

Cross-Modal Domain Adaptation for Cost-Efficient Visual Reinforcement Learning

Xiong-Hui Chen*, **Shengyi Jiang***, **Feng Xu**, **Zongzhang Zhang[†]**, **Yang Yu**

National Key Laboratory of Novel Software Technology

Nanjing University, Nanjing 210023, China

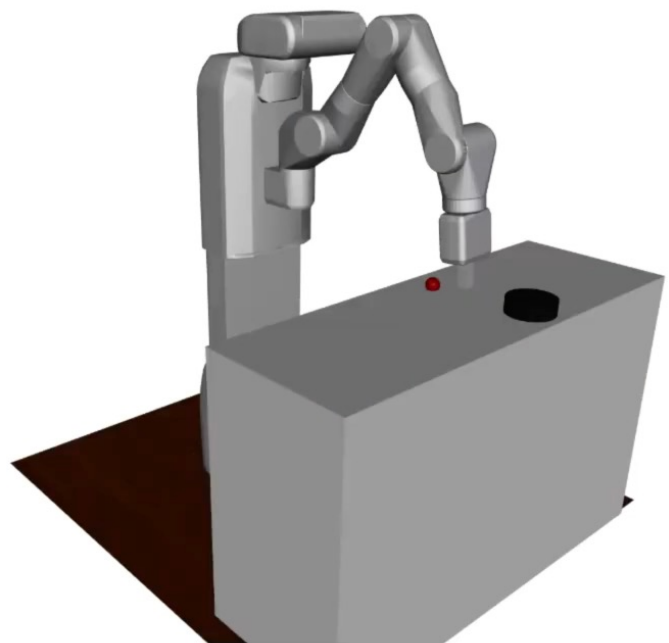
Pazhou Lab, Guangzhou 510330, China

{chenxh, jiangsy, xufeng}@lamda.nju.edu.cn, {zzzhang, yuy}@nju.edu.cn

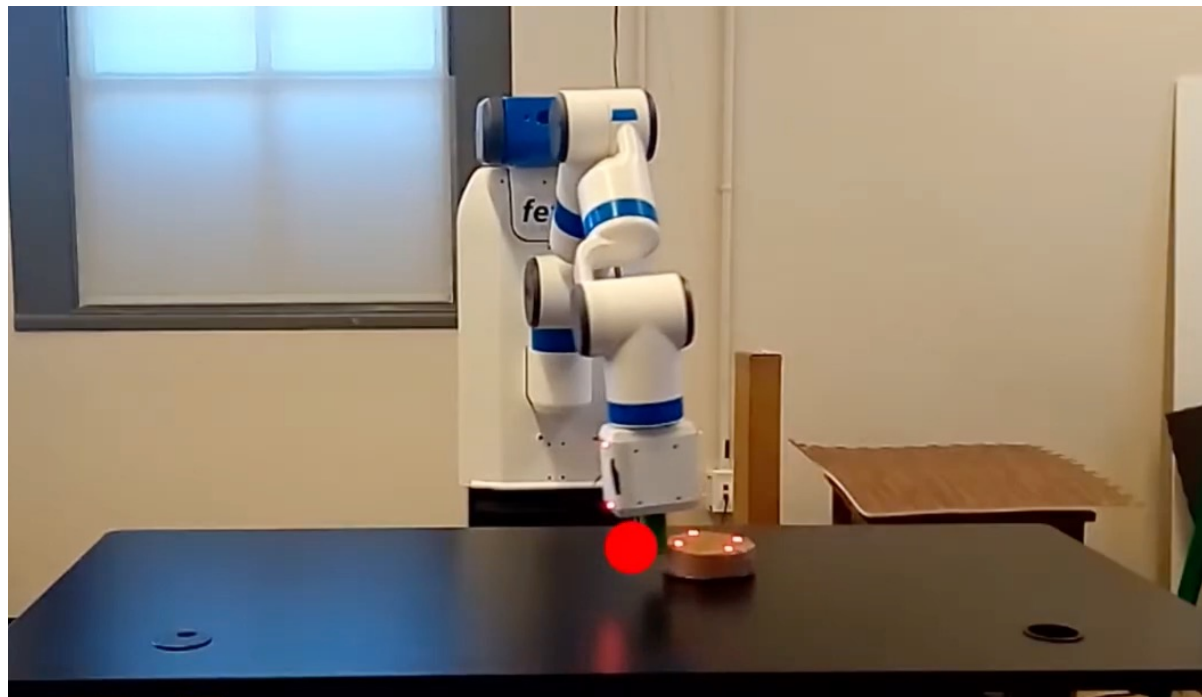
Table of Contents

1. **Background and Motivation**
2. Cross-Modal Domain Adaptation with Sequential structure (CODAS)
3. Experiment
4. Take-home Messages

