)

Exploring computational OT

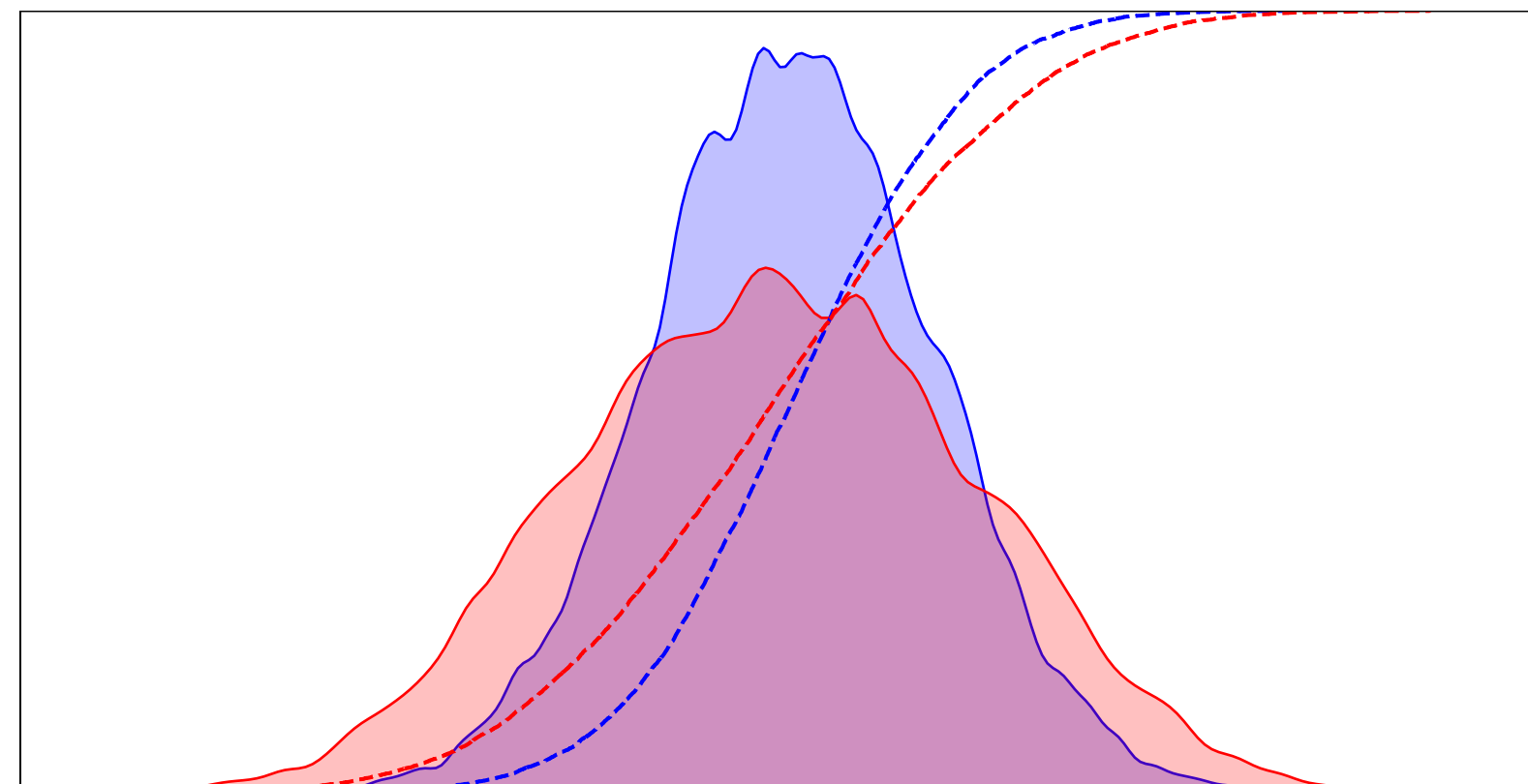
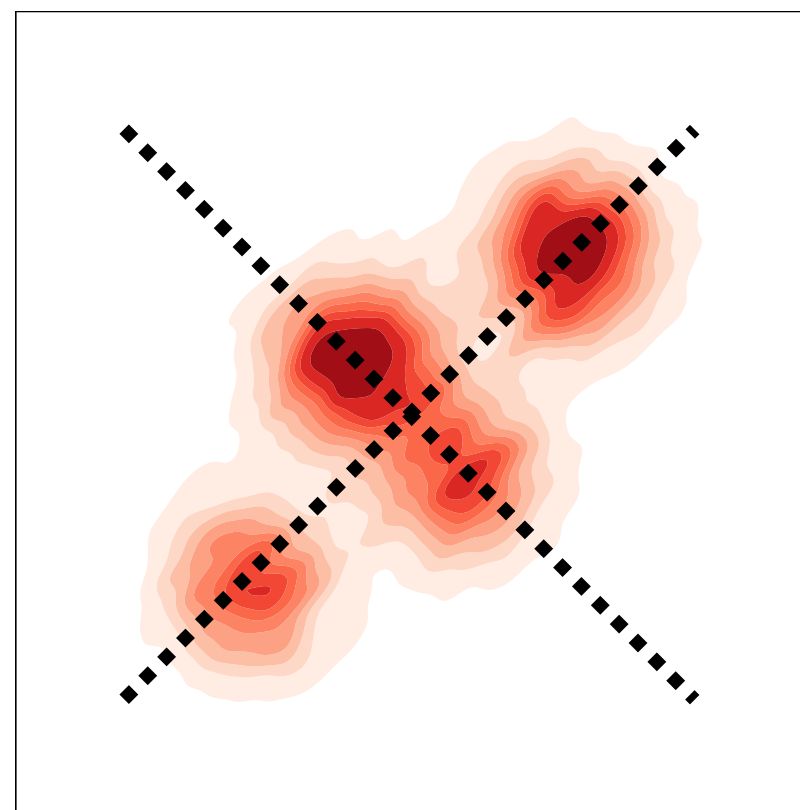
Sliced optimal transport



Leverage closed form of 1-D Wasserstein Distance computation

$$W_2(\rho, \rho') = \int_0^1 |F_\rho^{-1}(\mu) - F_{\rho'}^{-1}(\mu)|^2 d\mu$$

Efficient sorting-based computation



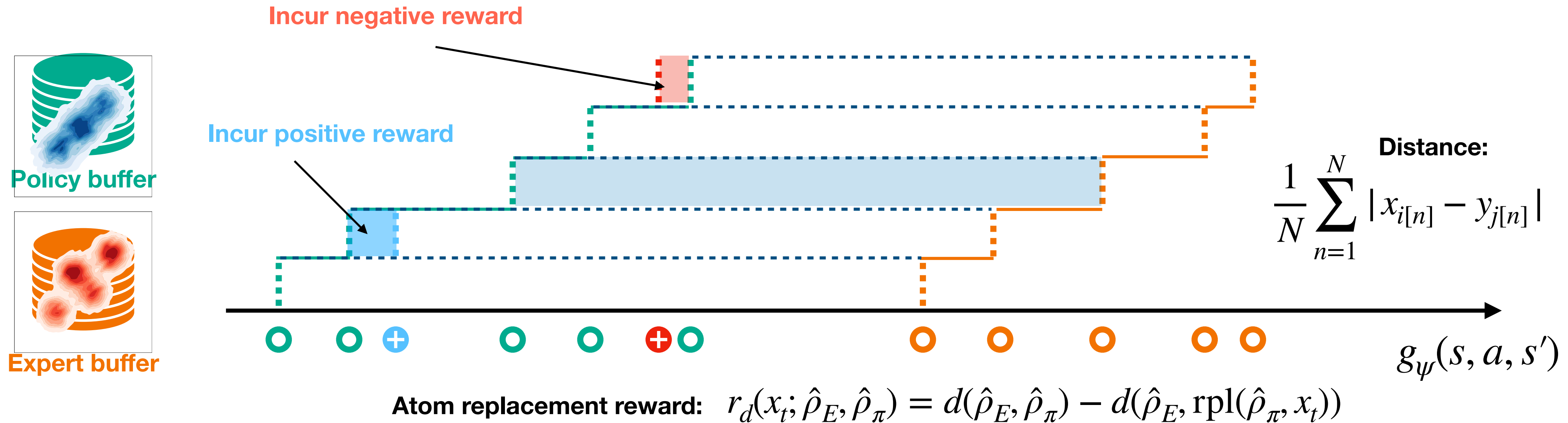
Sliced Wasserstein Distance

$$SW_2(\rho, \rho') = \int_{\mathcal{S}^{d-1}} W_2(\psi_{\#}\rho, \psi_{\#}\rho') d\psi$$

Imitation Learning: No access to dynamics

Sliced Wasserstein Distances

Proxy reward structure



Procedure

1. Project buffer samples
2. Sample policy and project atom
3. Find closest atom in ranking and measure replacement distance
4. Use measure replacement distance as reward

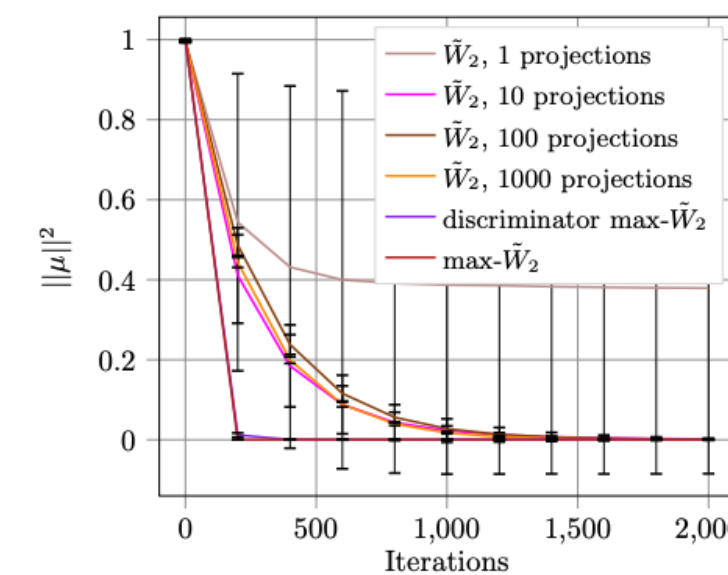


Policy optimizing reward reduces the 1-D Wasserstein distance

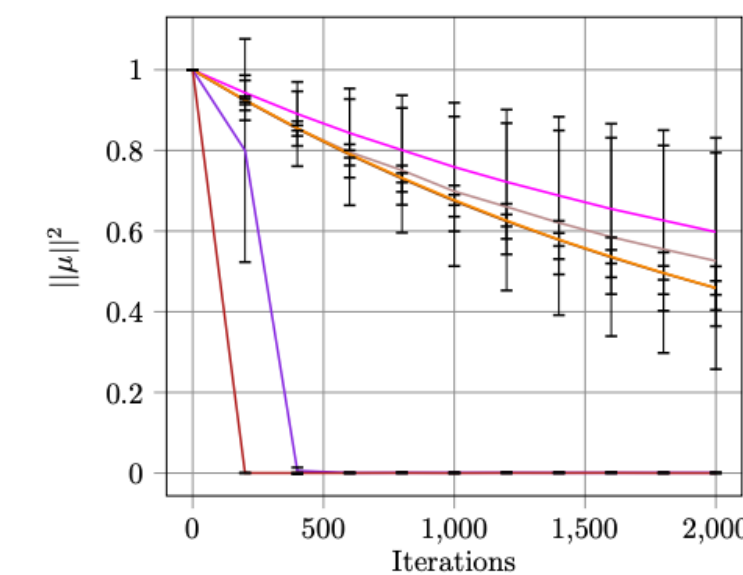
Choice of projections

Generalized Sliced Wasserstein (GSW) Distances

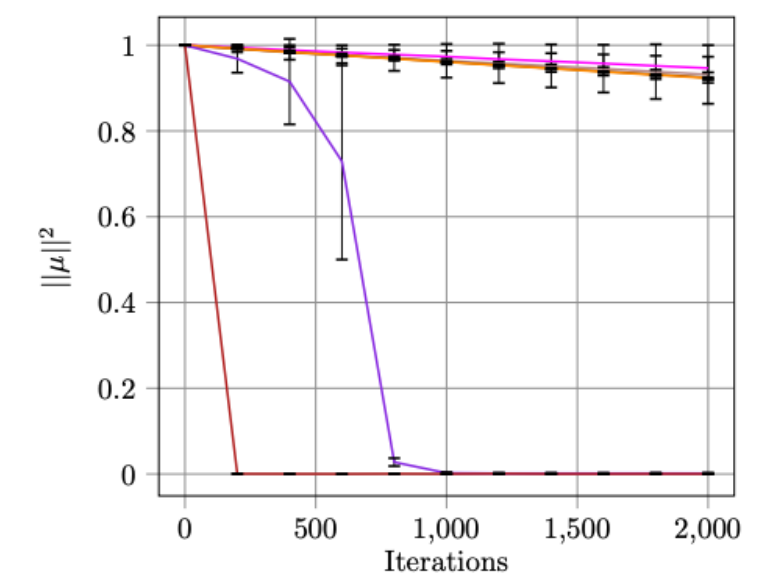
Linear projections scale poorly with state space dimensionality



(a) $d = 10$



(b) $d = 100$

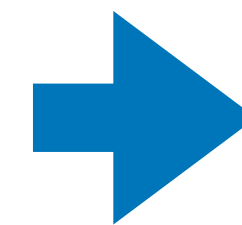


(c) $d = 1000$

Deshpande et al., 2019

Choose most informative (nonlinear) projection

$$SW_2(\rho, \rho') = \int_{\mathcal{S}^{d-1}} W_2(\psi_{\#}\rho, \psi_{\#}\rho') d\psi$$



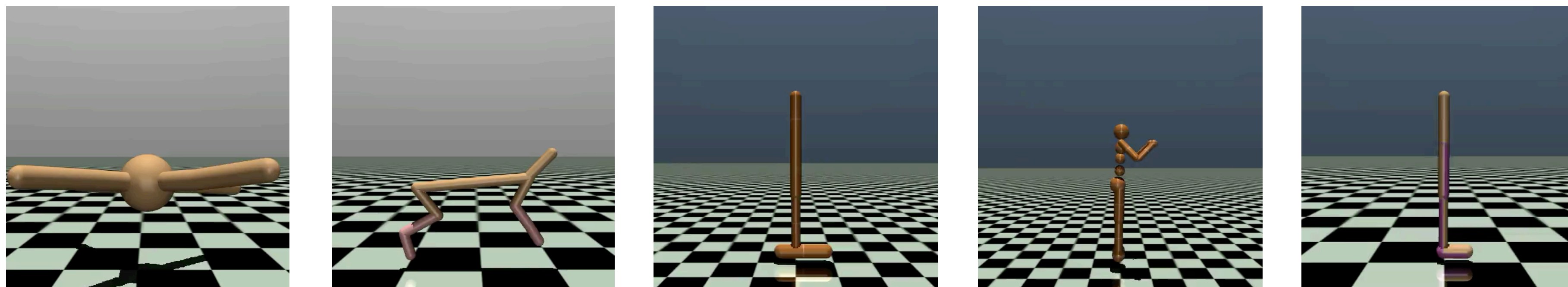
$$MSW_2(\rho, \rho') = \sup_{\psi \in \mathcal{S}^{d-1}} W_2(\psi_{\#}\rho, \psi_{\#}\rho')$$

$$MGSW_2(\rho, \rho') = \sup_{\psi \in \Psi} W_2(g_*^{(\psi)}\rho, g_*^{(\psi)}\rho')$$

$$\inf_{\pi \in \Pi} MGSW_2(\hat{\rho}_E, \hat{\rho}_\pi) = \inf_{\pi \in \Pi} \sup_{\psi \in \Psi} W_2(g_*^{(\psi)}\hat{\rho}_E, g_*^{(\psi)}\hat{\rho}_\pi)$$

Imitation Learning using GSW

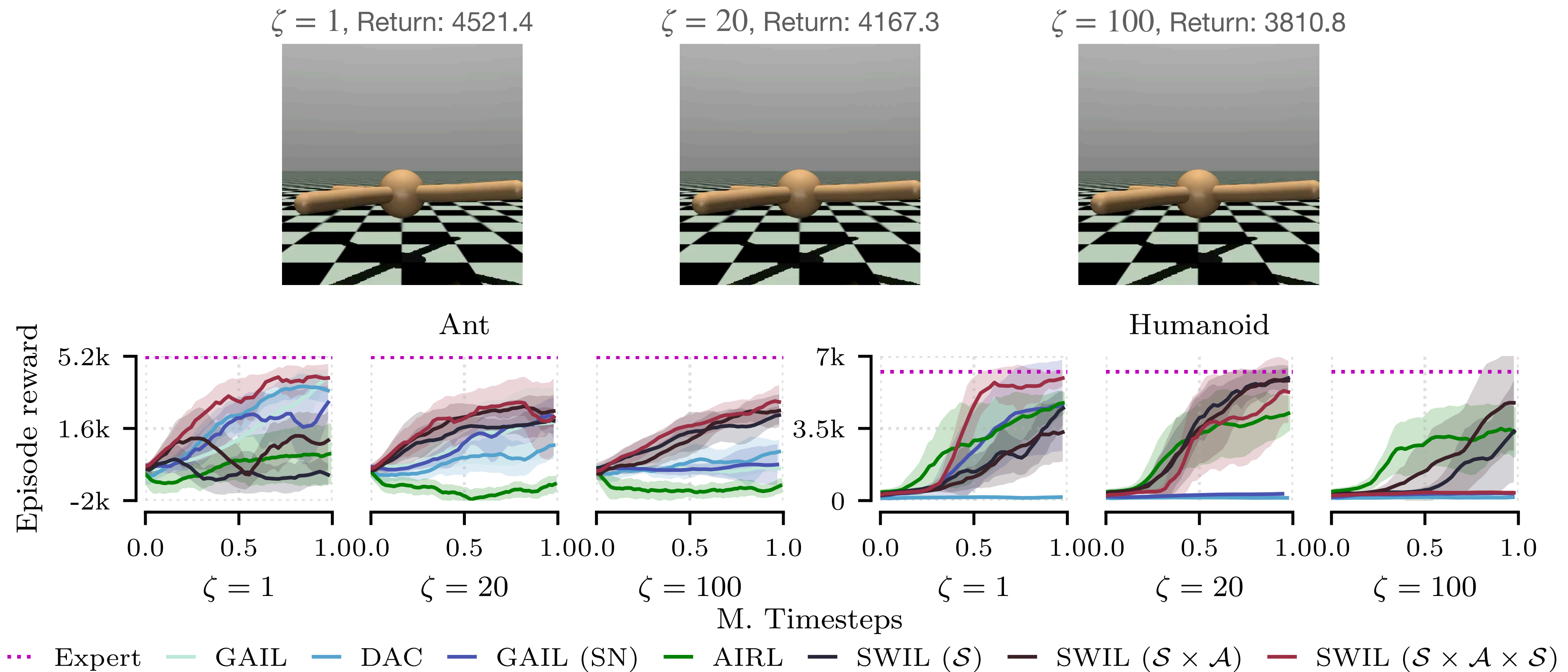
Recovers expert performance from a single trajectory



