

Model-based Meta Reinforcement Learning using Graph Structured Surrogate Models and Amortized Policy Search

Qi Wang & Herke van Hoof

Amsterdam Machine Learning Lab, UvA



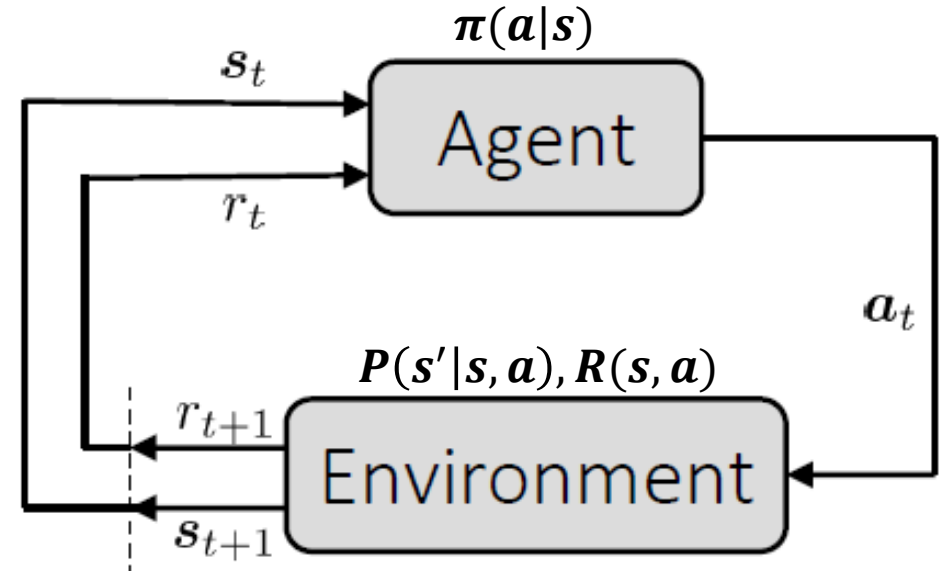
UNIVERSITY OF AMSTERDAM



Research Motivations

Sequential Decision-Making in MDP =.=

- **Learn policies from Feedbacks** → RL to solve MDPs
- **Crucial elements in MDPs** $M = \langle S, A, P, R, \gamma \rangle$:
 - ✓ State Space $s \in S$
 - ✓ Action Space $a \in A$
 - ✓ Policy Function $\pi(a|s)$
 - ✓ Transition Dynamics $P(s'|s, a)$
 - ✓ Reward Function $R(s, a)$
 - ✓ Step Reward Discount γ

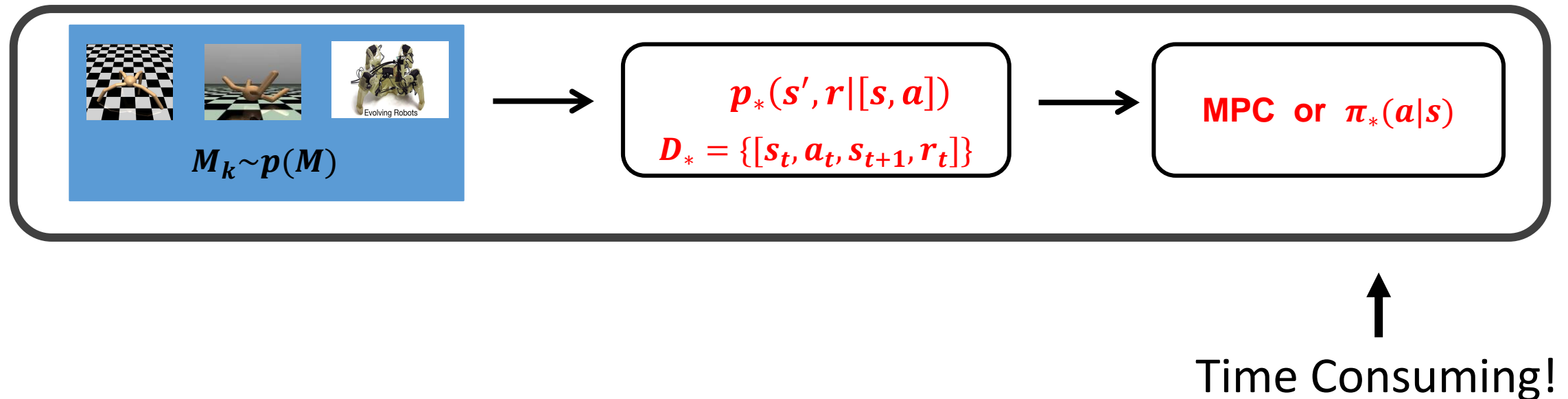


Learning Objective in MDPs:

$$\max_{\pi} J(\pi) = \mathbb{E}_{\mu_0, P, \pi} [\sum_{t=0}^{H-1} \gamma^t r_t]$$

Model-based Meta RL Might Make Sense! 😊

- RL in a collection of tasks $M_k = \langle S, A, P_k, R_k, \gamma \rangle$, $M_k \sim p(M)$:
 - ✓ Learn to adapt to a New Environment M_* within a few shots of interactions.
 - ✓ Execute time-consuming MPC or retrain policies in an approximate MDP $p_*(s', r|[s, a])$.