An example of Sim2Real Reinforcement Learning



Simulation



Real world

An example of reality-gap which is the core challenge in Sim2Real RL

The Framework of Unsupervised Domain Adaptation in Sim2Real RL

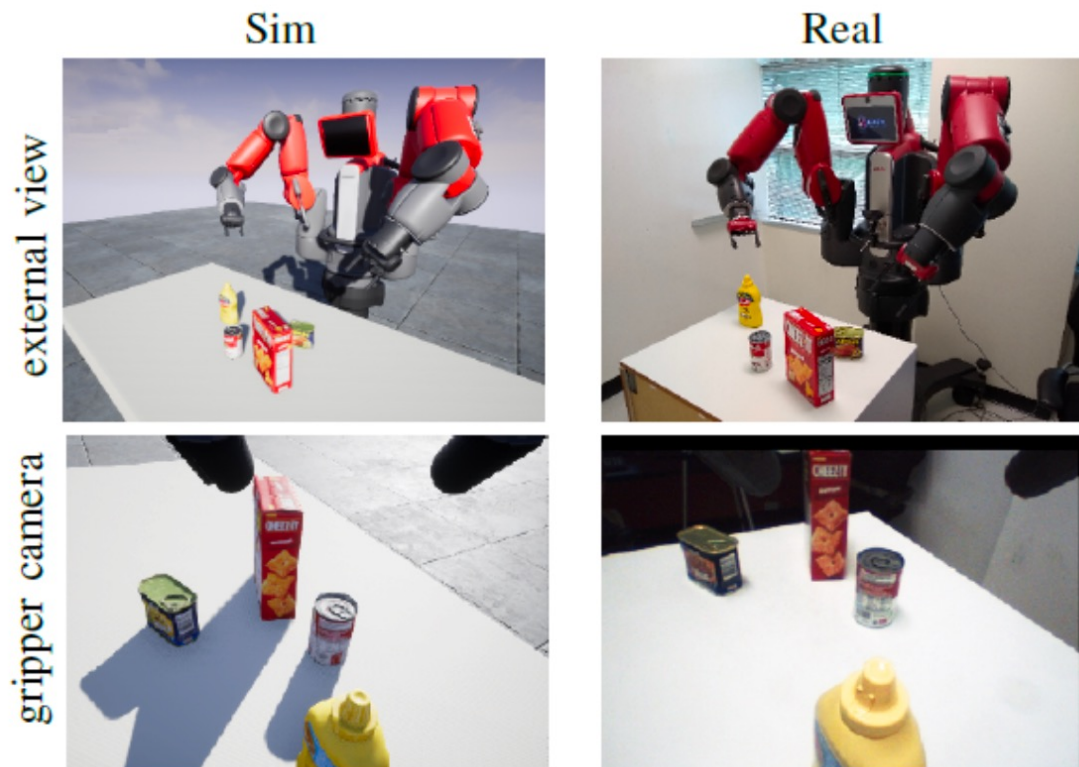


Fig. An example of reality-gaps in observation-space in Sim2Real RL [1]

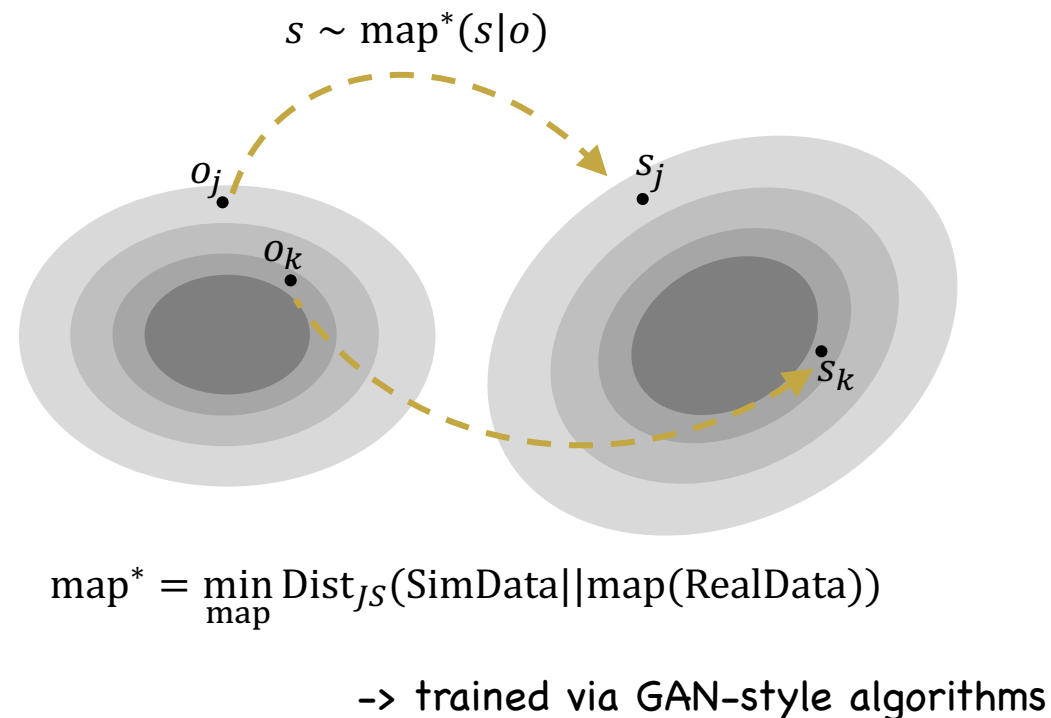


Fig. The framework of Unsupervised Domain Adaptation for Sim2Real RL

Unsupervised domain adaptation (UDA) learns a mapping function to align the data distribution of the source and the target domain to handle the challenge of reality-gap on observation-space in Sim2Real RL

Cross-Modal UDA: A cost-efficient Framework for Sim2Real RL



Cross-Modal UDA: A cost-efficient Framework for Sim2Real RL



Image-to-image UDA introduce three extra cost, which is ignored in discussion in previous work.

1. human labor of building a visual simulator
2. huge computation resource required by running the simulator
3. inferior policy training on visual simulator.