Environment	ANT	HALFCHEETAH	HOPPER	HUMANOID	WALKER2D
Expert	5159.77	8900.85	3607.17	6249.60	4063.81
BC	978.76 \pm 10.23	221.12 \pm 48.71	504.56 \pm 154.59	347.26 \pm 32.13	299.42 \pm 15.63
DAC	3553.38 \pm 1580.64	3292.43 \pm 1365.91	3400.87 \pm 247.88	169.06 \pm 40.84	3433.62 \pm 362.73
GAIL (SAC)	4261.13 \pm 415.56	2569.26 \pm 1179.47	2866.05 \pm 760.34	147.39 \pm 59.98	2774.17 \pm 563.45
GAIL (SAC-SN)	2956.94 \pm 697.31	804.24 \pm 1441.09	3263.01 \pm 540.13	4643.93 \pm 2203.33	3296.62 \pm 417.84
AIRL (SAC)	361.58 \pm 1424.61	-9.47 \pm 382.93	3197.09 \pm 614.25	4778.57 \pm 404.82	3257.18 \pm 867.42
PWIL	1289.23 \pm 697.31	1089.76 \pm 923.19	2890.14 \pm 430.12	5252.23 \pm 156.22	2452.17 \pm 856.22
GAIL (PPO-SN)	2143.35 \pm 556.98	2863.42 \pm 1155.38	2859.98 \pm 1114.94	425.11 \pm 125.65	2648.65 \pm 1128.32
SWIL	4338.65\pm555.83	7481.79\pm769.22	3585.28\pm66.63	5952.09\pm315.92	3936.35\pm365.69

Imitation Learning using GSW

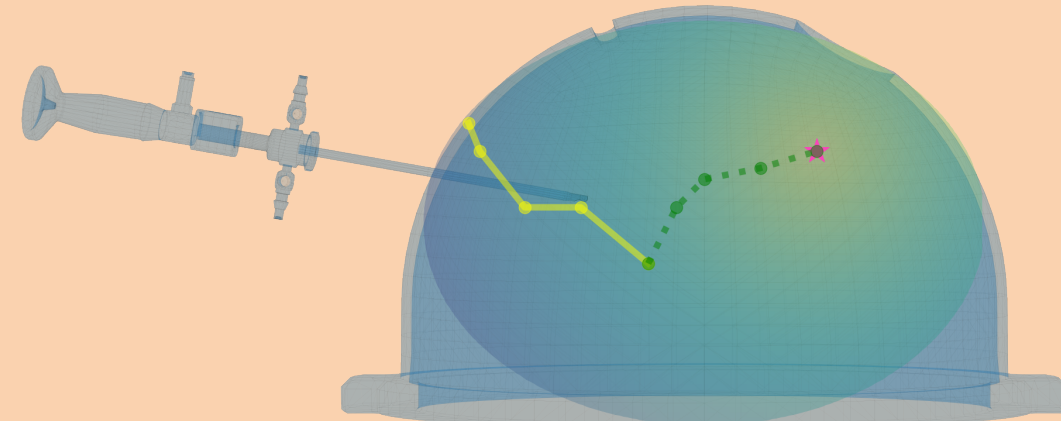
Evaluation: data sparse setting



Performance deteriorates more gracefully under trajectory sparsity

Conclusion

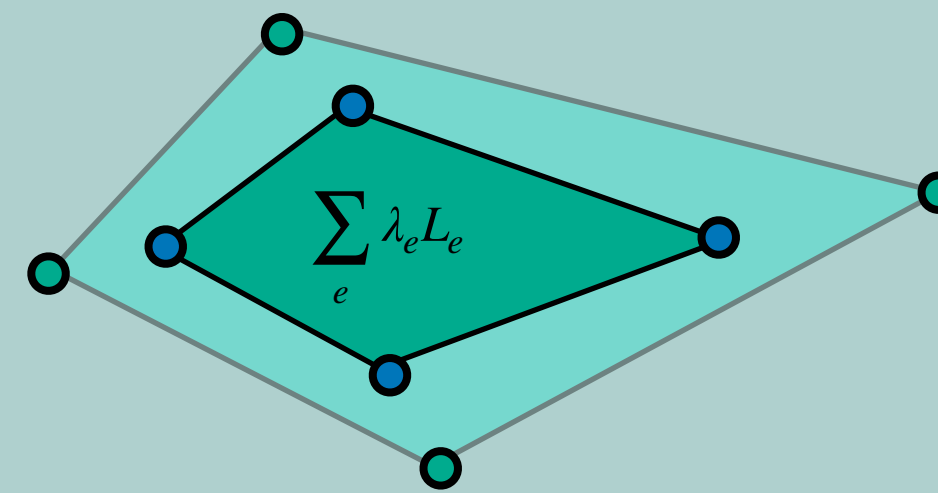
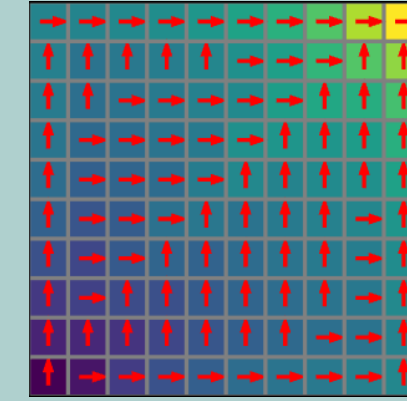
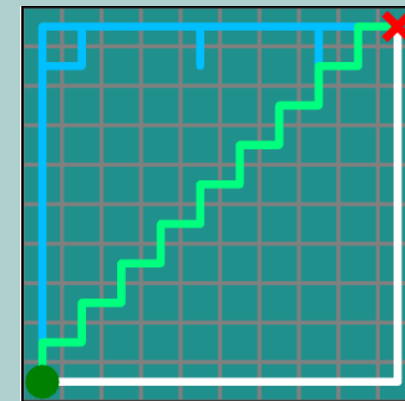
**Application:
Algorithmic RL pipeline in
Surgical Digital Twins**



Ivan Ovinnikov*, Ami Beuret*,
Flavia Cavaliere, Joachim Buhmann

**FASTRL and beyond: a reinforcement learning
pipeline for fundamentals of arthroscopic surgery training**
IJCARs, <https://doi.org/10.1007/s11548-024-03116-z>

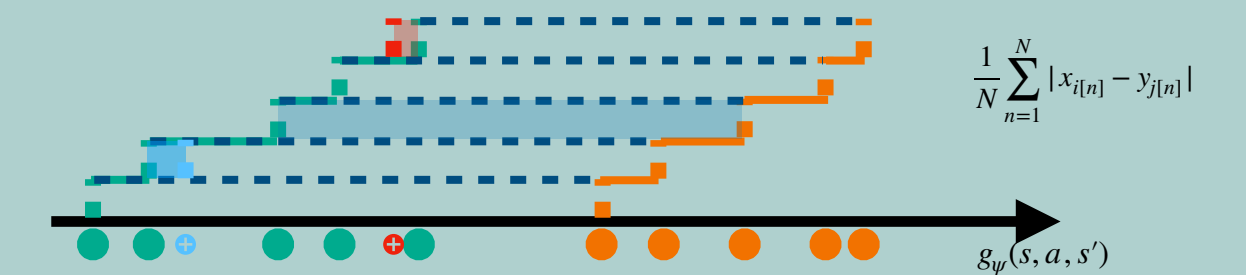
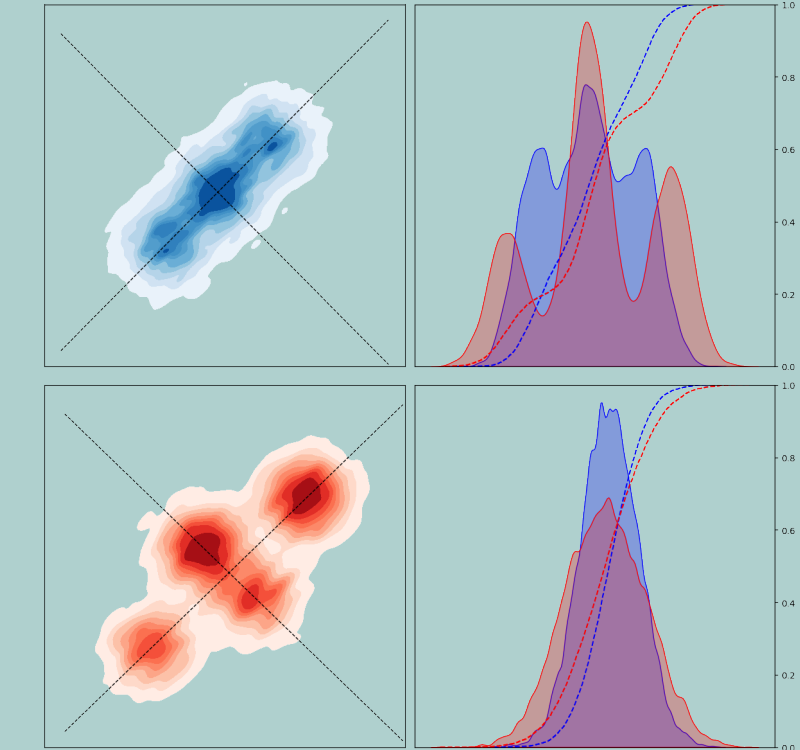
**Method I:
Addressing Reward
Generalization using Causal
Invariance**



Ivan Ovinnikov, Eugene Bykovets,
Joachim Buhmann

**Learning Causally Invariant Reward
Functions from Diverse Demonstrations**
TMLR, in review

**Method II:
Addressing Sample Efficiency
in Imitation Learning using
Sliced Optimal Transport**



**Imitation Learning via Generalized
Sliced Wasserstein Distances**
Preprint

Acknowledgements



Prof. Dr. Joachim Buhmann
Prof. Dr. Andreas Krause
Dr. Raimundo Sierra
Prof. Dr. Markus Gross

Rita Klute

Alexander Terenin
Ami Beuret
Eugene Bykovets
Flavia Cavaliere
Imant Daunhauer
Luis Haug