Proposed Method & Contributions

A New Optimization Framework of MBMRL

- **Meta Dynamics Model (GSSM):**

- ✓ System identification with L.V.s
→ Task representation
- ✓ Boost NPs with Message Passing
→ Minimize model discrepancy


- **Amortized Policy Search (APS):**

- ✓ Amortize Task-specific policies
→ Policy fast adaptations (no gradient updates)

$$\max_{\theta} \mathbb{E}_{\mathcal{M} \sim p(M)} \mathbb{E}_{([s,a],s') \sim \mathcal{M}} \ln [p_{\theta_{\mathcal{M}}}(s'|[s,a])] \quad \text{s.t. } p_{\theta_{\mathcal{M}}} = u(\theta, \mathcal{D}_{\mathcal{M}}^{\text{tr}})$$
$$\max_{\varphi} \mathbb{E}_{\mathcal{M} \sim p(M)} \mathbb{E}_{\substack{s' \sim p_{\theta_{\mathcal{M}}}(s'|[s,a]) \\ a \sim \pi_{\varphi_{\mathcal{M}}}}} \left[\sum_{t=0}^{H-1} \gamma^t r_{\mathcal{M}}(s_t, a_t) \right] \quad \text{s.t. } \pi_{\varphi_{\mathcal{M}}} = v(\varphi, \mathcal{E}_{\mathcal{M}}^{\text{tr}})$$

□ Graph Structured Surrogate Model
 $p_{\theta}(s_{t+1} | [s_t, a_t], z_t)$

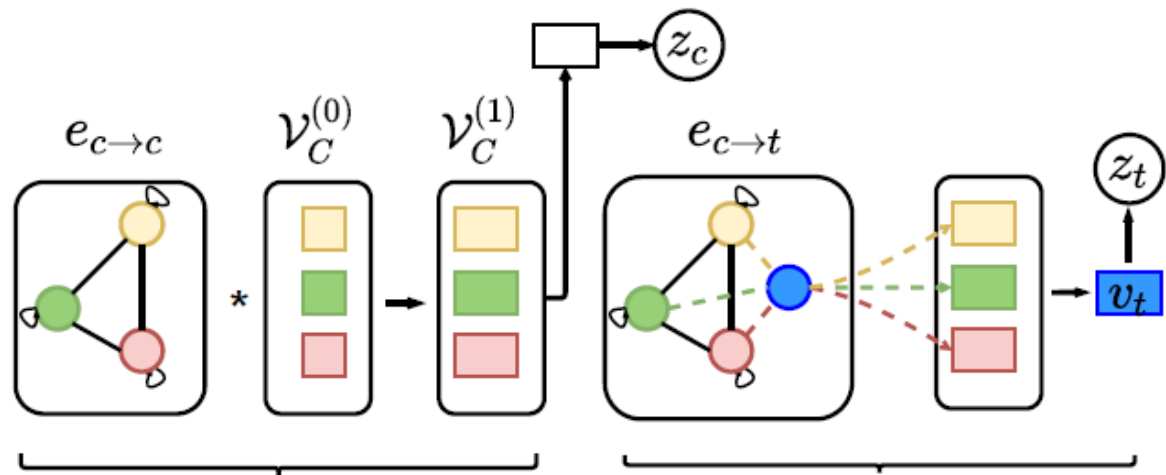
□ Amortized Policy Search
 $\pi_{\varphi}(a | s, z_c)$



Training Framework: *Unify two independent phases together using L.V.s.*
This is from a variant of *Thompson/Posterior Sampling* perspective.

$$c = [x_c, y_c] \xrightarrow{z_c} x_t \xrightarrow{v_t} z_t$$

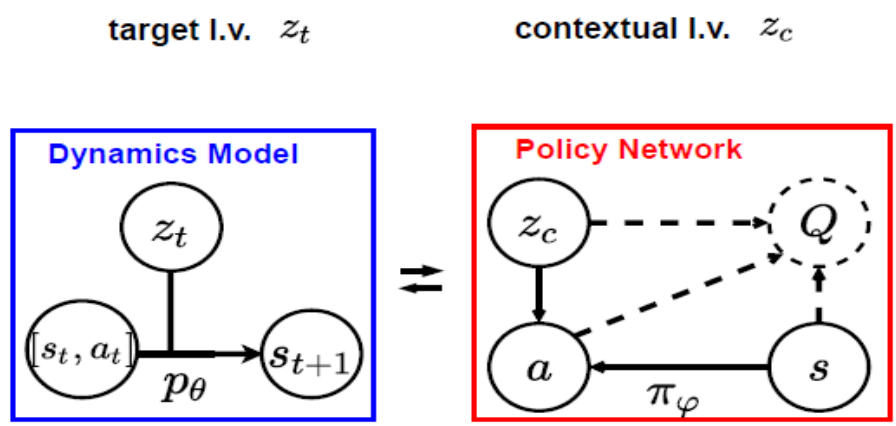
$$x = [s, a], y = s' \text{ or } y = \Delta s$$