Cross-Modal UDA: A cost-efficient Framework for Sim2Real RL

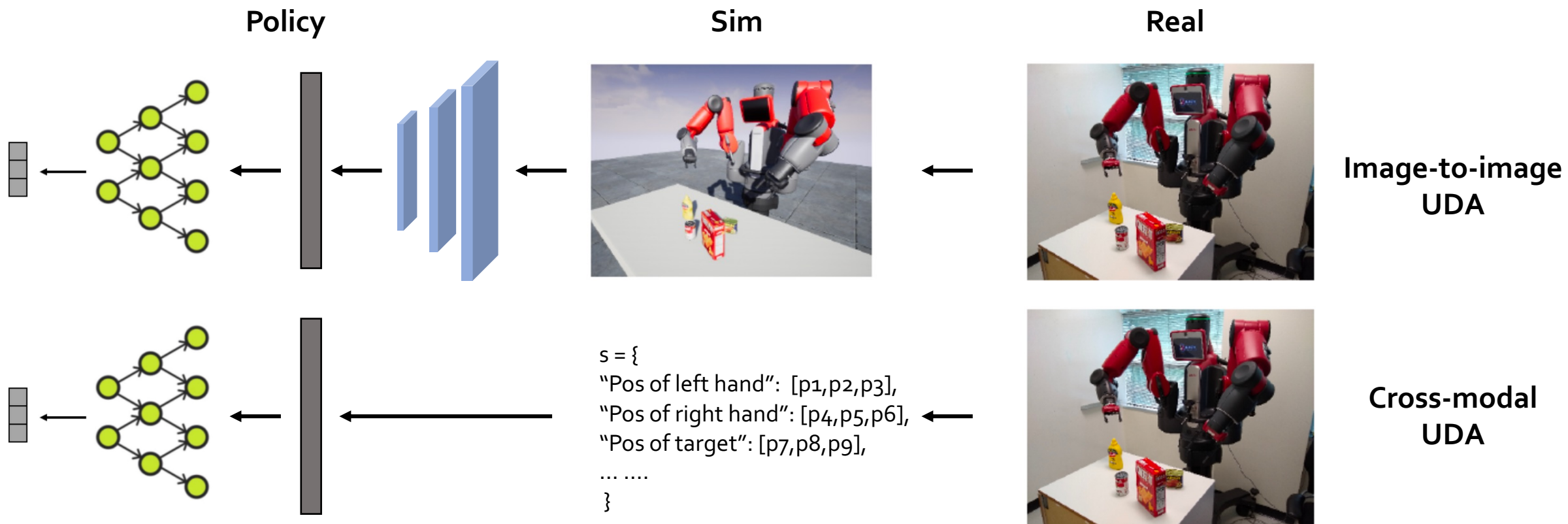


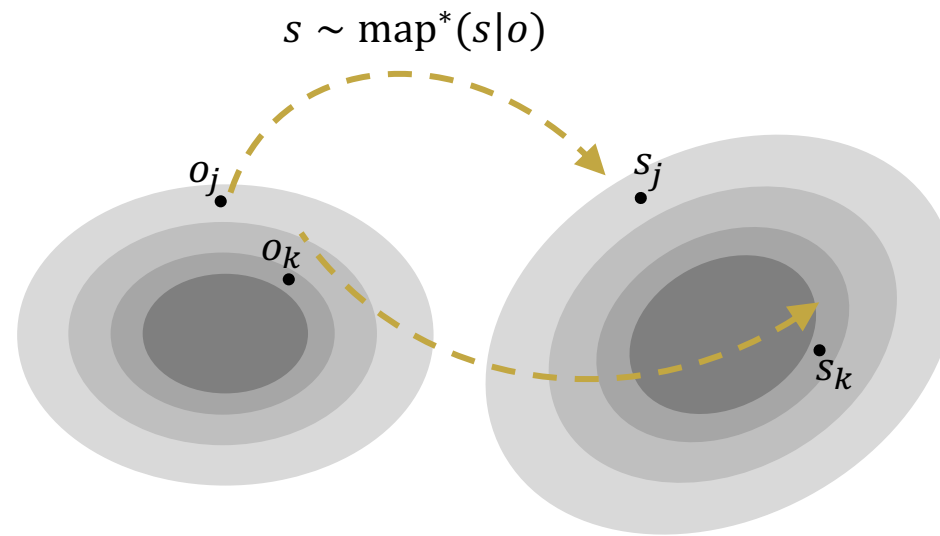
Image-to-image UDA introduce three extra cost, which is ignored in discussion in previous work.

1. human labor of building a visual simulator
2. huge computation resource required by running the simulator
3. inferior policy training on visual simulator.



Can be solved in Cross-modal UDA

Ill-posedness of the raw objective in current UDA solutions



$$\text{map}^* = \min_{\text{map}} \text{Dist}_{JS}(\text{SimData} || \text{map}(\text{RealData}))$$

-> trained via GAN-style algorithms

We cannot adopt the previous Unsupervised Domain Adaptation to Cross-Modal UDA setting.

Ill-posedness of the raw objective in current UDA solutions

$$\text{map}^* = \min_{\text{map}} \text{Dist}_{JS}(\text{SimData} || \text{map}(\text{RealData}))$$

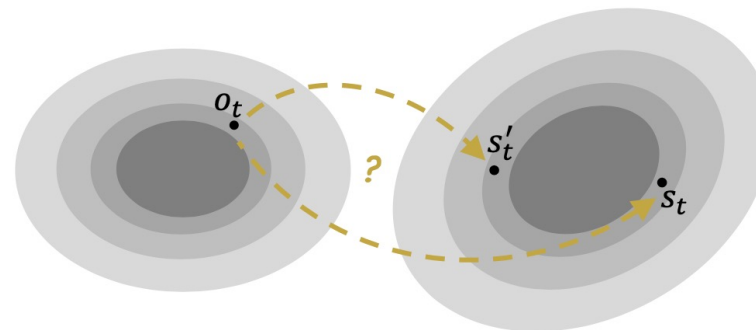
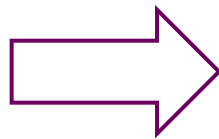


Fig. An example of ill-posedness of the distribution minimization objective in current UDA algorithms

Since s_t and s'_t have similar probabilities, mapping an instance o_t to anywhere of a similar probability in the source domain is “reasonable” if we only consider distribution matching.

In image-to-image UDA, current methods rely on additional constraints on modality consistency to handle the problem implicitly

- special model structure [2], e.g. U-Net, Cycle-GAN;
- auxiliary losses, e.g. geometry consistency [3];

However, these constraints cannot hold anymore in Cross-modal UDA setting.

Ill-posedness of the raw objective in current UDA solutions

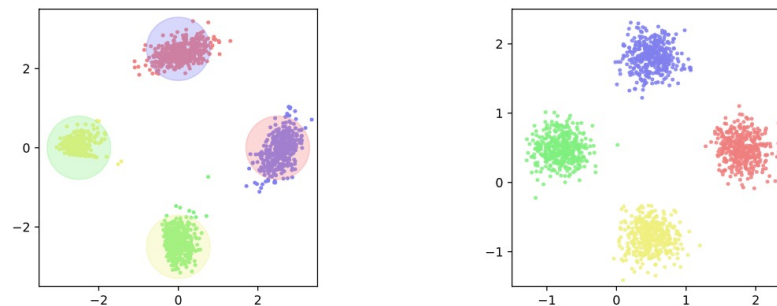
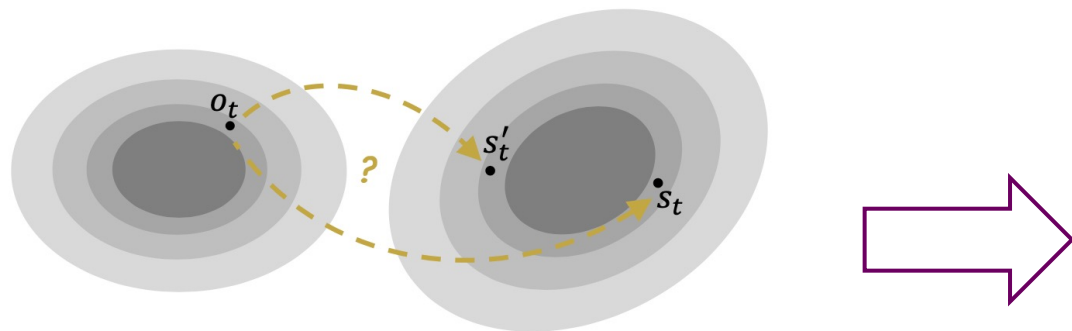