Martina Vitz
Basil Fierz
Francesco Maggolino

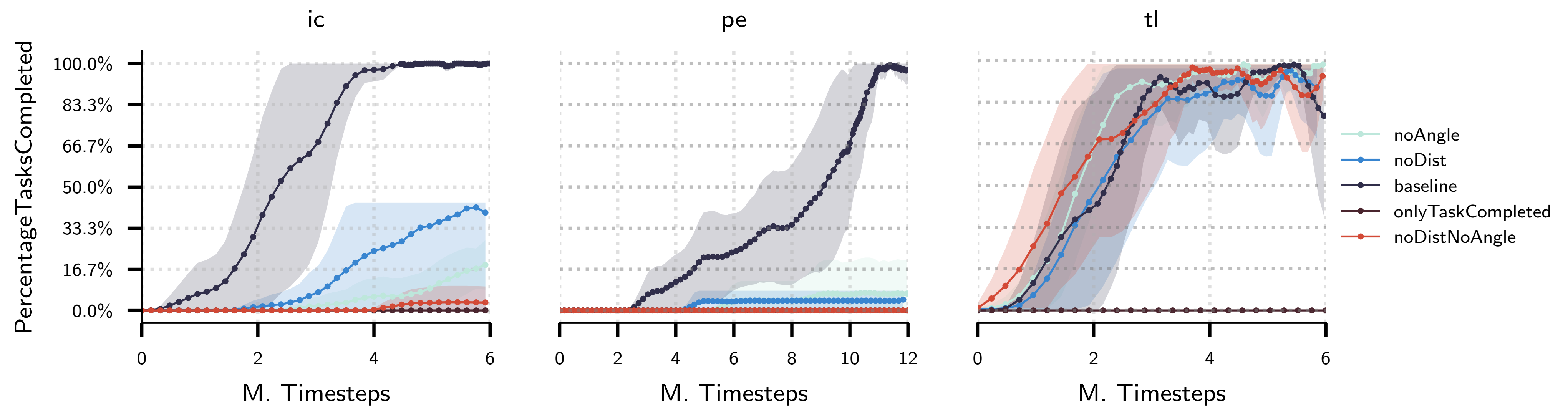
Institute of Machine Learning

Friends and Family

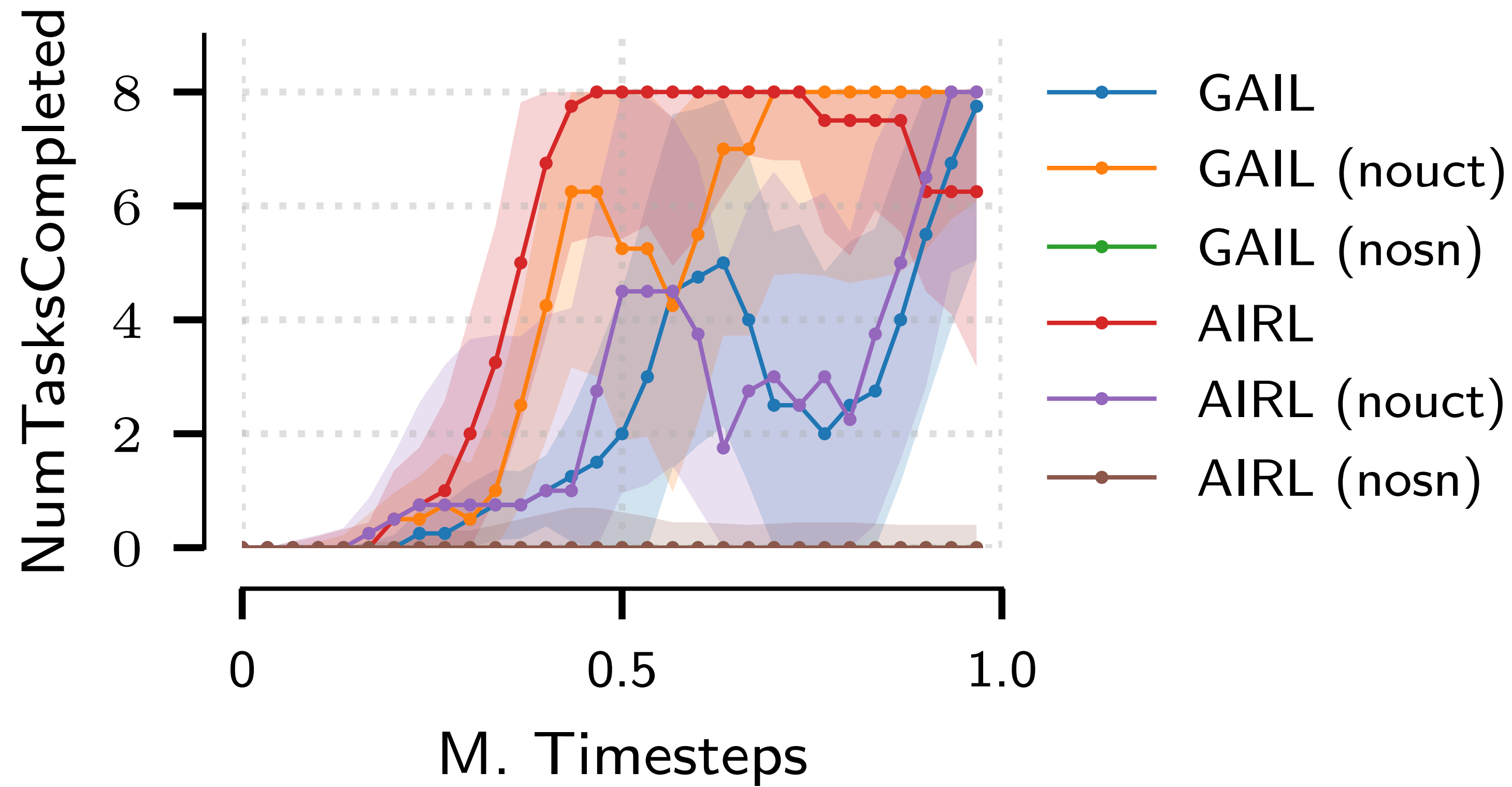
Questions

Appendix

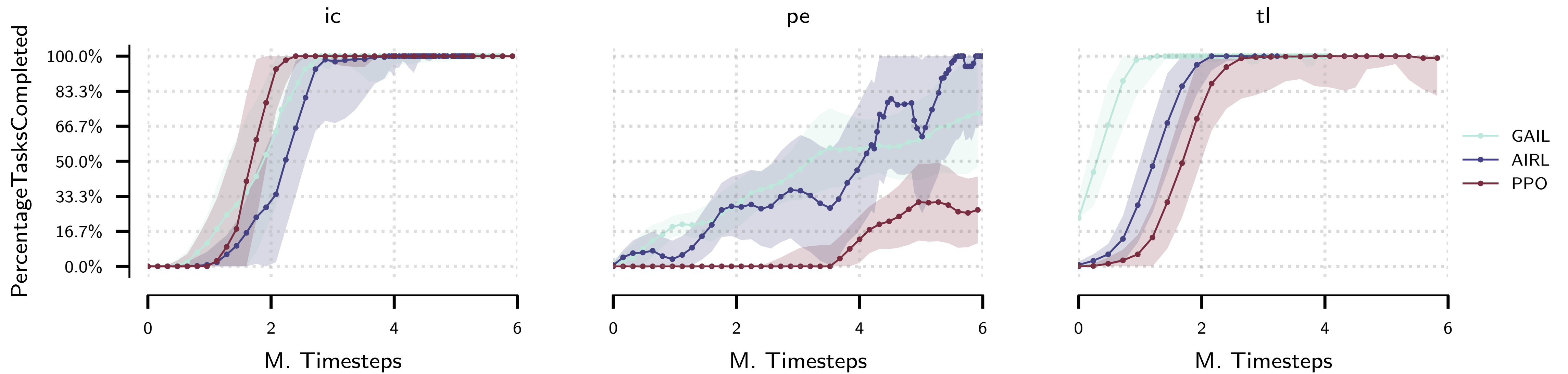
FASTRL Ablations



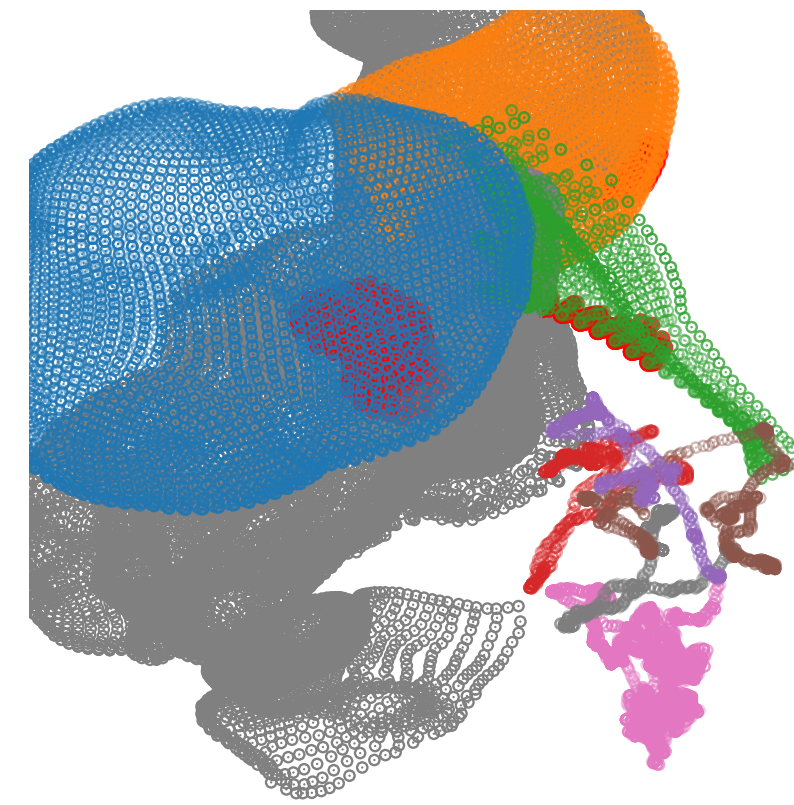
FASTRL : TraceLines SAC



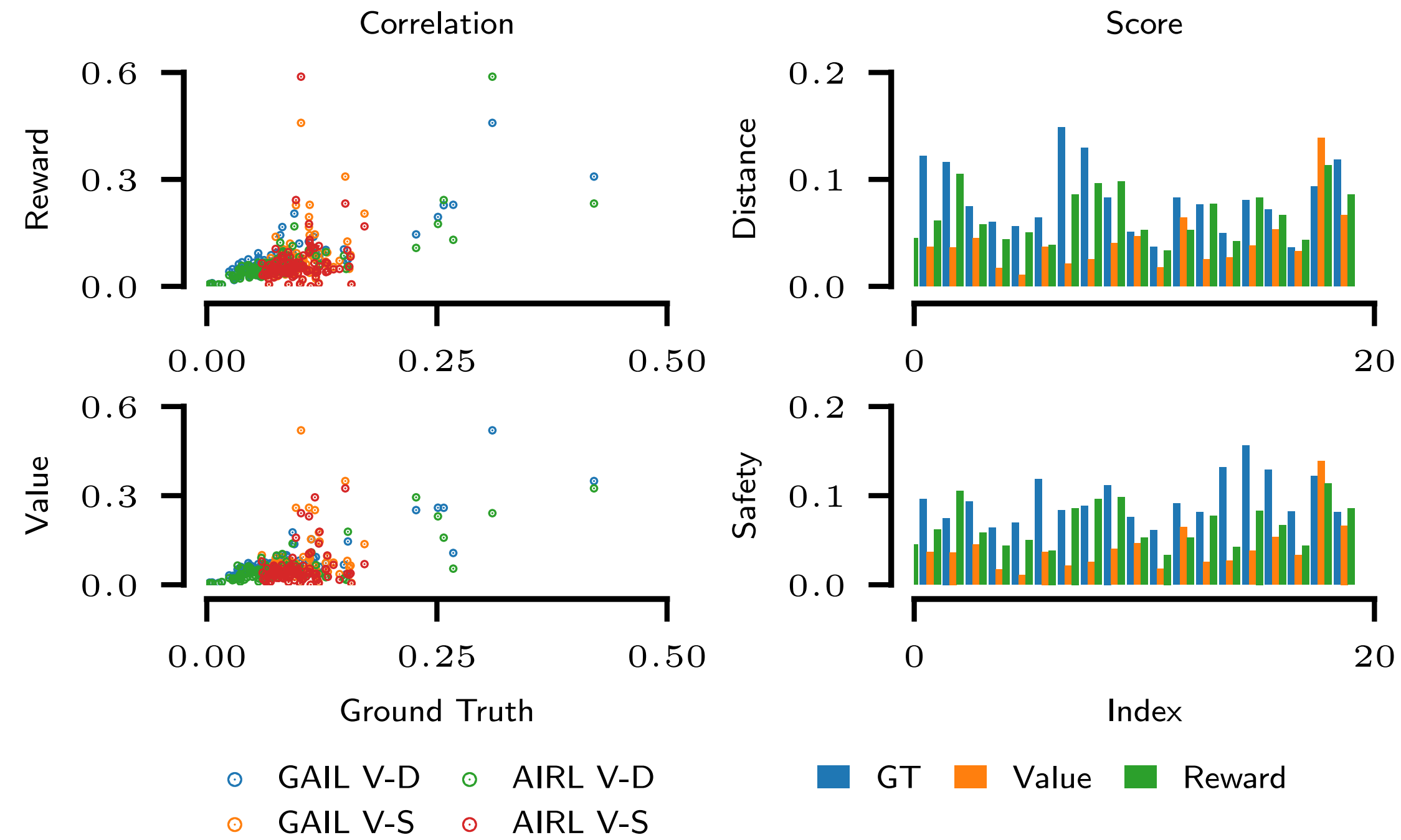
FASTRL : algorithm comparison



Laparos: Davos Dataset Evaluation



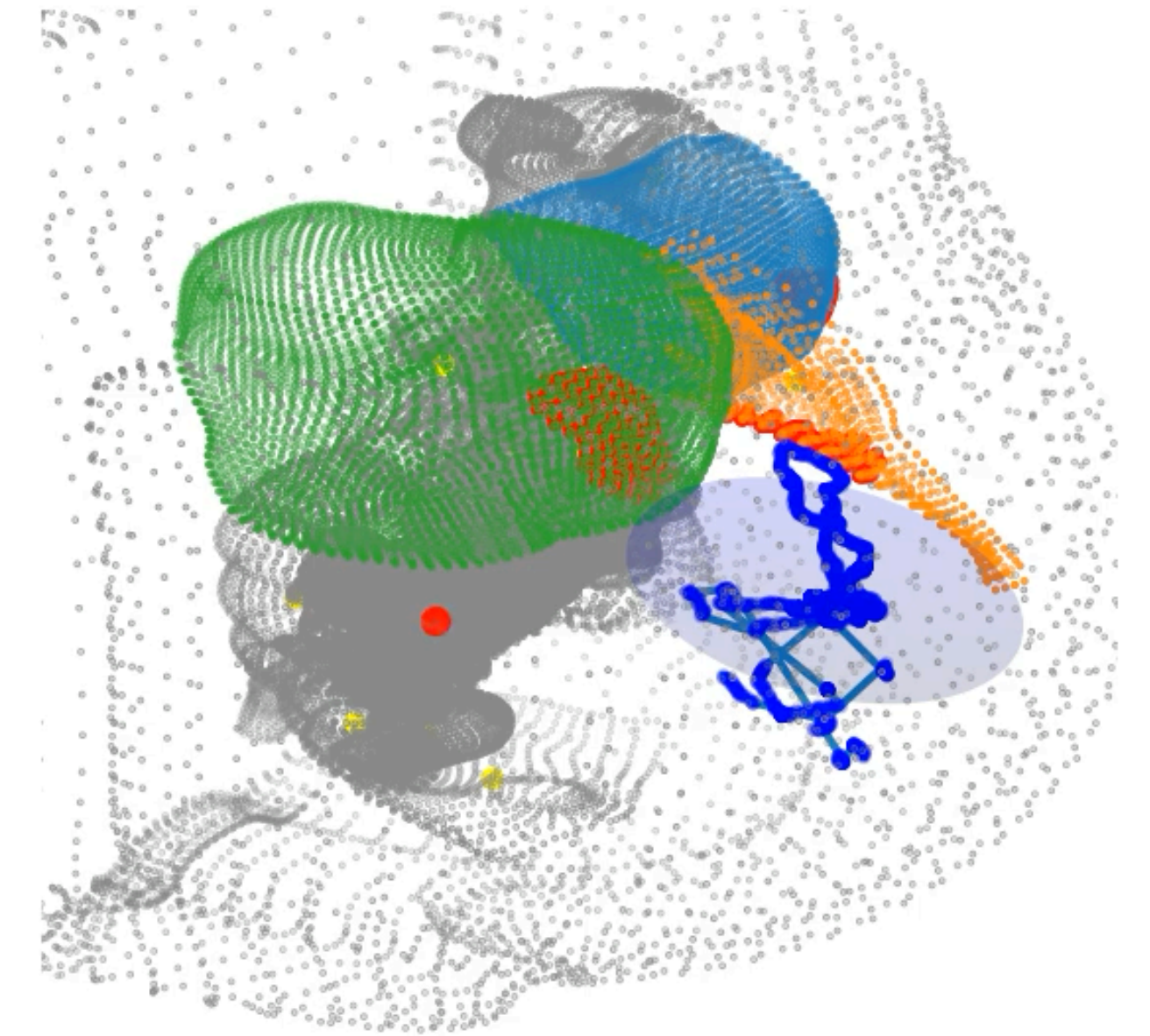
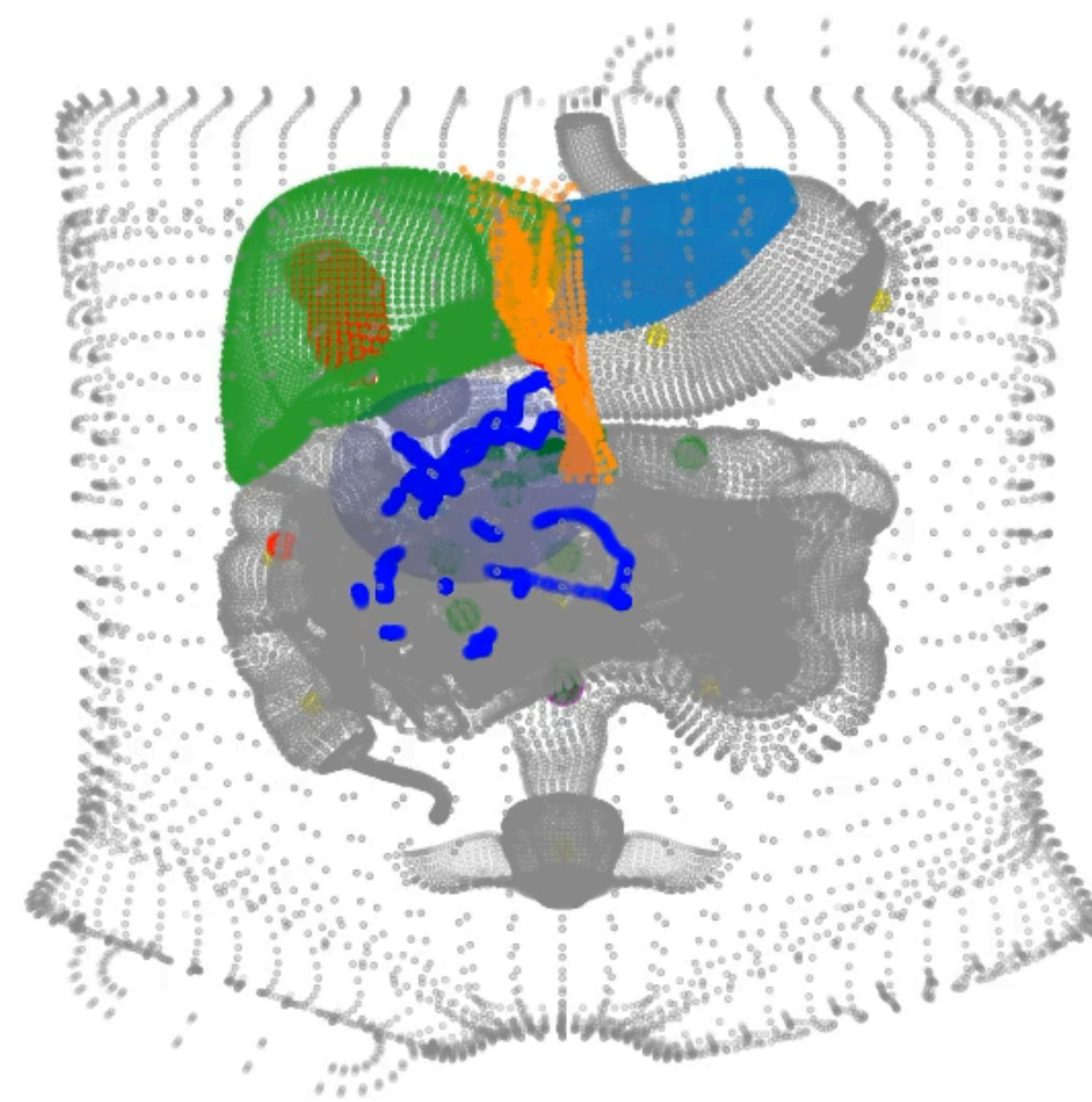
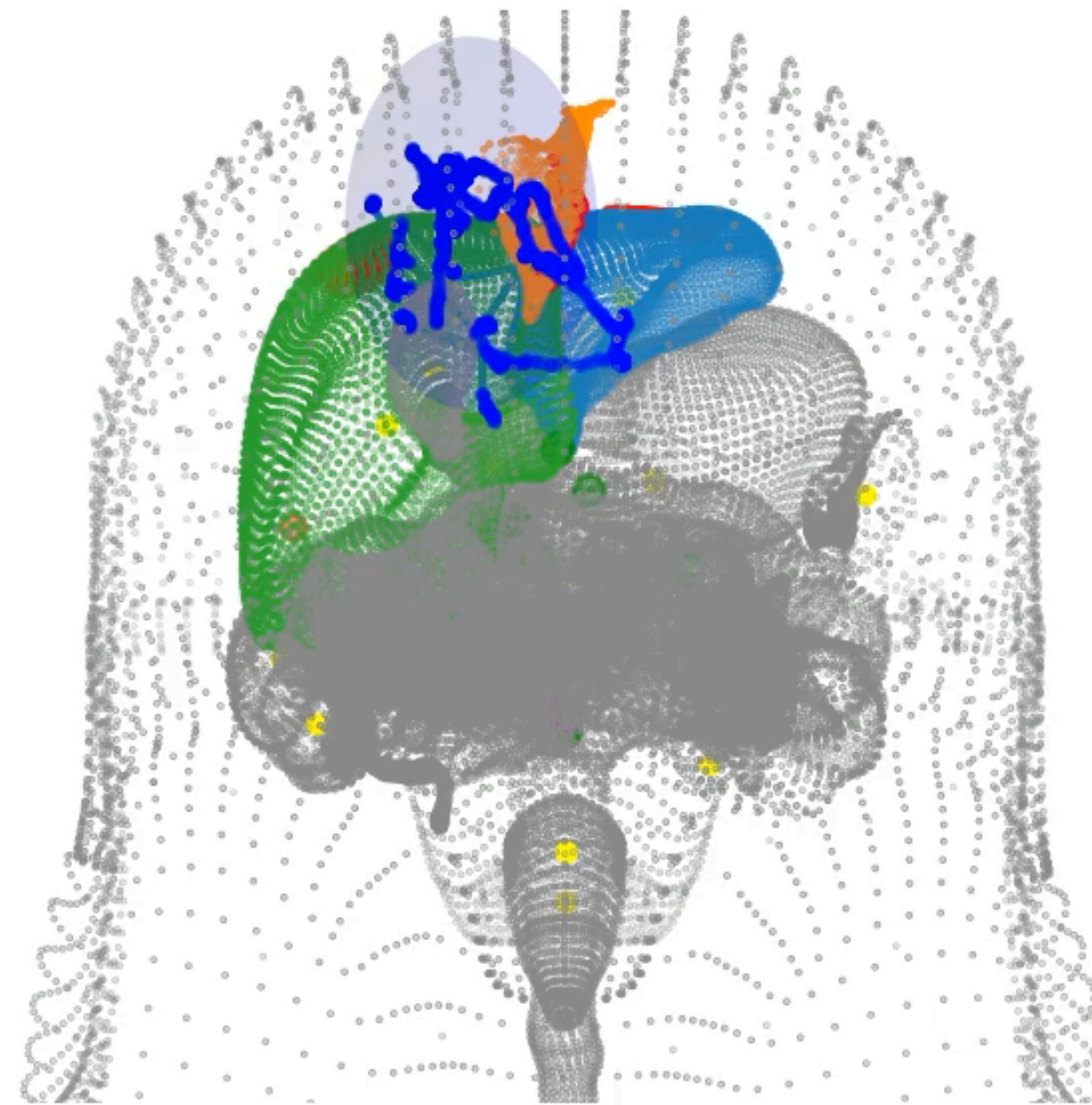
• Target area • Expert traj.