Fig. An example of ill-posedness of the distribution minimization objective in current UDA algorithms

Fig. A toy example of ill-posed UDA

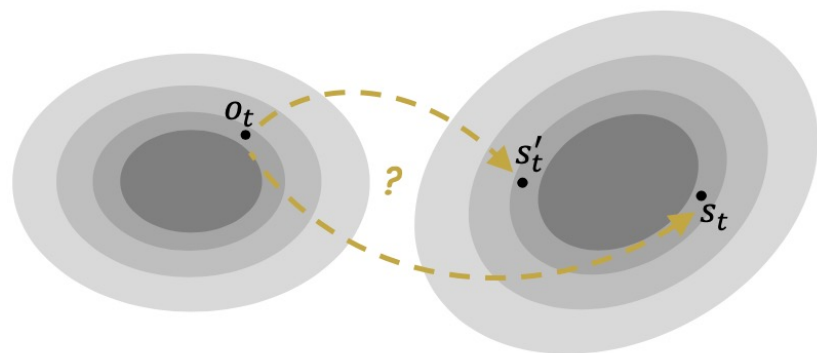
Since s_t and s'_t have similar probabilities, mapping an instance o_t to anywhere of a similar probability in the source domain is “reasonable” if we only consider distribution matching.

Our research question: Can we handle the ill-posedness of the objective directly?

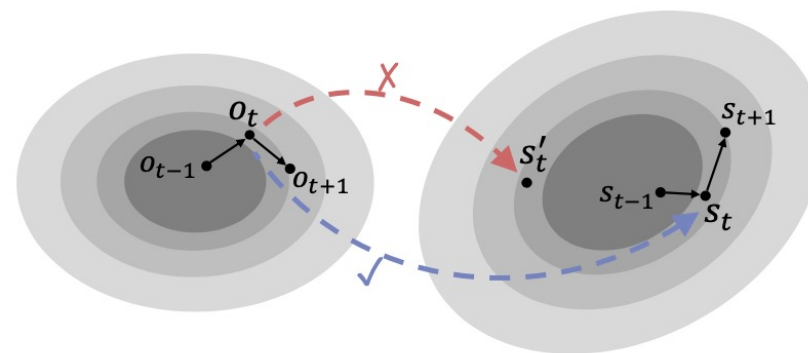
Table of Contents

1. Background and Motivation
2. Cross-Modal Domain Adaptation with Sequential structure (CODAS)
3. Experiment
4. Take-home Messages

Any other potential way to handle the ill-posedness of the objective?



(a) Mapping only considering state-distribution matching

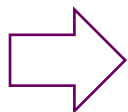
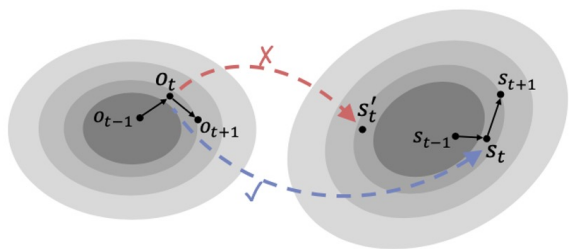


(b) Mapping with sequential structure

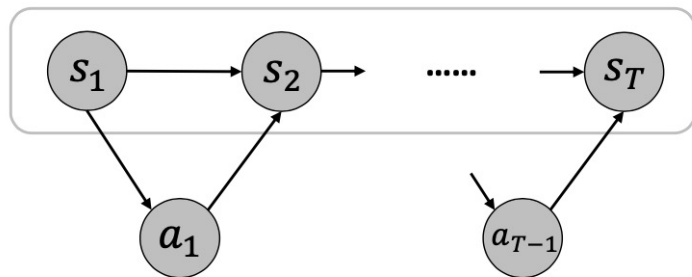
If we can make use of the **sequential structure in the Markov Decision Processes**, the historical information will give us the ability to identify the difference between s'_t and s_t ,

then the proposed ill-posedness of the objective will be fixed.

Reformulate the objective of UDA in RL based on the framework of variational inference

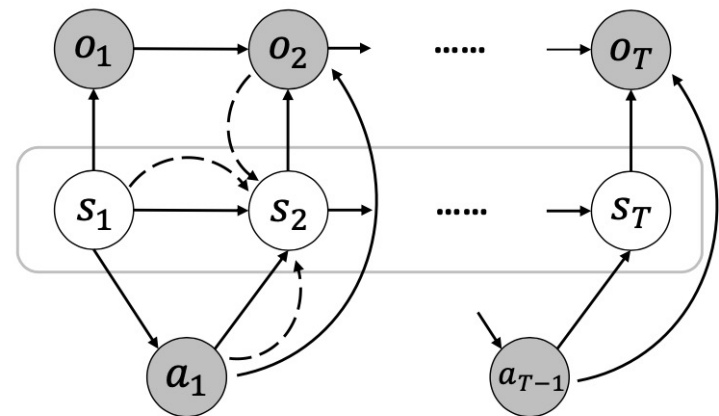


Shaded nodes: Observable
White nodes: Unobservable



(a) Generation process in the source domain

Infer



(b) Generation and inference process in the target domain

Fig. The generation and inference process of UDA based on the framework of variational inference

Reformulate the objective of UDA in RL based on the framework of variational inference

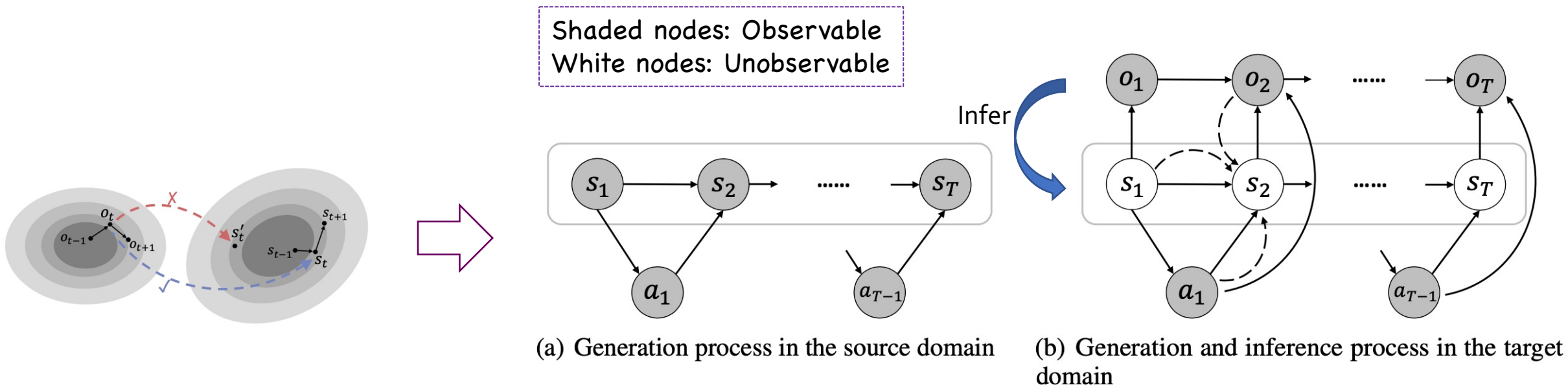


Fig. The generation and inference process of UDA based on the framework of variational inference