Model	Reward-Dst. (p)	Reward-Sft. (p)	Value-Dst. (p)	Value-Sft. (p)
GAIL	0.73 ± 0.03 (0.00)	0.25 ± 0.02 (0.01)	0.78 ± 0.01 (0.00)	0.24 ± 0.01 (0.02)
AIRL	-0.82 ± 0.01 (0.00)	-0.25 ± 0.01 (0.01)	0.80 ± 0.00 (0.00)	0.24 ± 0.00 (0.02)

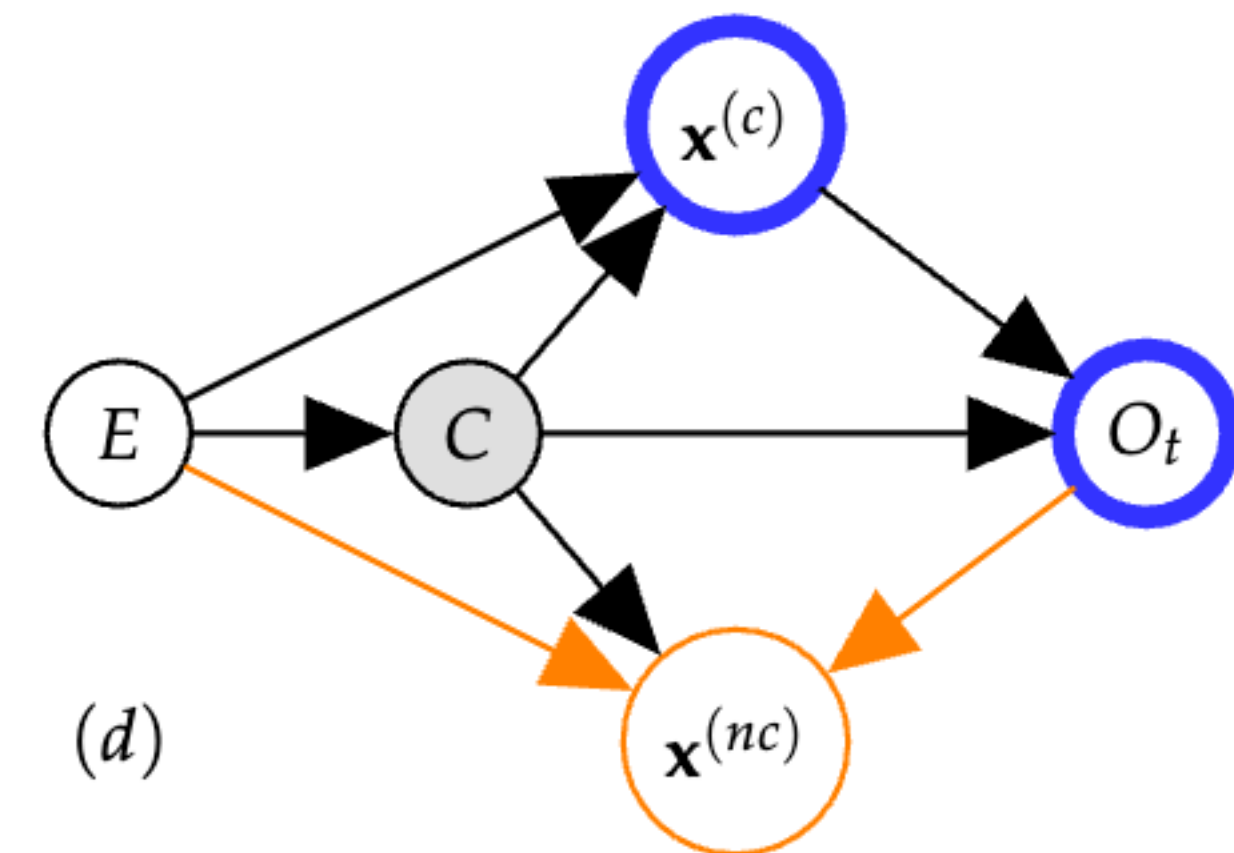
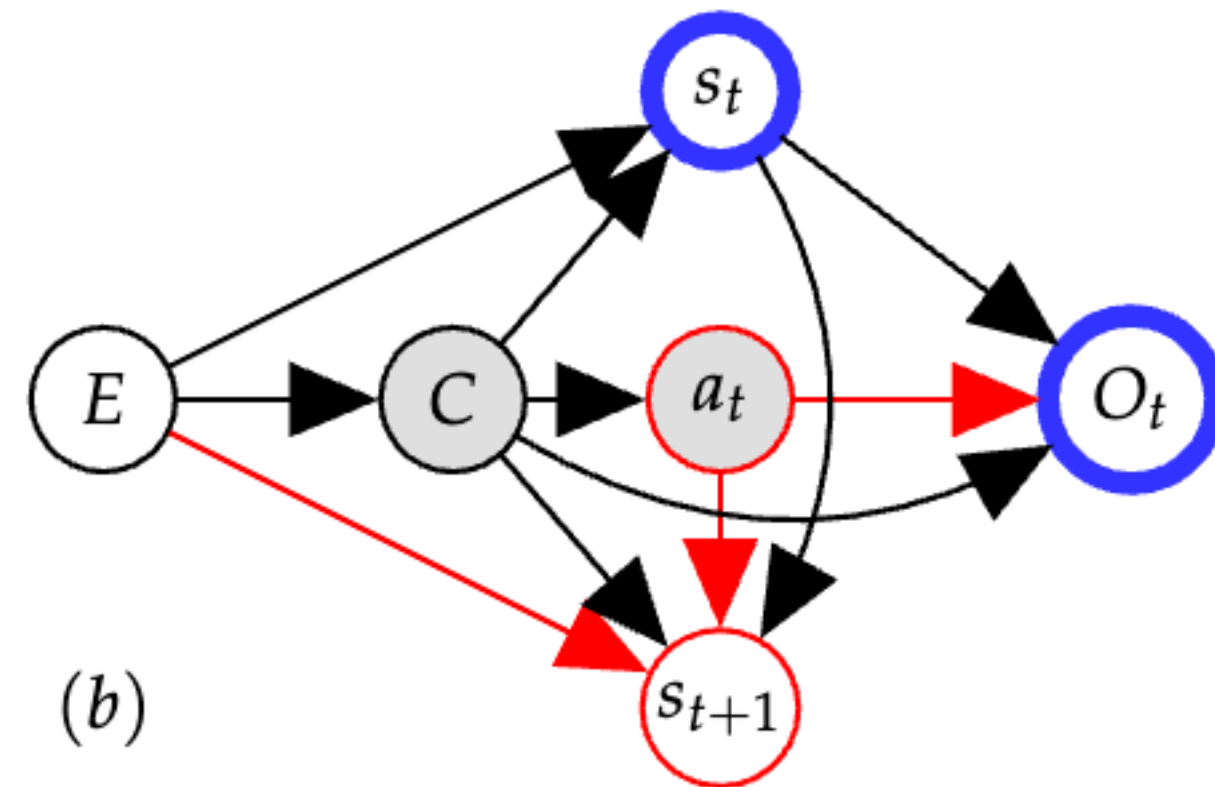
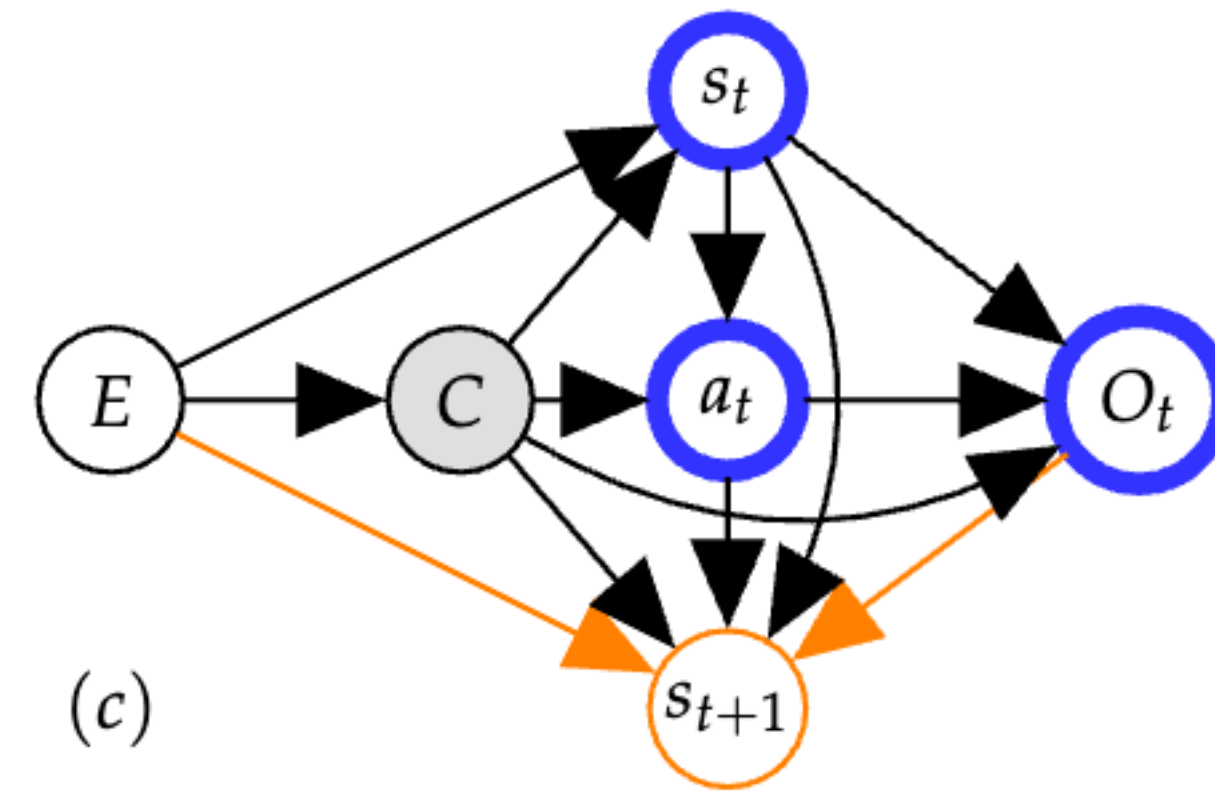
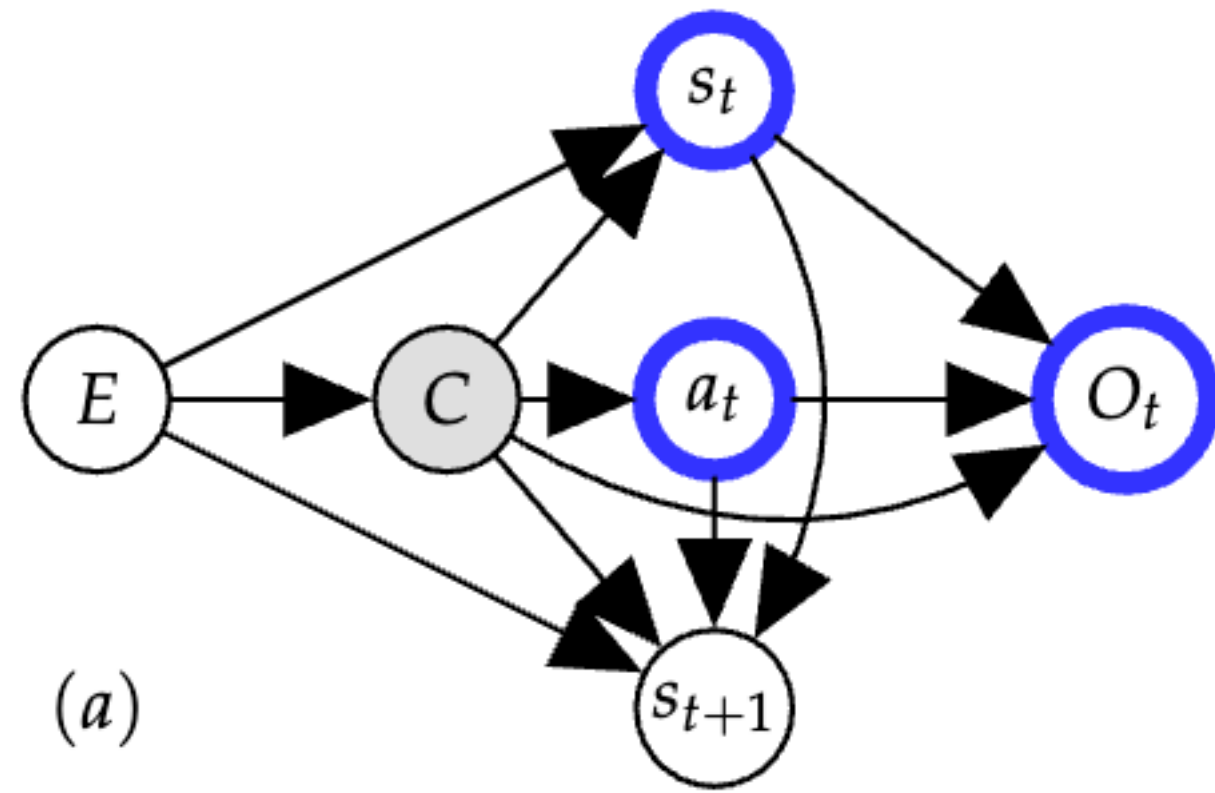
Table 3: Correlation coefficients (p-values in brackets) between laparoscopic trajectories evaluated using recovered rewards and two ground truth metrics (*Dst*: total instrument distance and *Sft*: total safety distance).

SimpleLap

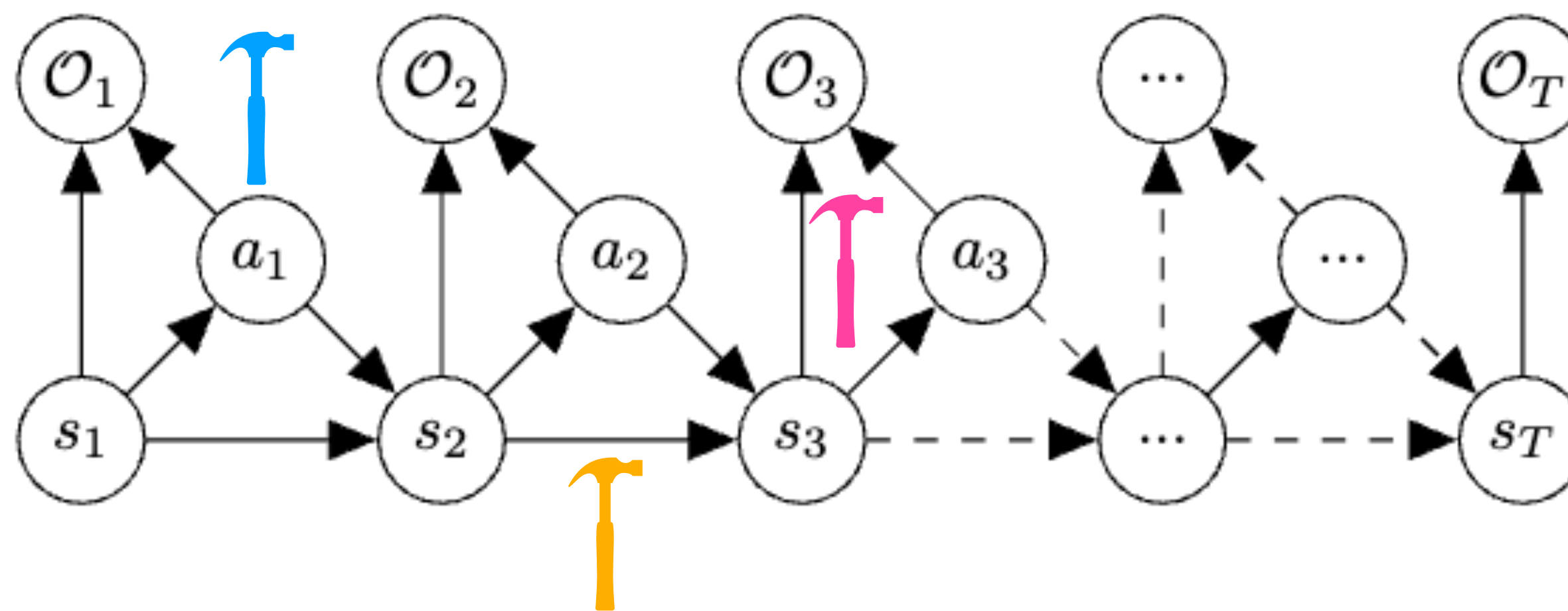


Transition SCM

Spurious correlations



Trajectory SCM



Policy intervention
Dynamics intervention
Reward preference intervention

(a) Trajectory Model

$$\begin{aligned} p(\tau | \mathcal{O}_{1:T}) &\propto p(\tau, \mathcal{O}_{1:T}) = p_0(\mathbf{s}_1) \prod_{t=1}^T p(\mathcal{O}_t = 1 | \mathbf{s}_t, \mathbf{a}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \\ &= \left(p_0(\mathbf{s}_1) \prod_{t=1}^T p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \right) \exp \left(\sum_{t=1}^T r_\psi(\mathbf{s}_t, \mathbf{a}_t) \right) \end{aligned}$$

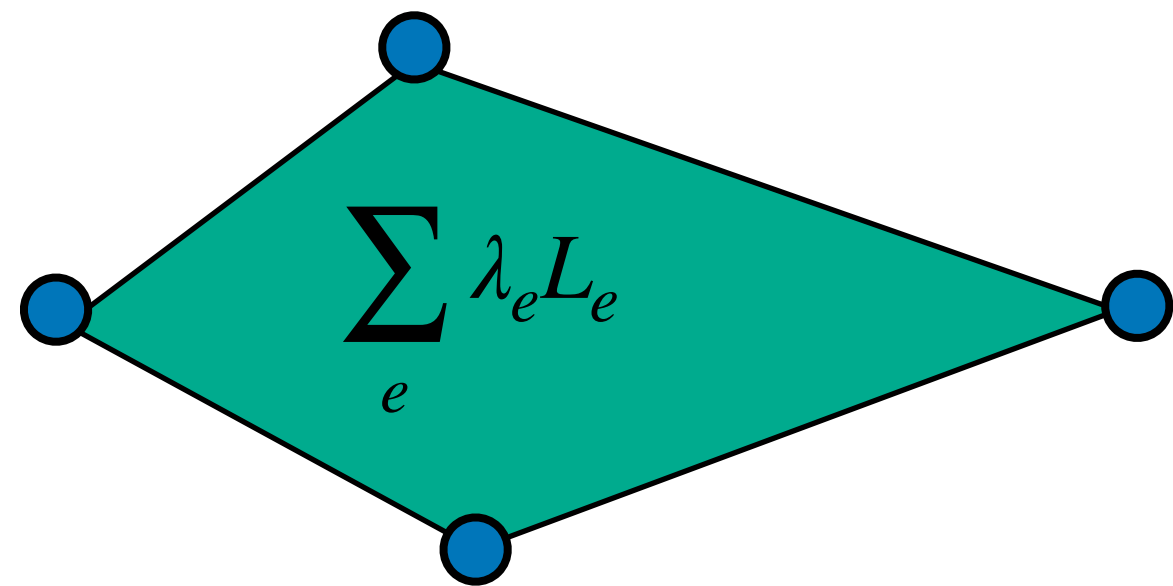
Reward generalization

Leveraging diversity

Goal: recover reward functions which provide meaningful training signal to policies trained across a variety of dynamics

Distributionally robust optimization

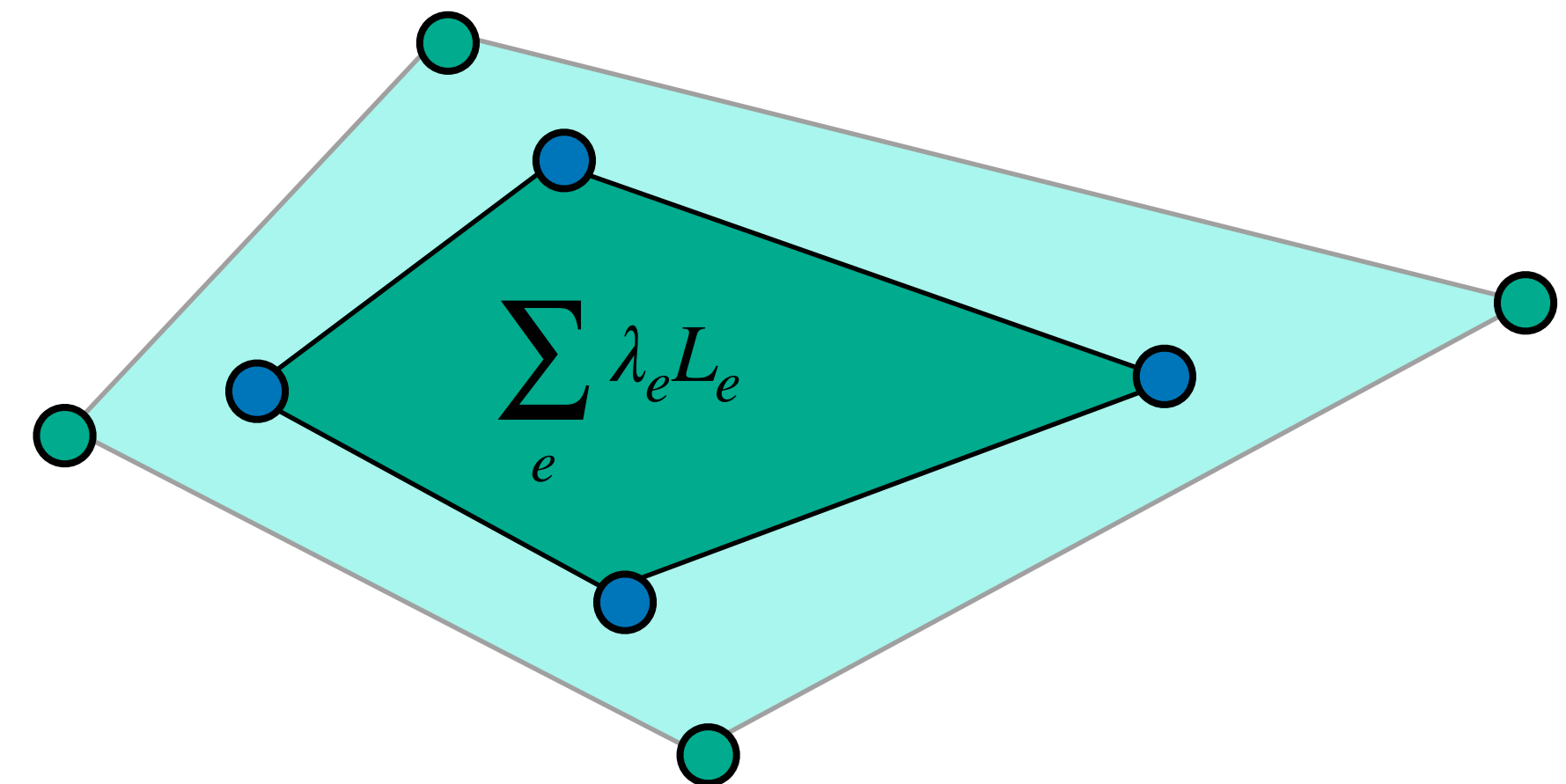
$$\min_f \max_{e \in \mathcal{E}_{tr}} \mathbb{E}_{\xi \sim \mathcal{D}_e} [\mathcal{L}_e(f, \xi)]$$



● Training setting

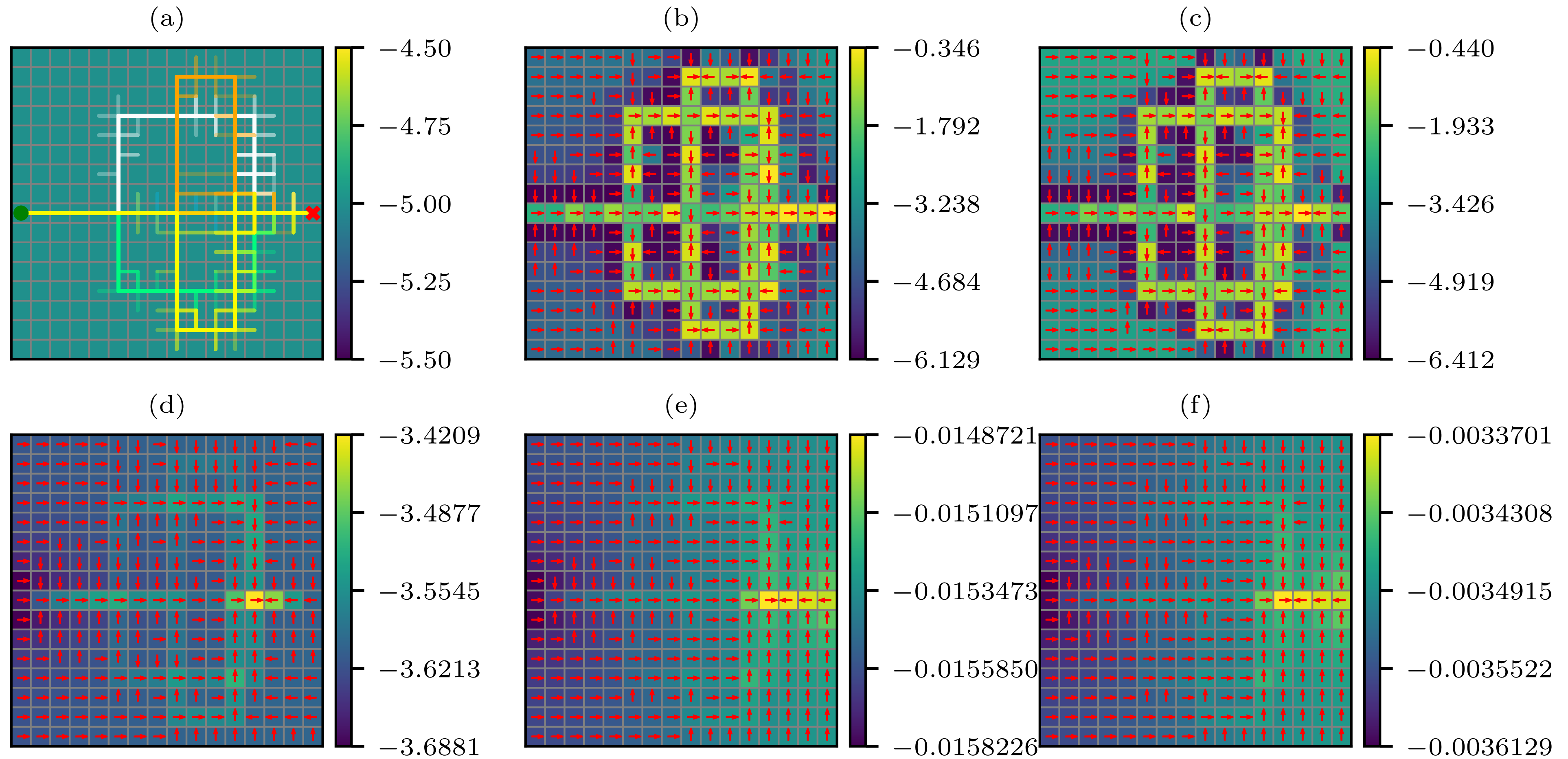
Causal invariance

$$\min_f \mathbb{E}_{\xi \sim \mathcal{D}_e} [\mathcal{L}_e(f, \xi)] \quad \forall e \in \mathcal{E}_{tr}$$



● Generalization setting

CI-IRL: gridworld



CI-IRL: adversarial training

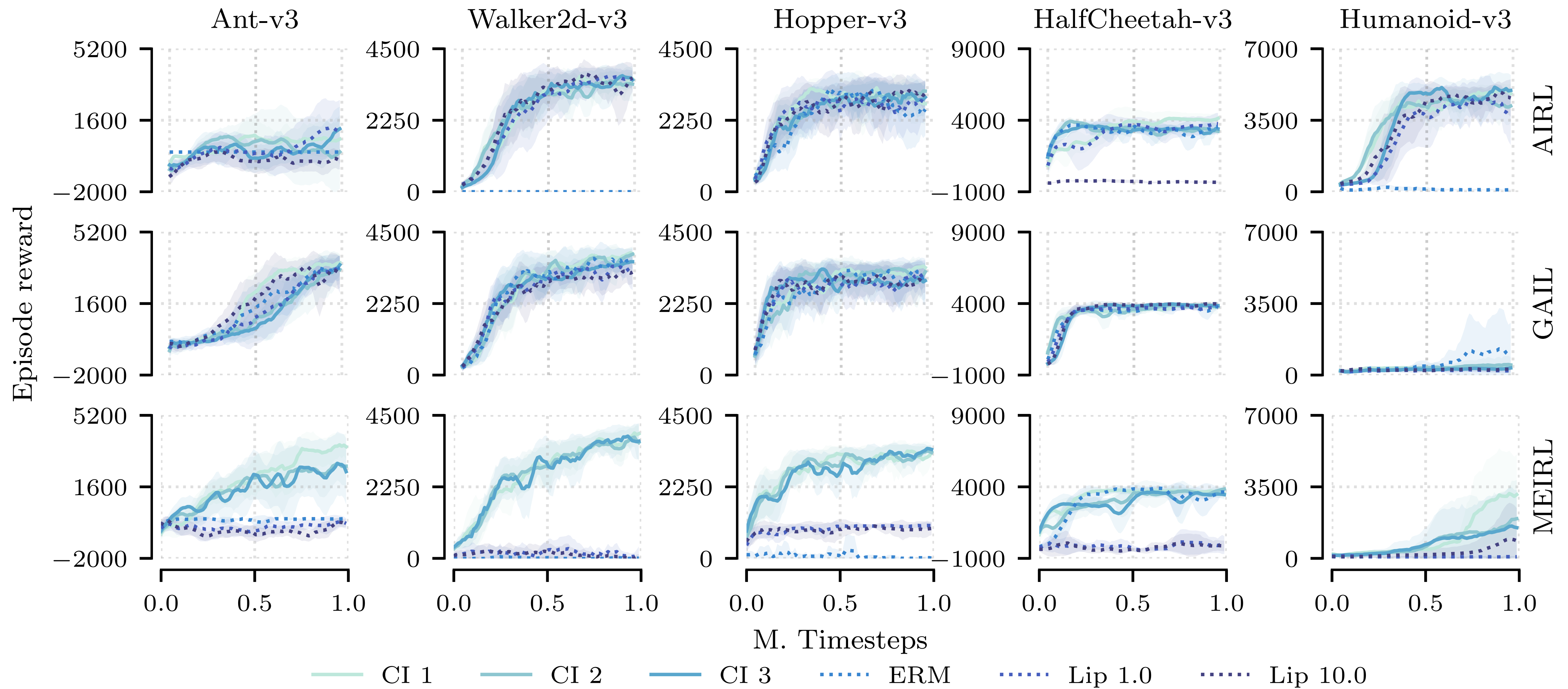


Table: body mass

Environment	ANT-V3	WALKER2D-V3	HOPPER-V3	HALFCHEETAH-V3	HUMANOID-V3
Expert	3168.49 \pm 1715.68	3565.33 \pm 527.40	3119.54 \pm 524.36	4340.61 \pm 2020.14	4774.17 \pm 2063.52
AIRL (ERM)	580.78 \pm 1048.73	-3.29 \pm 0.71	77.64 \pm 88.27	2046.29 \pm 460.98	4451.74 \pm 1759.31
AIRL (Lip)	1194.04 \pm 1583.08	3388.48 \pm 1586.45	3382.91 \pm 234.02	4388.94 \pm 726.69	1788.85 \pm 1643.00
AIRL (CI)	1880.42 \pm 935.15	4162.70 \pm 517.13	3334.91 \pm 221.80	4477.97 \pm 532.72	5107.54 \pm 119.31
GAIL (ERM)	-746.31 \pm 468.03	328.44 \pm 66.02	1637.50 \pm 1419.59	886.25 \pm 404.82	122.62 \pm 71.53
GAIL (Lip)	220.97 \pm 524.83	553.36 \pm 277.24	1832.39 \pm 832.32	1403.77 \pm 1282.75	77.24 \pm 4.66
GAIL (CI)	230.43 \pm 565.68	1172.57 \pm 539.86	2636.65 \pm 1114.94	2365.55 \pm 1679.64	549.63 \pm 1692.08
MEIRL (ERM)	-66.66 \pm 112.03	169.11 \pm 344.87	3.22 \pm 0.22	-177.39 \pm 211.43	55.99 \pm 3.45
MEIRL (Lip)	-365.41 \pm 143.70	917.14 \pm 132.05	1045.40 \pm 54.76	-335.10 \pm 84.66	1001.49 \pm 1889.60
MEIRL (CI)	153.43 \pm 1134.46	2520.24 \pm 994.27	2351.07 \pm 679.37	1371.59 \pm 1469.01	3099.51 \pm 2411.21

TABLE 4.1: **Policy rollout results using ground truth reward for perturbed MuJoCo environments** after being trained for 1M timesteps using the rewards recovered from the different discriminators in section 4.3.2. Here, the *body mass* parameter is perturbed with a noise magnitude of $\varepsilon = 0.2$. The results are averaged over 10 rollouts and obtained by training the model using five different random seeds.

Table: actuator control range

Environment	ANT	WALKER2D	HOPPER	HALFCHEETAH	HUMANOID
Expert	3168.49 \pm 1715.68	3565.33 \pm 527.40	3119.54 \pm 524.36	4340.61 \pm 2020.14	4774.17 \pm 2063.52
AIRL (ERM)	1279.87 \pm 1281.66	-0.41 \pm 3.07	37.62 \pm 62.71	2536.90 \pm 212.27	357.19 \pm 135.66
AIRL (Lip)	809.60 \pm 1425.76	2779.73 \pm 1303.91	2784.28 \pm 510.50	4175.49 \pm 918.92	312.24 \pm 90.05
AIRL (CI)	2166.58 \pm 1471.70	3897.58 \pm 831.24	2884.34 \pm 130.60	4470.17 \pm 731.81	2730.60 \pm 982.13
GAIL (ERM)	-641.93 \pm 284.65	1180.57 \pm 1413.21	491.27 \pm 565.63	1862.67 \pm 1026.74	1219.94 \pm 1784.26
GAIL (Lip)	-92.74 \pm 363.44	1672.06 \pm 1263.95	2028.07 \pm 1004.96	3638.51 \pm 1164.42	93.01 \pm 24.01
GAIL (CI)	2486.00 \pm 2078.11	2660.74 \pm 866.36	2985.44 \pm 280.70	3979.90 \pm 2494.31	2986.70 \pm 2389.78
MEIRL (ERM)	-10.63 \pm 3.35	-3.83 \pm 0.45	3.17 \pm 0.27	-36.66 \pm 489.57	58.40 \pm 0.15
MEIRL (Lip)	-411.65 \pm 244.20	832.80 \pm 311.20	1073.22 \pm 138.00	-261.74 \pm 164.78	1064.41 \pm 2011.91
MEIRL (CI)	133.50 \pm 969.39	2286.80 \pm 1040.59	2551.59 \pm 1131.76	3303.82 \pm 2332.89	3058.78 \pm 2286.41

TABLE 4.3: Policy rollout results using ground truth reward for perturbed MuJoCo environments after being trained for 1M timesteps using the rewards recovered from the different discriminators in section 4.3.2. Here, the *actuator control range* parameter is perturbed with a noise magnitude of $\varepsilon = 0.2$. The results are averaged over 10 rollouts and obtained by training the model using five different random seeds.