$$\min_{\phi} \mathbb{E}_{\tau^o} [D_{KL} [q_{\phi}(\tau^s | \tau^o) || p(\tau^s | \tau^o)]]$$

↓ (ELBO)

$$\max_{\phi, \theta} \mathbb{E}_{\tau^o} \left[\underbrace{\mathbb{E}_{\hat{\tau}^s \sim q_{\phi}(\tau^s | \tau^o)} [\log p_{\theta}(\tau^o | \hat{\tau}^s)]}_{\text{Reconstruction Error}} - \underbrace{D_{KL} [q_{\phi}(\tau^s | \tau^o) || p(\tau^s)]}_{\text{Divergence between prior and inferred state distributions}} \right]$$

D_{KL} : KL Divergence
 τ : trajectory
 o : observation
 s : state
 ϕ : parameters to learn

Reformulate the objective of UDA in RL based on the framework of variational inference

$$\min_{\phi} \mathbb{E}_{\tau^o} [D_{\text{KL}} [q_{\phi}(\tau^s | \tau^o) || p(\tau^s | \tau^o)]]$$

↓ (ELBO)

$$\max_{\phi, \theta} \mathbb{E}_{\tau^o} \left[\underbrace{\mathbb{E}_{\hat{\tau}^s \sim q_{\phi}(\tau^s | \tau^o)} [\log p_{\theta}(\tau^o | \hat{\tau}^s)]}_{\text{Reconstruction Error}} - \underbrace{D_{\text{KL}} [q_{\phi}(\tau^s | \tau^o) || p(\tau^s)]}_{\text{Divergence between prior and inferred state distributions}} \right]$$

Reconstruction Error

Divergence between prior and inferred state distributions

D_{KL} : KL Divergence

τ : trajectory

o : observation

s : state

ϕ : parameters to learn



$$\max_{\phi, \theta} \mathbb{E}_{\tau^o \sim \mathcal{D}^o} \left[\sum_{t=1}^T \mathbb{E}_{\hat{s}_t \sim q_{\phi}(s_t | \hat{s}_{t-1}, a_{t-1}, o_t)} \left[\log p_{\theta}(o_t | \hat{s}_t, o_{t-1}, a_{t-1}) \right. \right. \right]$$

Reconstruction loss

$$\left. - \lambda_D \log \left(1 - D_{\omega^*}(\hat{s}_t, a_t, h_{t-1}) \right) \right],$$

Trajectory-distribution mismatch loss

$$\text{s.t. } \omega^* = \arg \max_{\omega} \mathbb{E}_{\tau^s \sim \mathcal{D}^s} \left[\sum_{t=1}^T \log D_{\omega}(s_t, a_t, h_{t-1}) \right]$$

Discriminator loss

$$+ \mathbb{E}_{\tau^o \sim \mathcal{D}^o} \left[\sum_{t=1}^T \mathbb{E}_{\hat{s}_t \sim q_{\phi}(s_t | \hat{s}_{t-1}, a_{t-1}, o_t)} \log \left(1 - D_{\omega}(\hat{s}_t, a_t, h_{t-1}) \right) \right];$$

Embedded Dynamics Model for Stable Training

Inference Function only outputs a small Δs_t . The main part is from **Embedded DM** (p_ϕ).

$$\hat{s}_t = p_\phi(s_{t-1}, a_{t-1}) + \alpha \Delta s_t,$$

where $\Delta s_t \sim q_\phi(\Delta s \mid s_{t-1}, a_{t-1}, o_t)$

The parameters of **Embedded DM** are copied from a DM trained by:

$$\min_{\phi} \mathbb{E}_{(s,a,s') \sim D^s \cup D^s} [(p_\phi(s, a) - s')^2]$$

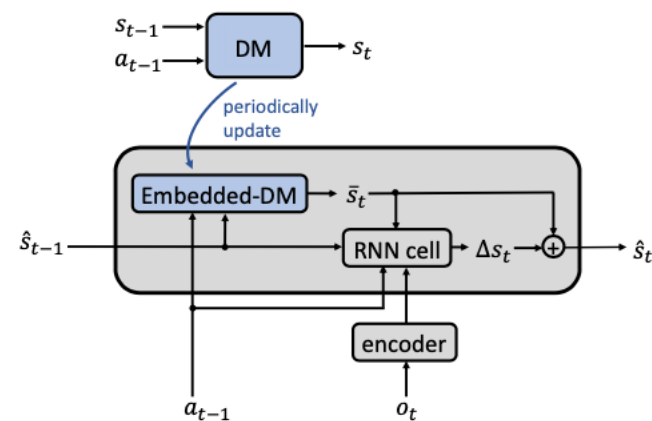
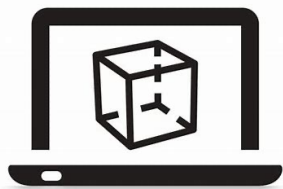


Figure: Detailed Structure of the Inference Function with Embedded DM

Overall Model Structure

(s_{t-1}, a_{t-1}, s_t)



A simulator of the source domain

(o_{t-1}, a_{t-1}, o_t)