Table: geometry friction

Environment	ANT	WALKER2D	HOPPER	HALFCHEETAH	HUMANOID
Expert	3168.49 \pm 1715.68	3565.33 \pm 527.40	3119.54 \pm 524.36	4340.61 \pm 2020.14	4774.17 \pm 2063.52
AIRL (ERM)	603.00 \pm 909.86	-3.87 \pm 0.44	149.17 \pm 151.96	2141.85 \pm 942.19	4507.23 \pm 659.04
AIRL (Lip)	283.73 \pm 1294.15	3429.17 \pm 372.29	3311.25 \pm 128.82	4659.44 \pm 533.32	1432.57 \pm 952.87
AIRL (CI)	1434.01 \pm 1530.65	4167.15 \pm 721.26	3288.86 \pm 149.76	4737.72 \pm 749.52	4756.93 \pm 368.15
GAIL (ERM)	-421.27 \pm 752.40	910.12 \pm 951.80	939.05 \pm 959.38	1563.17 \pm 1245.51	871.61 \pm 1002.21
GAIL (Lip)	-172.69 \pm 196.92	1065.02 \pm 1672.12	2541.08 \pm 906.67	4795.59 \pm 1018.65	89.51 \pm 16.87
GAIL (CI)	1148.32 \pm 1938.45	2395.43 \pm 1282.70	3068.00 \pm 459.50	4037.08 \pm 983.32	3385.58 \pm 2279.28
MEIRL (ERM)	-103.10 \pm 188.90	-3.43 \pm 0.15	3.36 \pm 0.19	191.23 \pm 630.81	56.41 \pm 2.24
MEIRL (Lip)	-252.29 \pm 184.32	865.84 \pm 308.25	1128.52 \pm 125.06	-284.92 \pm 204.61	991.70 \pm 1871.72
MEIRL (CI)	-191.76 \pm 912.46	2546.37 \pm 1073.22	2425.38 \pm 1070.78	1724.44 \pm 2057.07	2155.55 \pm 2281.06

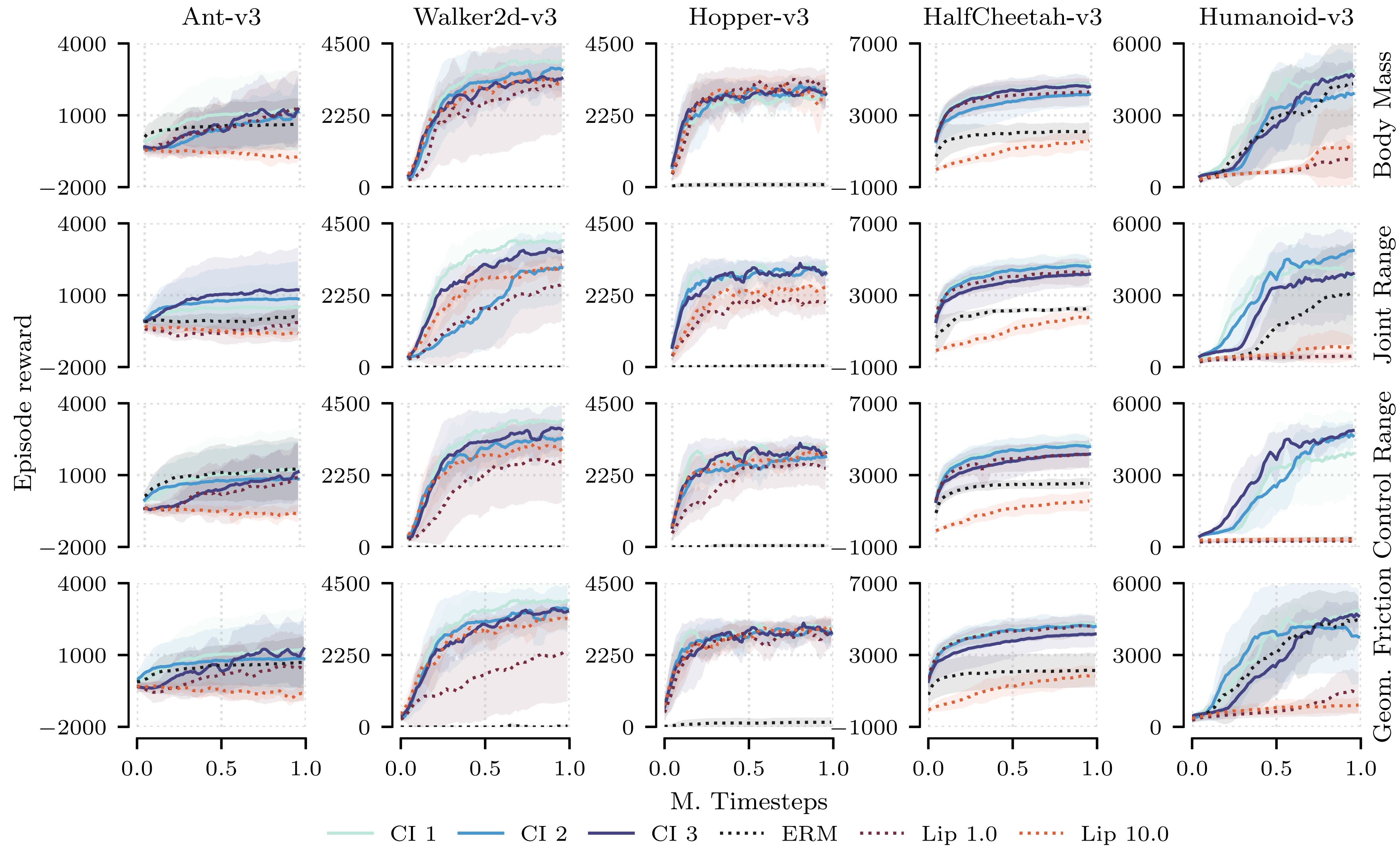
TABLE 4.4: Policy rollout results using ground truth reward for perturbed MuJoCo environments after being trained for 1M timesteps using the rewards recovered from the different discriminators in section 4.3.2. Here, the *geometry friction* parameter is perturbed with a noise magnitude of $\varepsilon = 0.2$. The results are averaged over 10 rollouts and obtained by training the model using five different random seeds.

Table: joint range

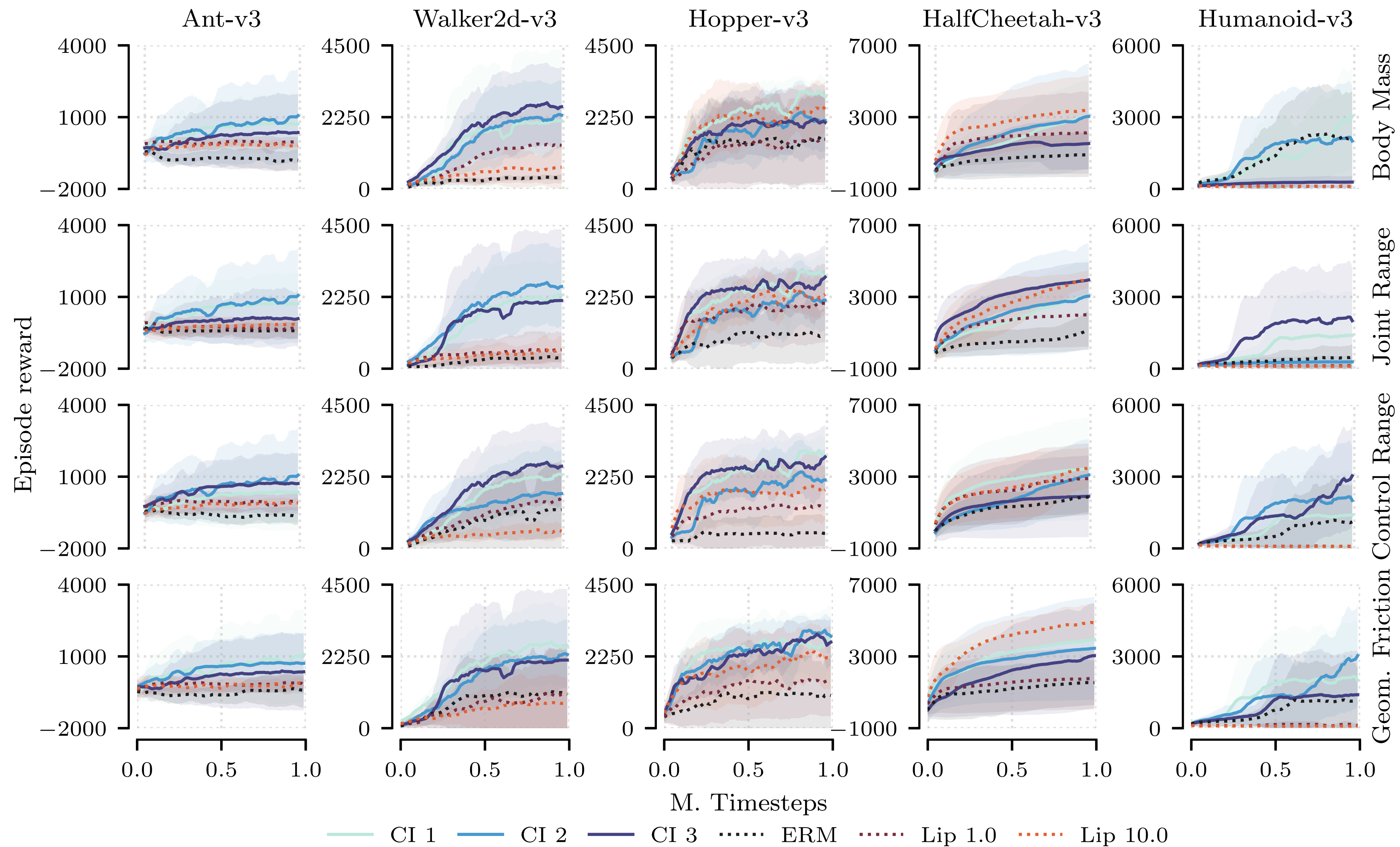
Environment	ANT	WALKER2D	HOPPER	HALFCHEETAH	HUMANOID
Expert	3168.49 \pm 1715.68	3565.33 \pm 527.40	3119.54 \pm 524.36	4340.61 \pm 2020.14	4774.17 \pm 2063.52
AIRL (ERM)	-18.73 \pm 255.23	-3.50 \pm 1.66	48.82 \pm 76.24	2310.93 \pm 118.75	3261.17 \pm 2272.68
AIRL (Lip)	-213.89 \pm 738.75	3202.09 \pm 185.98	2544.28 \pm 445.86	4293.70 \pm 666.33	710.88 \pm 356.04
AIRL (CI)	155.62 \pm 875.01	3670.21 \pm 599.00	2906.77 \pm 490.33	4653.89 \pm 762.09	4022.43 \pm 671.37
GAIL (ERM)	-486.05 \pm 388.06	359.23 \pm 254.85	1047.76 \pm 871.98	1160.78 \pm 1134.26	508.86 \pm 601.99
GAIL (Lip)	-208.97 \pm 252.99	577.40 \pm 550.86	2339.87 \pm 465.97	4168.13 \pm 472.22	122.49 \pm 82.43
GAIL (CI)	1021.26 \pm 1845.56	3479.99 \pm 1242.39	2976.49 \pm 417.33	5581.43 \pm 1442.39	2170.96 \pm 2425.51
MEIRL (ERM)	-57.30 \pm 93.23	-3.69 \pm 0.18	3.17 \pm 0.52	347.55 \pm 790.54	54.70 \pm 1.92
MEIRL (Lip)	-554.25 \pm 82.05	622.45 \pm 464.25	1136.95 \pm 209.97	-297.84 \pm 71.23	1014.42 \pm 1915.24
MEIRL (CI)	-337.17 \pm 1310.57	2292.94 \pm 1521.43	2800.37 \pm 666.09	1650.50 \pm 1683.38	2935.90 \pm 2262.13

TABLE 4.2: Policy rollout results using ground truth reward for perturbed MuJoCo environments after being trained for 1M timesteps using the rewards recovered from the different discriminators in section 4.3.2. Here, the *joint range* parameter is perturbed with a noise magnitude of $\varepsilon = 0.2$. The results are averaged over 10 rollouts and obtained by training the model using five different random seeds.

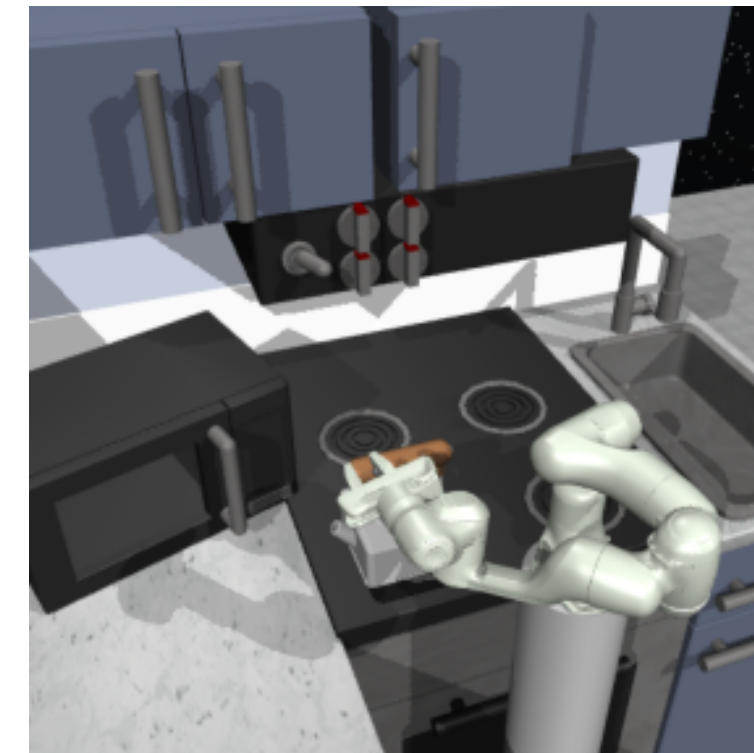
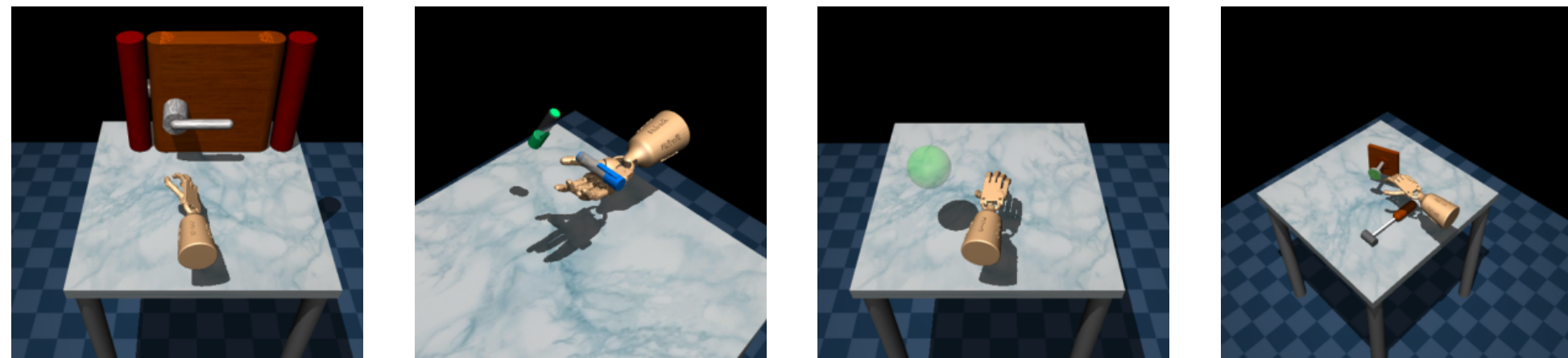
AIRL training dynamics I



GAIL Training Dynamics I

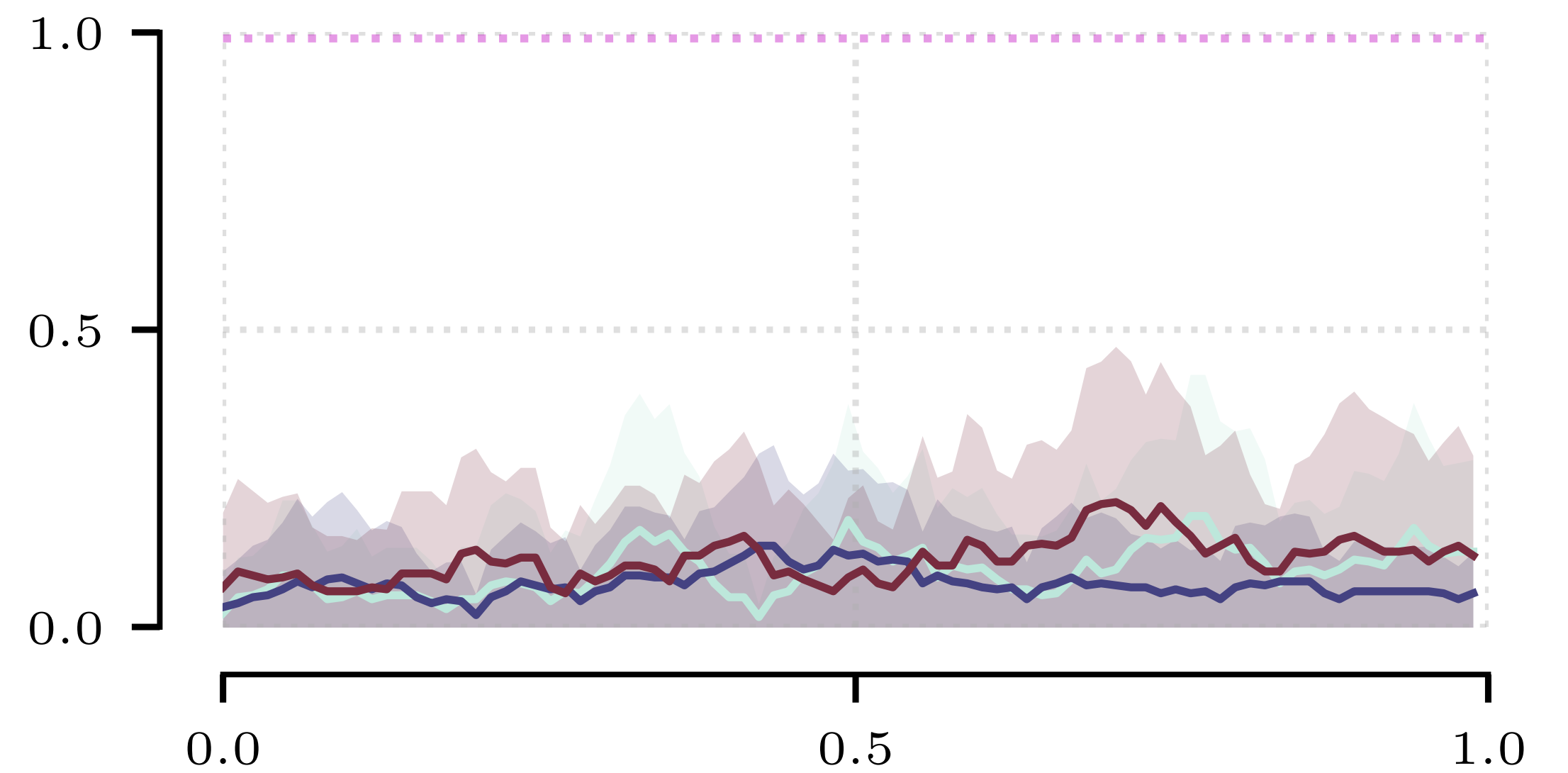
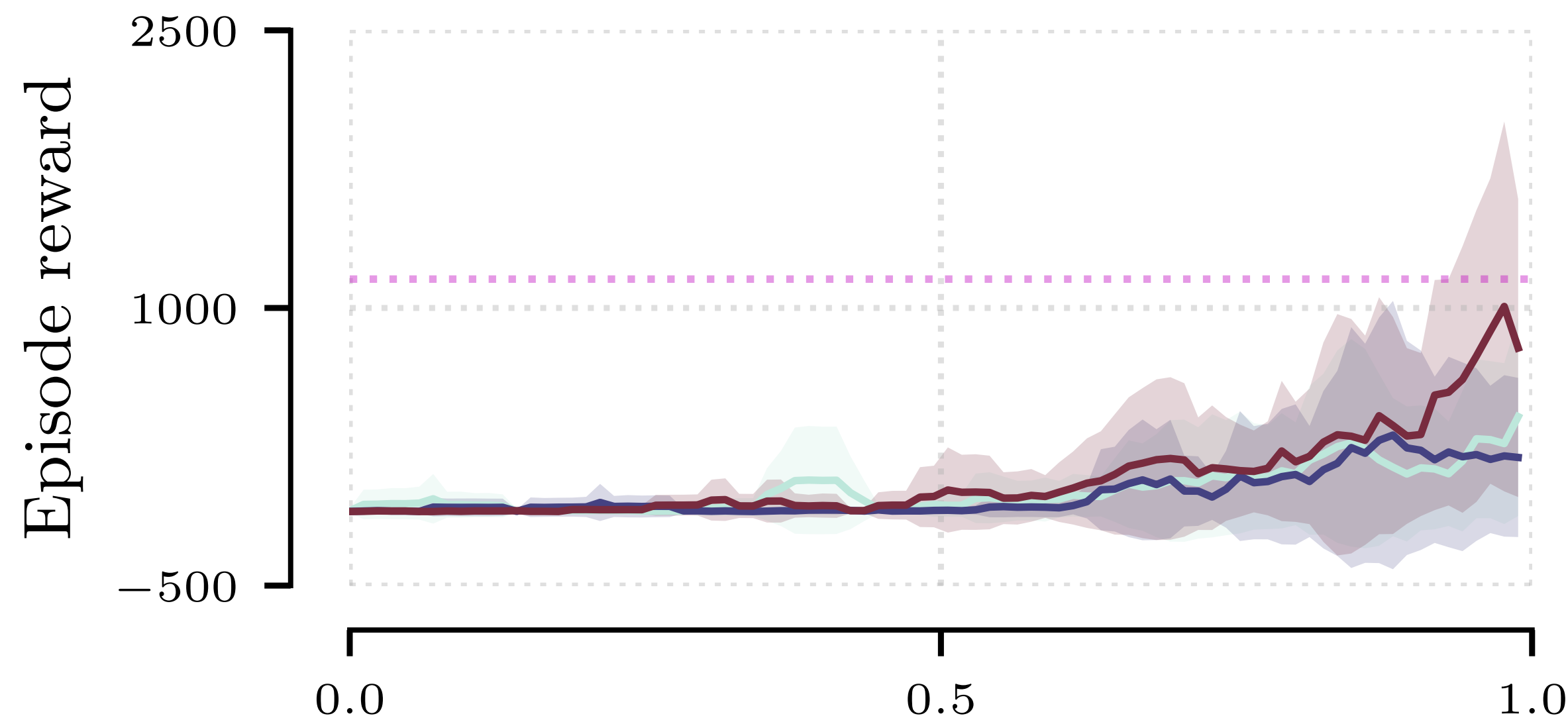