Pre-collected dataset in the target domain

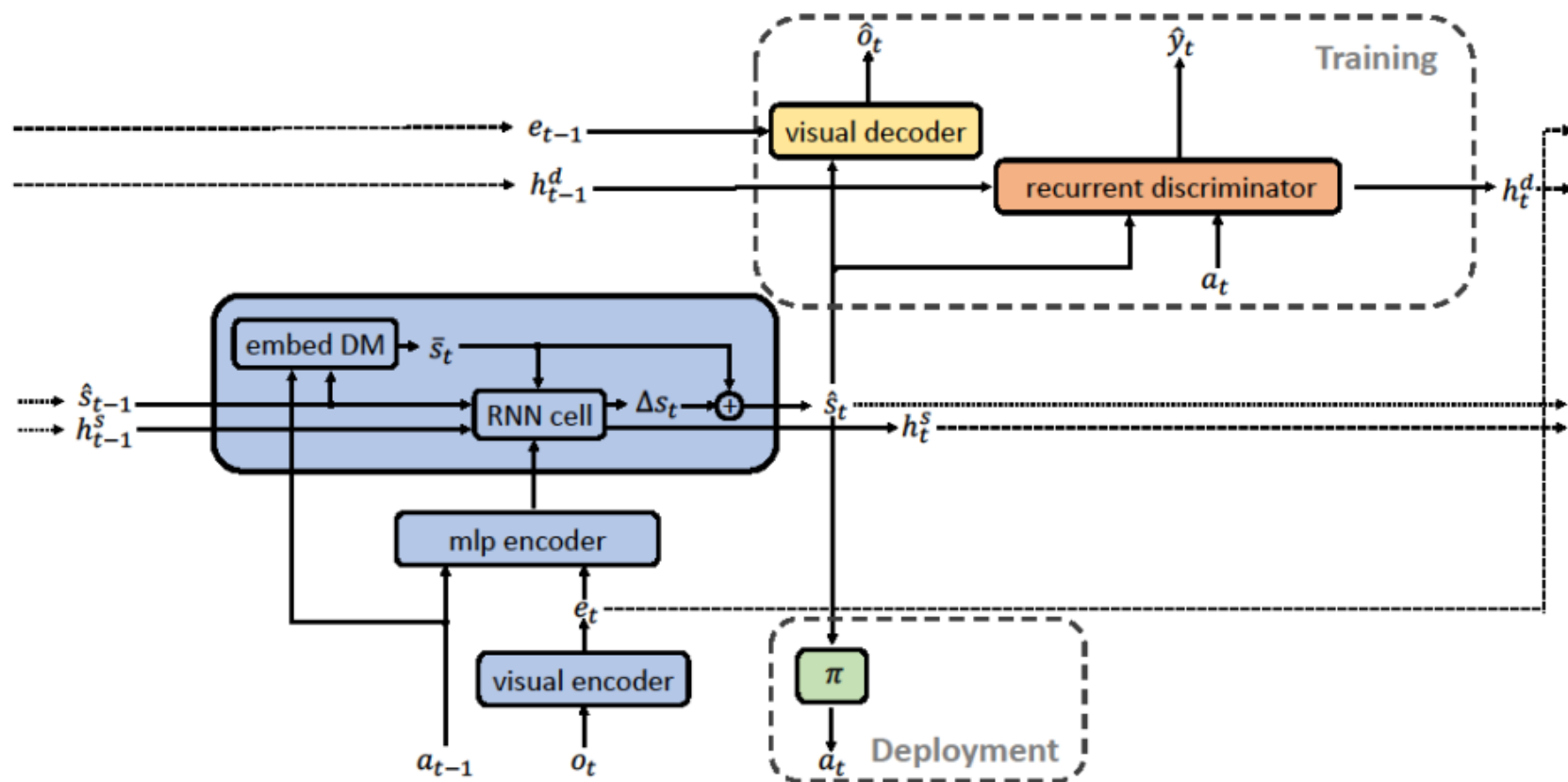
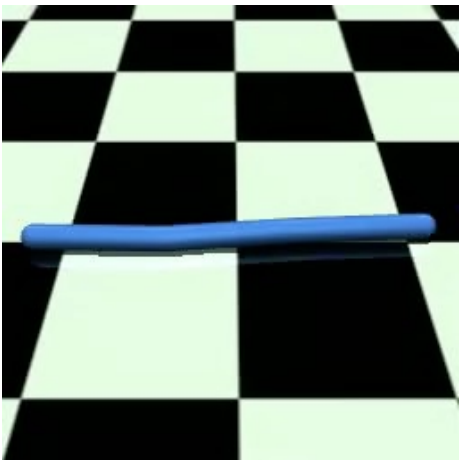


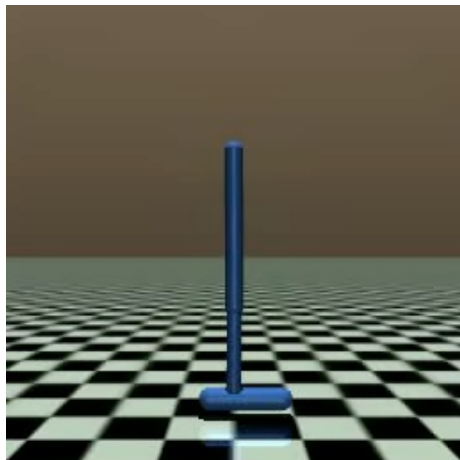
Table of Contents

1. Background and Motivation
2. Cross-Modal Domain Adaptation with Sequential structure (CODAS)
- 3. Experiment**
4. Take-home Messages

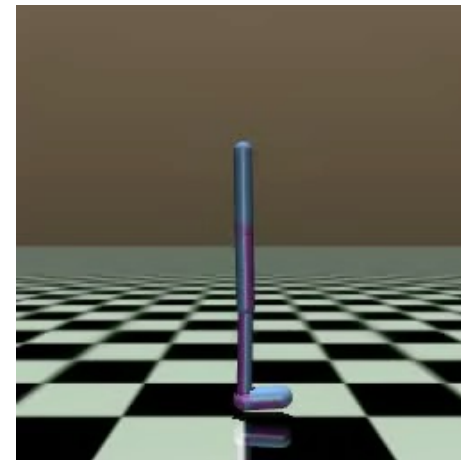
Performance on MuJoCo Tasks



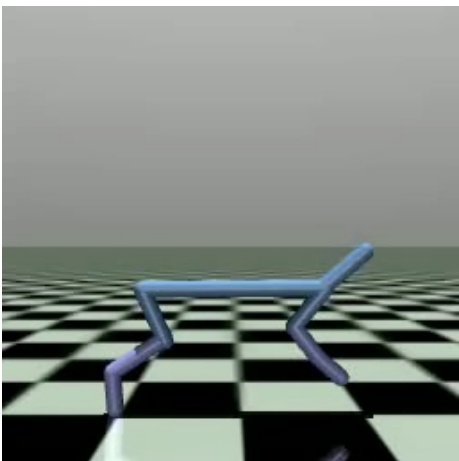
Swimmer



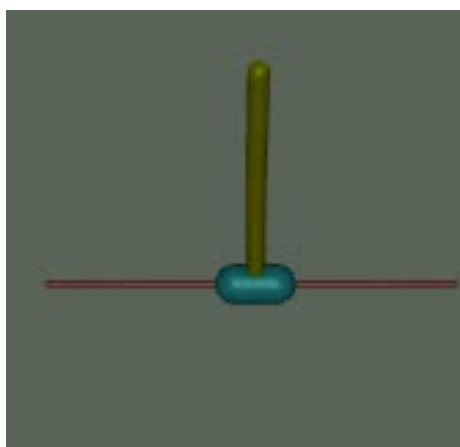
Hopper



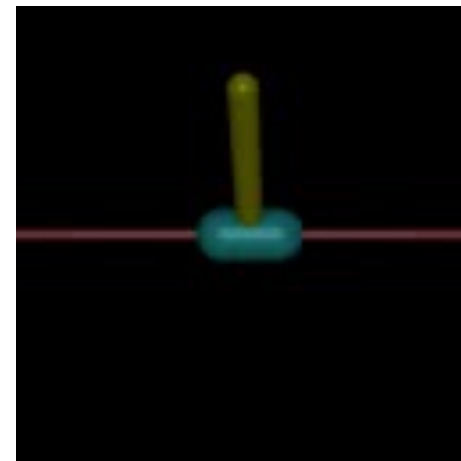
Walker2d



Half Cheetah



Inverted Double
Pendulum



Inverted Pendulum

Comparative Evaluation in MuJoCo Tasks

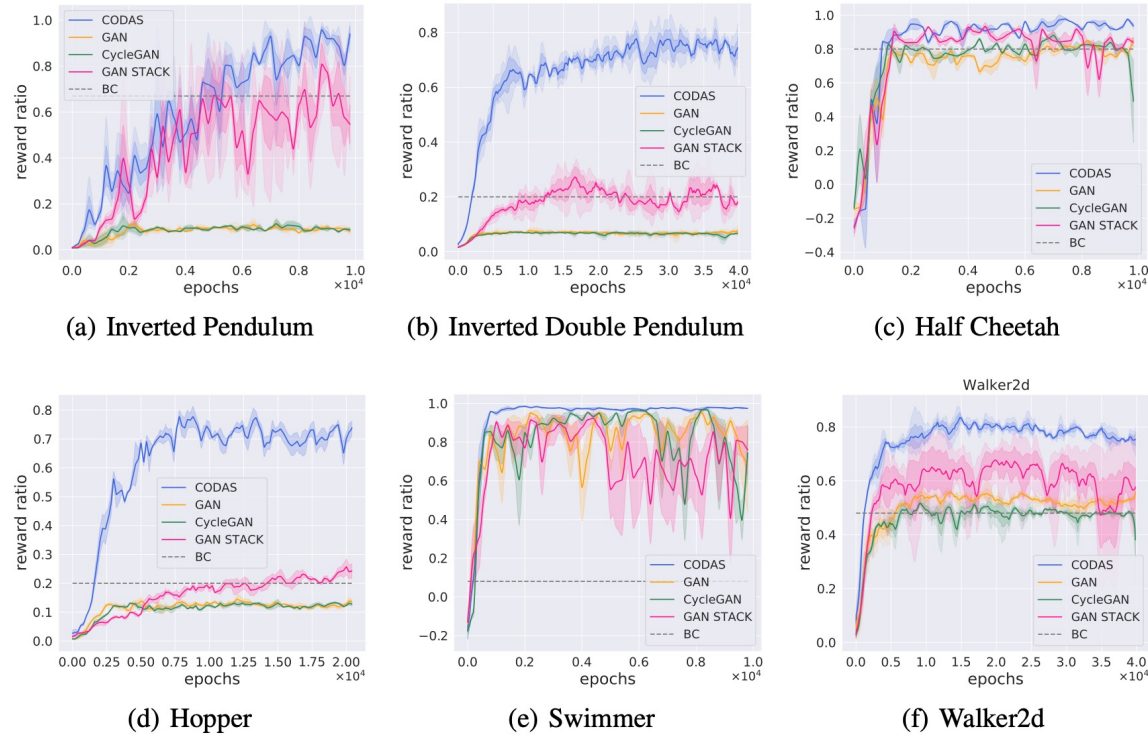


Figure 4: Training curves of different methods on MuJoCo.