Adroit / Franka Results



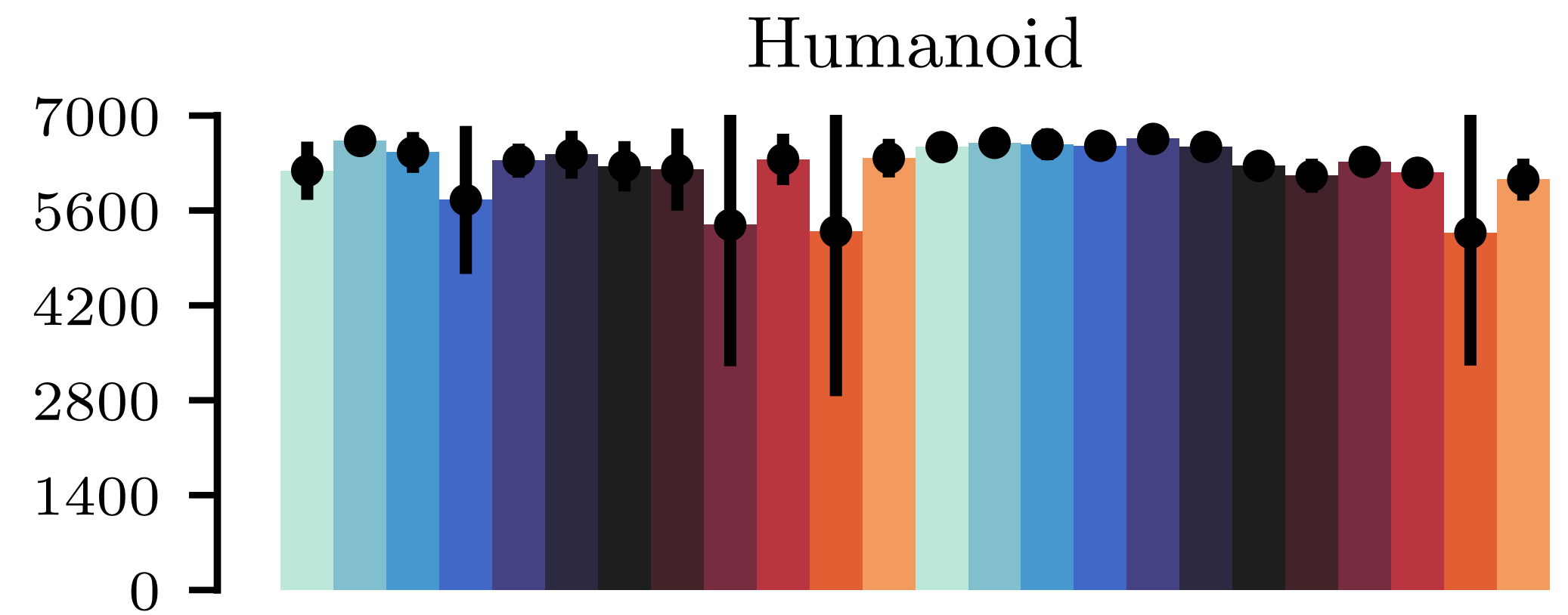
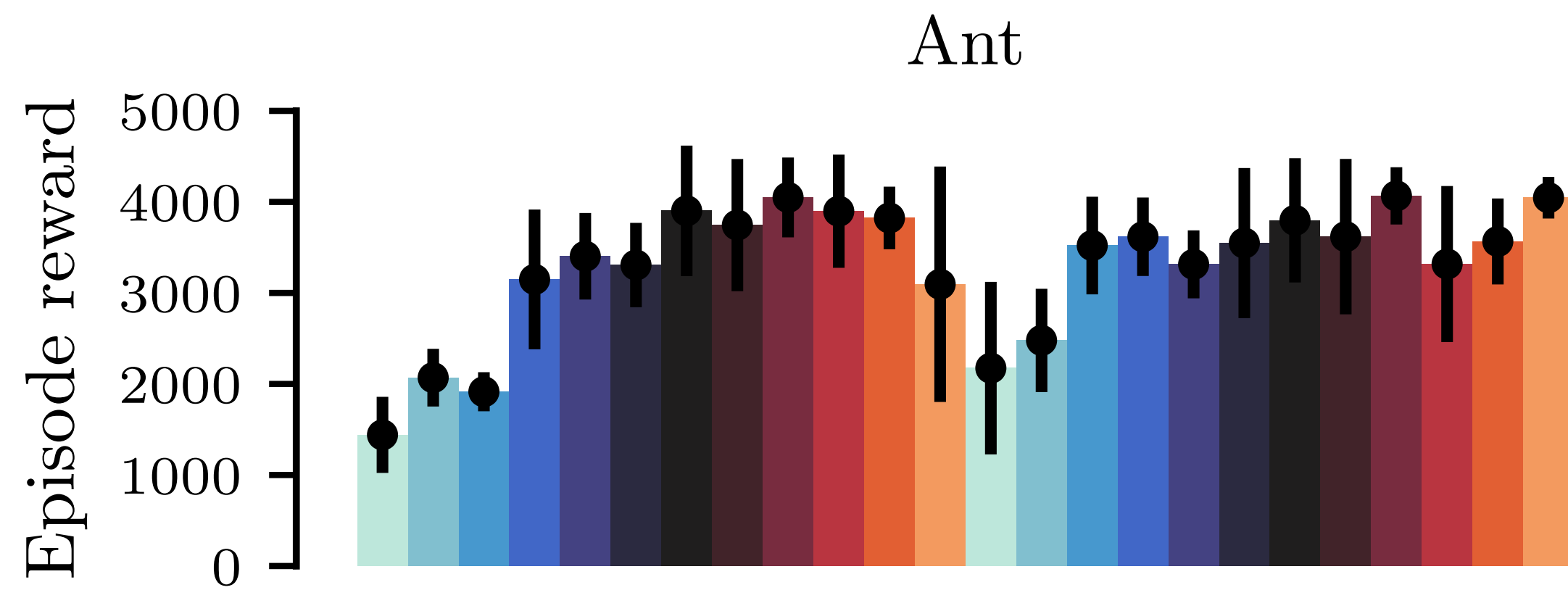
AdroitHandDoor

FrankaKitchen



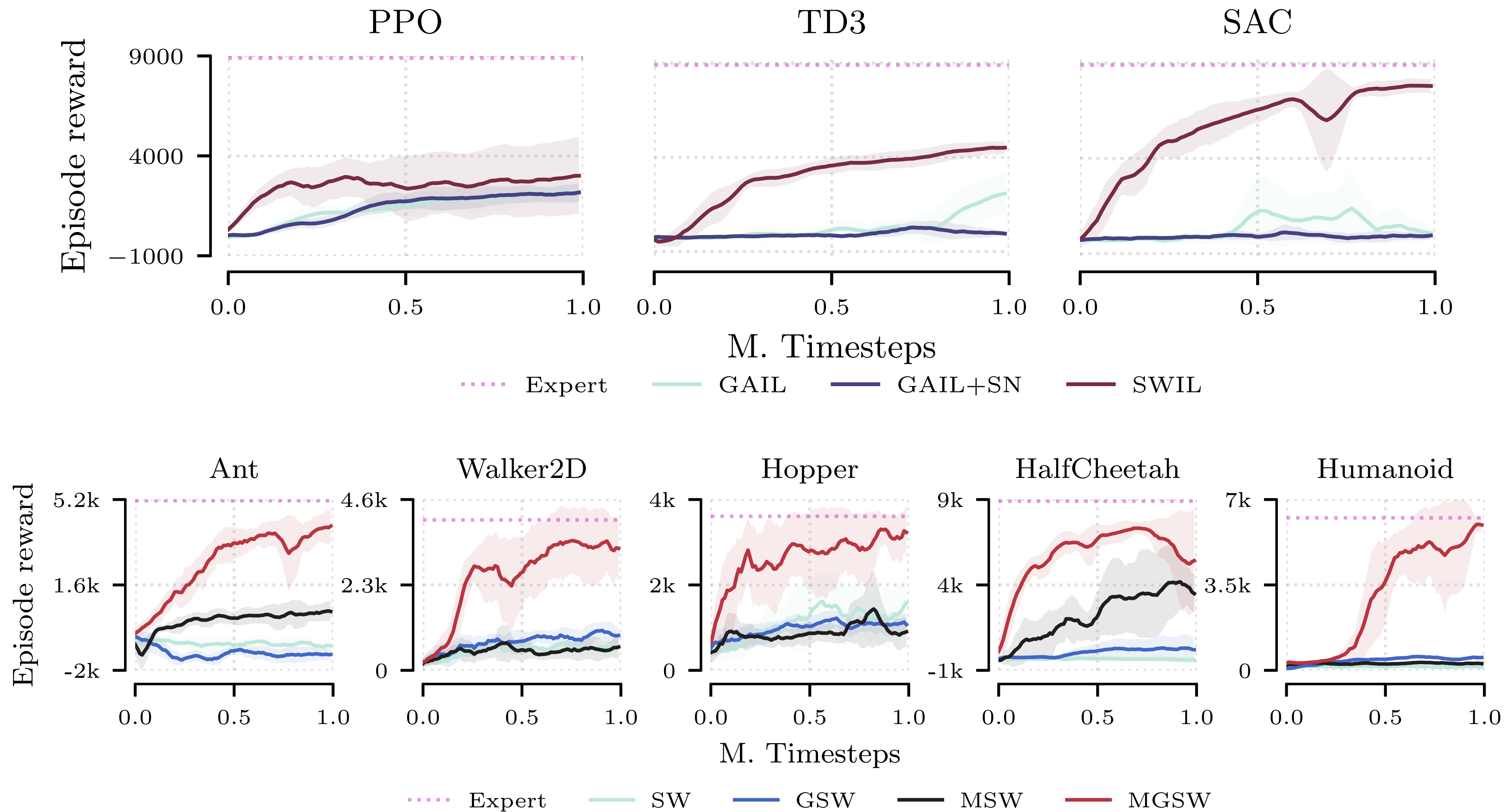
Expert GAIL AIRL SWIL

Ablation SWIL design choices



- | | | | | | |
|---------------|---------------|-----------------|-------------|-------------|---------------|
| nora_QL:1_SO | nora_QL:5_SA | nora_QL:10_SAS | ra_QL:1_SO | ra_QL:5_SA | ra_QL:10_SAS |
| nora_QL:5_SO | nora_QL:10_SA | nora_QL:1_SASD | ra_QL:5_SO | ra_QL:10_SA | ra_QL:1_SASD |
| nora_QL:10_SO | nora_QL:1_SAS | nora_QL:5_SASD | ra_QL:10_SO | ra_QL:1_SAS | ra_QL:5_SASD |
| nora_QL:1_SA | nora_QL:5_SAS | nora_QL:10_SASD | ra_QL:1_SA | ra_QL:5_SAS | ra_QL:10_SASD |

SWIL: comp RL / slicing ablation



SWIL: Algorithm

Algorithm 1 Sliced Wasserstein Imitation Learning (SWIL)

- 1: **require** Reinforcement learning algorithm with a policy improvement step for π_θ .
 - 2: **input** Expert trajectories $(\xi_i)_{i=1}^N$, initial policy π_θ (and any other initial state needed by the reinforcement learning algorithm), initial \mathcal{MGSW}_2 -critic $g^{(\psi)} : X \rightarrow \mathbb{R}_+$, policy and expert replay buffers with batch size B .
 - 3: **while** below maximum number of iterations **do**
 - 4: Compute a rollout $(x_R^{(t)})_{t=1}^T$ using the policy π_θ , add each element to the policy replay buffer.
 - 5: Sample a minibatch of state-action pairs $(x_\pi^{(b)})_{b=1}^B$ and $(x_E^{(b)})_{b=1}^B$ from the policy and expert replay buffers and compute $\mathcal{G}_\pi = (g^{(\psi)}(x_\pi^{(b)}))_{b=1}^B$ and $\mathcal{G}_E = (g^{(\psi)}(x_E^{(b)}))_{b=1}^B$, which we implicitly interpret as empirical measures.
 - 6: For each t , replace the closest atom in \mathcal{G}_π with x_t , obtaining $\mathcal{G}_R(x_t) = \text{rpl}(\mathcal{G}_\pi, x_t)$.
 - 7: Update $g^{(\psi)}$ via gradient-based supervised learning using the cross-entropy loss $\mathcal{L} = \sum_{b=1}^B \log g^{(\psi)}(x_\pi^{(b)}) - \sum_{b=1}^B (1 - \log g^{(\psi)}(x_E^{(b)}))$ with respect to the sampled minibatches.
 - 8: Perform a reinforcement learning step using the rewards $r_{\mathcal{MGSW}_2} = W_2(\mathcal{G}_E, \mathcal{G}_\pi) - W_2(\mathcal{G}_E, \mathcal{G}_R(x_t))$ for each $t = 1, \dots, T$.
 - 9: **return** learned policy π_θ .
-

CI-IRL: Algorithm 1

Algorithm 1 CI regularized Feature Matching IRL (CI-FMIRL)

Input: Expert trajectories \mathcal{D}_E^e assumed to be obtained from multiple experts *by intervening on* $p(\xi|\psi, \varphi)$

Init: Initialize reward estimate r_ψ and state feature network φ_θ

for setting e in $\{1, \dots, \mathcal{E}_{tr}\}$ **do**

while $r_{\psi, \varphi}$ not converged **do**

 Compute feature matching gradient $\nabla_\psi \mathcal{L}(\psi, \varphi; e) = \mathbb{E}_{\mathcal{D}_E^e}[\varphi(\xi)] - \mathbb{E}_{p(\xi|\psi)}[\varphi(\xi)]$ and *causal invariance* penalty gradient $\nabla_\varphi \mathbb{D}(\psi, \varphi; e)$ and backpropagate the weighted sum through feature network $\varphi_\theta(s)$

 Compute policy $\pi_{r_{\psi, \varphi}}$ using value iteration on the reward estimate $r_{\psi, \varphi}$

end for

end for

Return: Trained reward $r_{\varphi, \psi}$

CI-IRL: Algorithm 2

Algorithm 2 CI regularized Adversarial IRL (CI-AIRL)

Input: Expert trajectories \mathcal{D}_E^e assumed to be obtained from multiple experts *by intervening on* $p(\xi|\psi, \varphi)$

Init: Initialize actor-critic $\pi_\theta, \nu_\vartheta$ and discriminator $g_{\xi, \varphi}$

for setting e in $\{1, \dots, \mathcal{E}_{tr}\}$ **do**

 Collect trajectory buffer $\mathcal{D}_\pi = \{\xi_i\}_{i \leq |\mathcal{D}_\pi|}$ by executing the policy π_θ