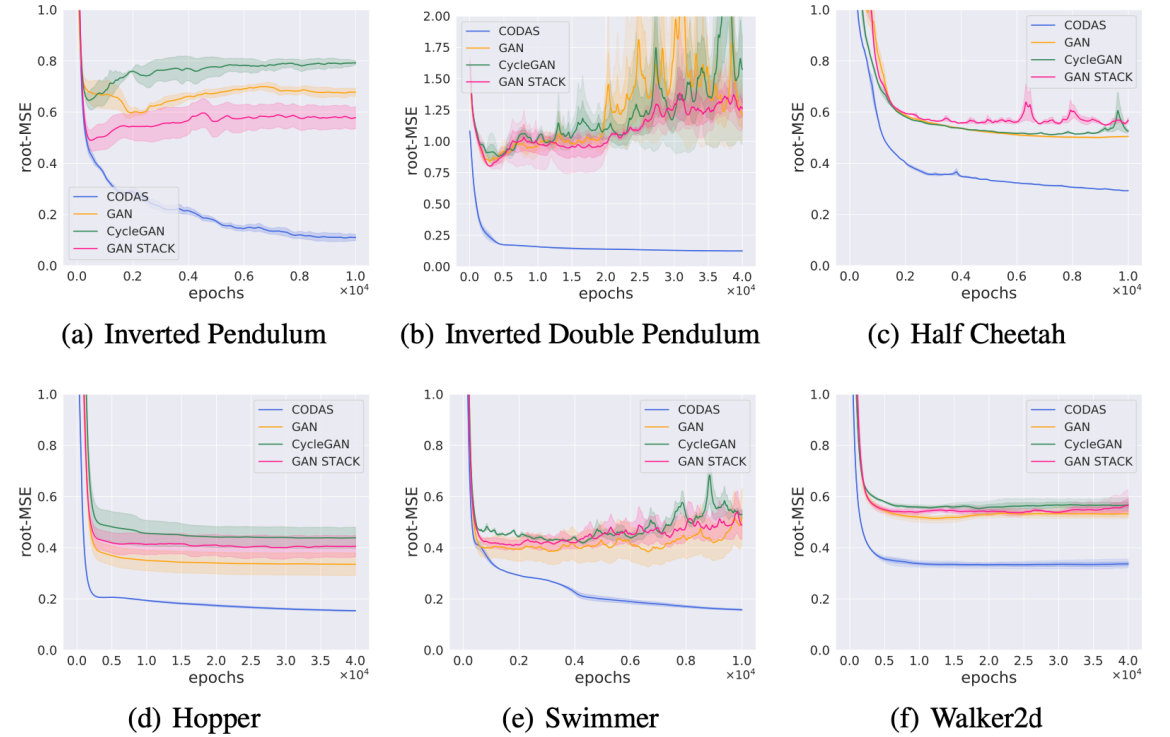


Figure 6: Root mean squared error between mapped states and ground-truth states. The solid lines denote the mean value. The shadows denote the standard deviation.

For all of the tasks, CODAS can map the correct states (i.e., with the smallest MSE-loss to the oracle states) and the performance of the deployment policies reach reasonable performance (75%~100%)

Visualization of the Learned Mapping on MuJoCo Tasks

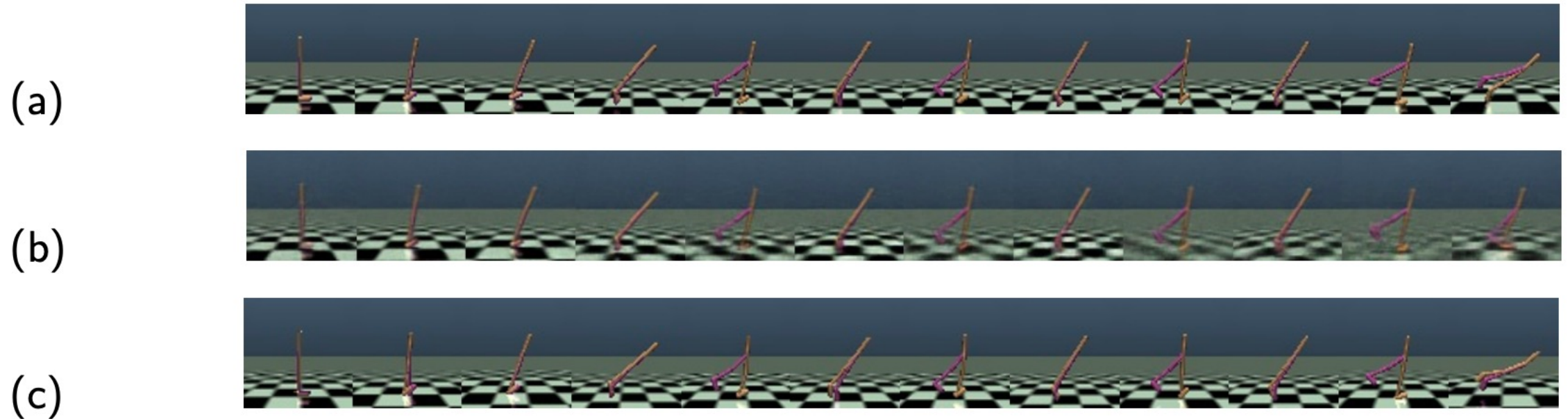


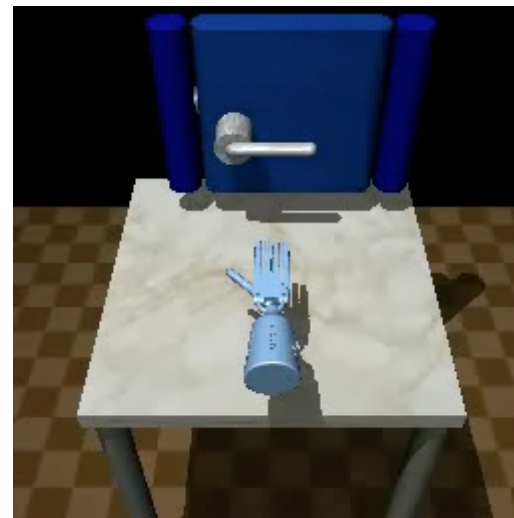
Fig. A visual illustration of (a) original images, (b) reconstructed images, and (c) re-rendered images of the mapped states in Hopper.

Both reconstructed images and re-rendered images match the original ones well. Re-rendered images can even match the original ones well in the last falling frames which are sparse in the dataset.

Performance on Robot Hand Manipulation Tasks



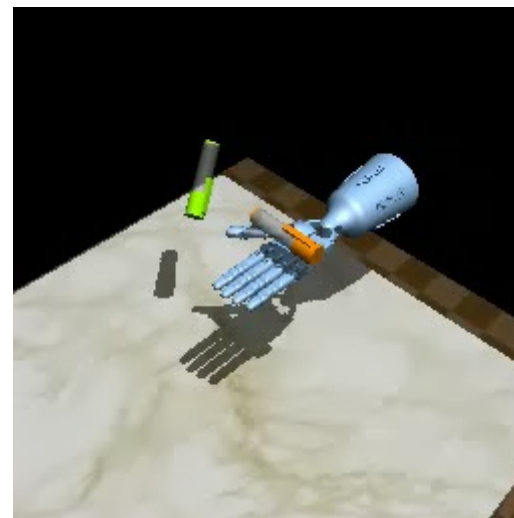
Relocate



Door



Hammer



Pen

Performance on Robot Hand Manipulation Tasks

Data collecting policy



Re-rendered video



Reconstructed video



| Tasks | hammer | pen | door | relocate |
|--------------|--------|-------|-------|----------|
| Reward Ratio | 0.820 | 0.701 | 0.886 | 0.090 |

In three out of four tasks, CODAS yields reasonable mapping functions for policy deployment.

Take-home Message

KEY point:

The Formulation of Variational Inference which considered the sequential structure in MDP can handle the ill-posedness of the objective and solve the UDA problem without relying on the knowledge of modality consistency.

Future work:

1. CODAS solve the UDA problem in a general way, it can be adopted to image-to-image UDA in theory and the practical adoption can be tried.
2. In the current formulation, we assume the policies/dynamics in the source and target domains are the same, which might not hold in real-world applications. By modeling the mismatching of dynamics models and data-collected policies into the CODAS framework, we can build a more practical UDA algorithm.

>> Thanks