 Update $g_{\varphi, \theta}(s, a)$ via binary logistic regression by maximizing

$$\mathcal{L}(\varphi, \psi; e) = \mathcal{L}_{\text{BCE}}(\xi, \varphi, \psi; e) + \lambda \|\nabla_{\psi|_{\psi=1.0}} \mathcal{L}_{\text{BCE}}(\xi, \varphi, \psi; e)\|^2$$

 using dataset tuple $(\mathcal{D}_E^e, \mathcal{D}_\pi)$

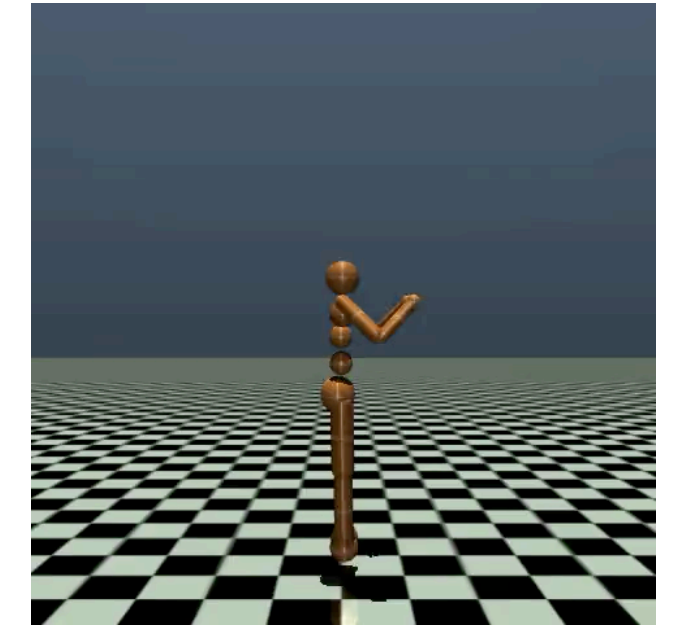
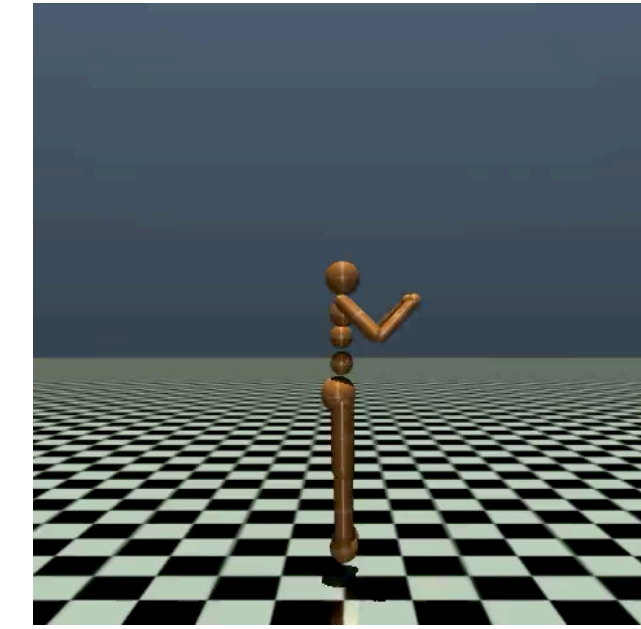
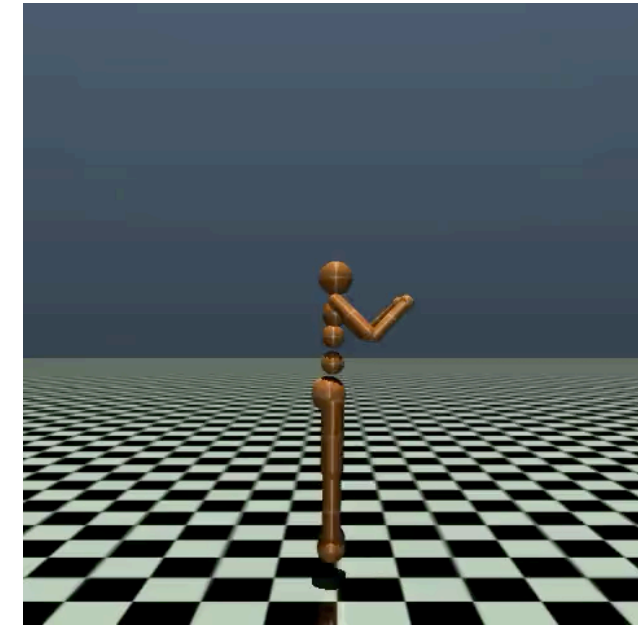
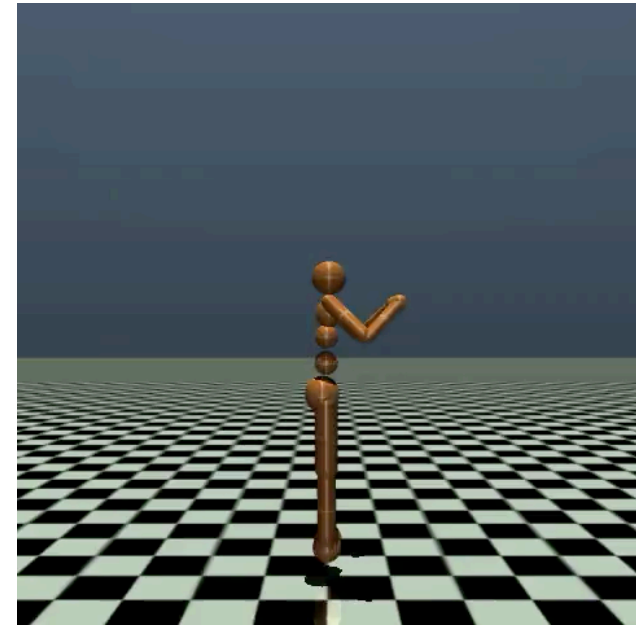
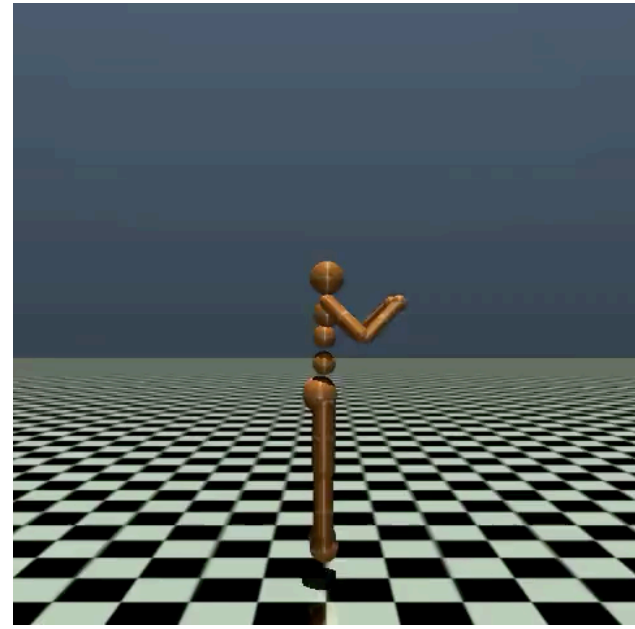
 Update actor-critic $(\pi_\theta, \nu_\vartheta)$ w.r.t. the reward function of the *regularized discriminator* using the *soft-actor-critic* RL procedure

end for

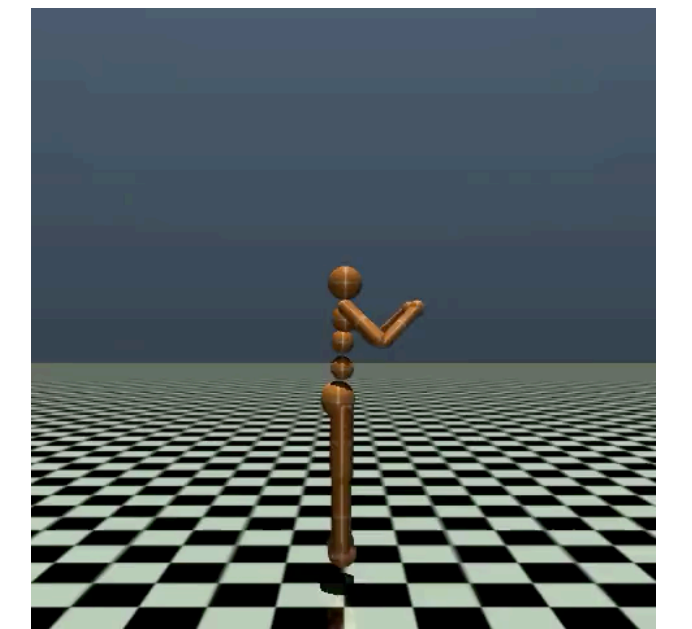
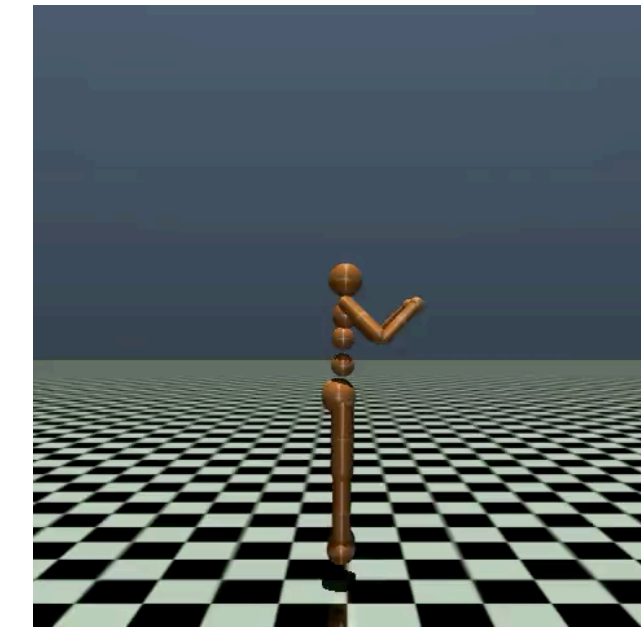
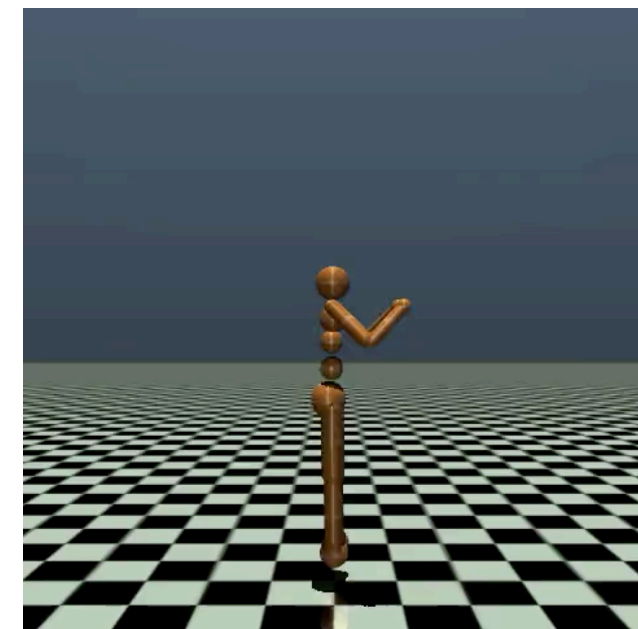
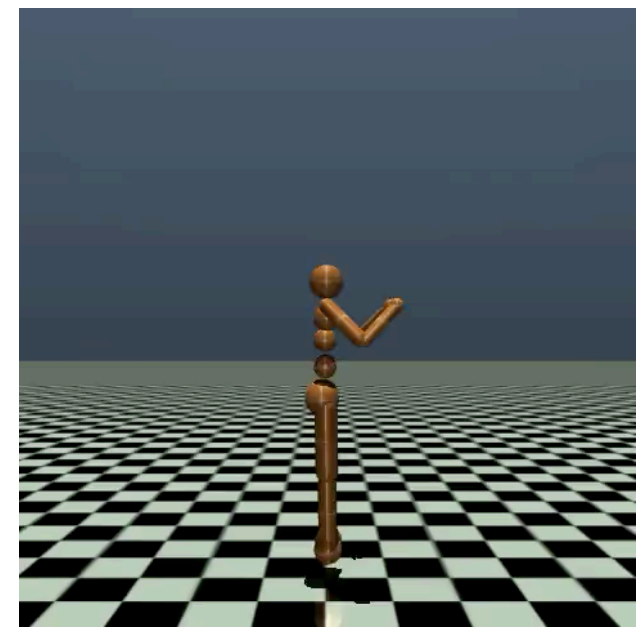
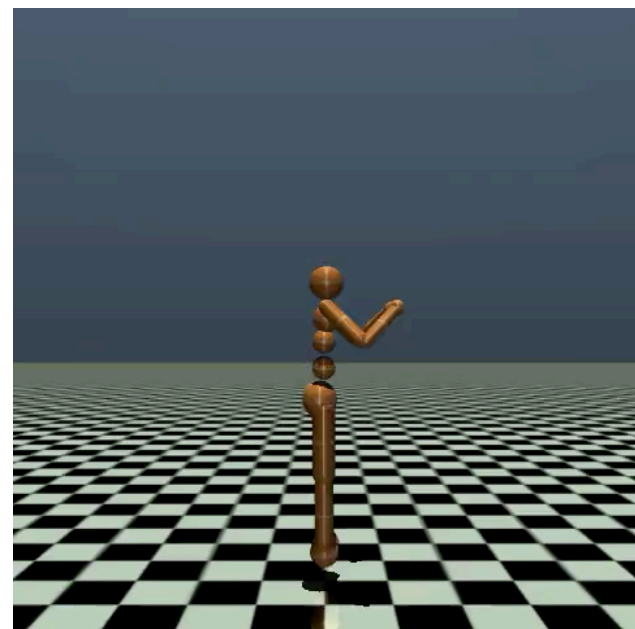
Return: Trained reward $r_{\varphi, \psi}$ and actor-critic $\pi_\theta, \nu_\vartheta$

More silly walks

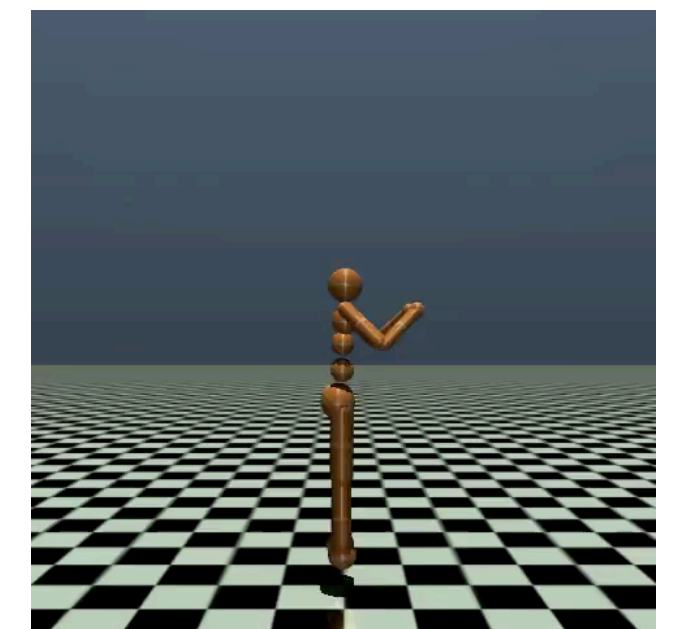
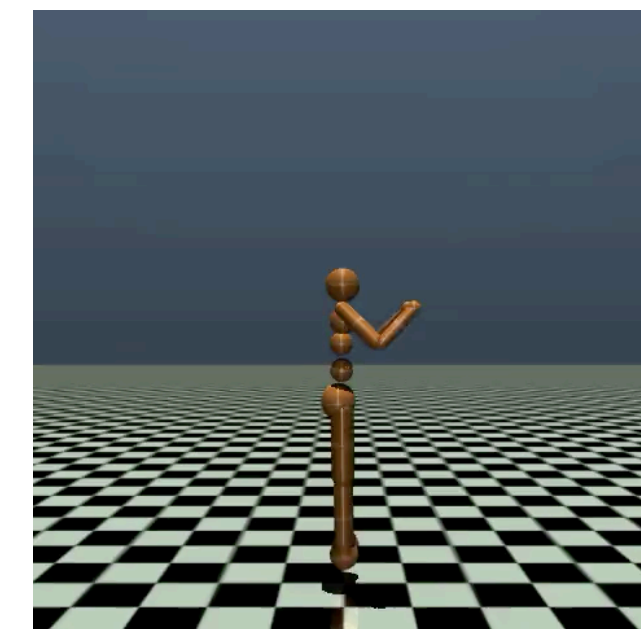
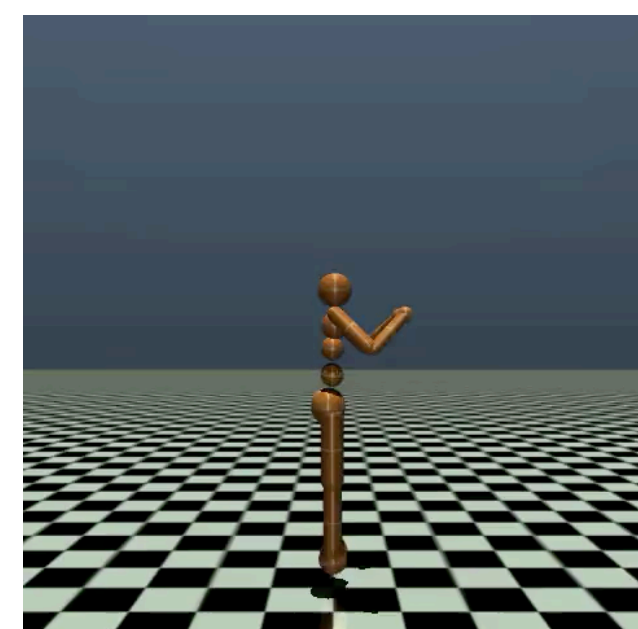
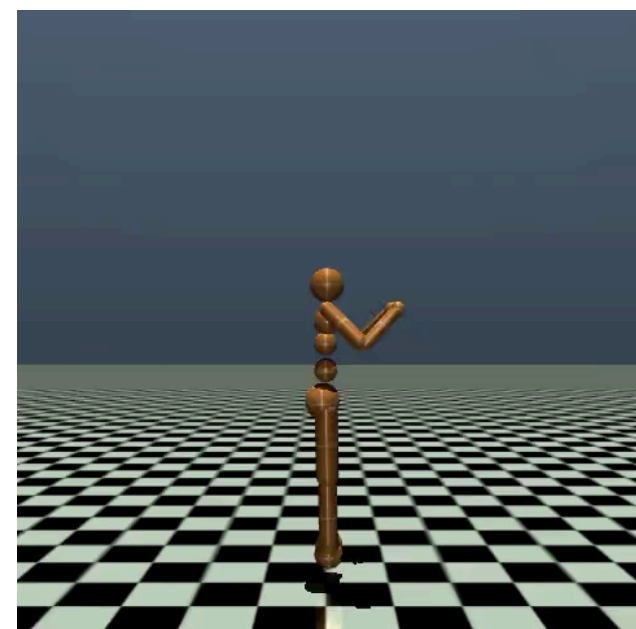
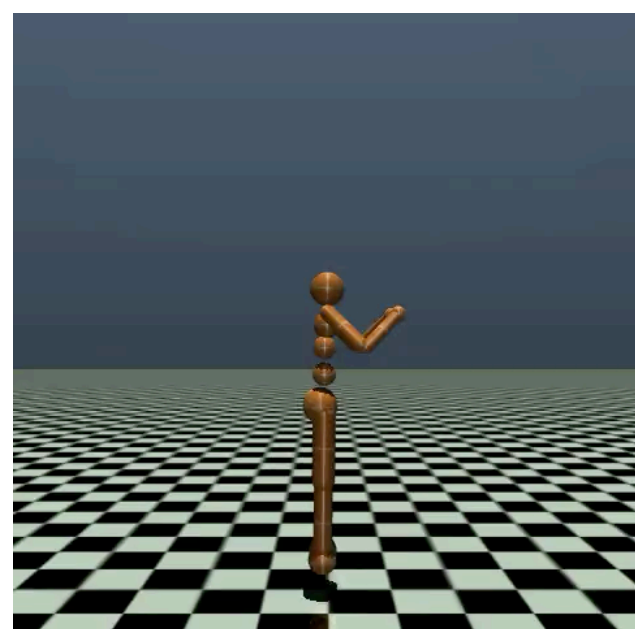
ERM



Lip



CI



Posterior Agreement in Soft Actor Critic