Meta Learning Dynamics Models

$$\mathbb{E}_{p(D)} [\ln p(y_t|x_t, x_c, y_c)] \geq \mathbb{E}_{p(D)} [\mathbb{E}_{q_{\phi_1}} [\ln p_{\theta}(y_t|x_t, z_t)] - D_{KL}[q_{\phi_1}(z_t|x_t, x_c, y_c) \parallel q_{\phi_2}(z_c|x_c, y_c)]] \quad (6)$$

Meta Model-based Policy Search



$$\nabla_{\varphi} \mathcal{J}(\varphi) = \mathbb{E}_{\substack{\mathcal{M} \sim p(\mathcal{M}; \theta, \phi) \\ \tau \sim p(\tau | \mathcal{M}; \varphi, \phi)}} \left[\sum_{t=0}^{T-1} \nabla_{\varphi} \ln \pi(a_t | [s_t, z_c]) \cdot A_t([s_t, z_c], a_t) \right] \quad (11)$$

Training Paradigm in GSSM in the Background of a Collection of MDPs

← Dyna-style

Results & Analysis

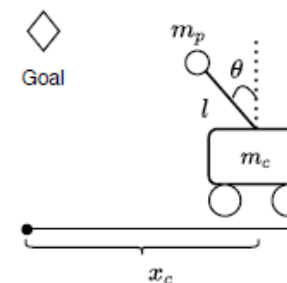
Generate a Collection of MDPs

- **Varying Dynamics via Adapting Hyper-Parameters**

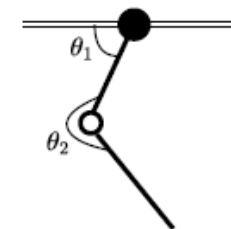
- ✓ Cart-Pole Swing-up/ Acrobot/ Mujoco-Tasks

- **Baselines**

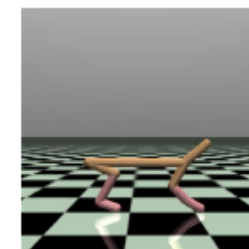
- ✓ M-DPILCO → BNN-based Meta-RL
- ✓ MLSM-v0 → Latent Variable-based Meta-RL
- ✓ MLSM-v1 → Latent Variable-based Meta-RL
- ✓ L2A → MAML-based Meta-RL
- ✓ DR-PPO → Train PPO across MDPs by Domain Randomization
- ✓ PE-PPO → Train PPO with Probabilistic Embedding Contexts $\pi(a|s, z)$



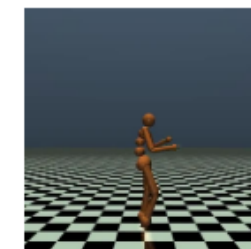
(a) Cart-Pole



(b) Acrobot



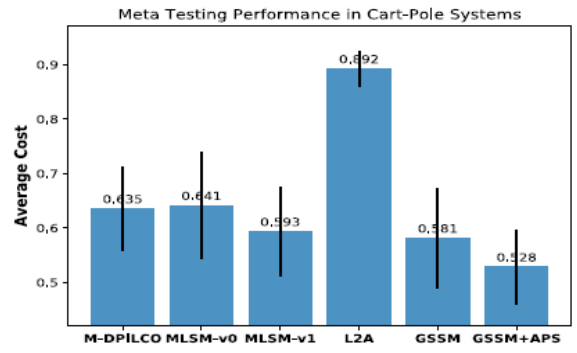
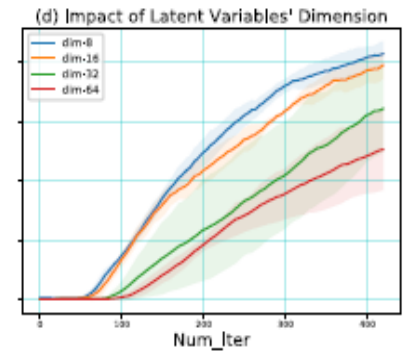
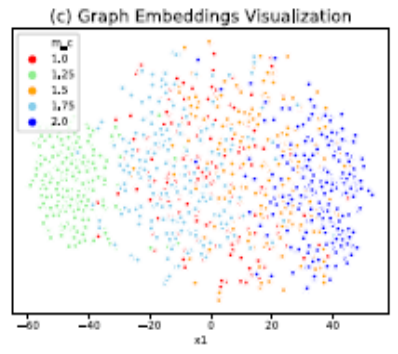
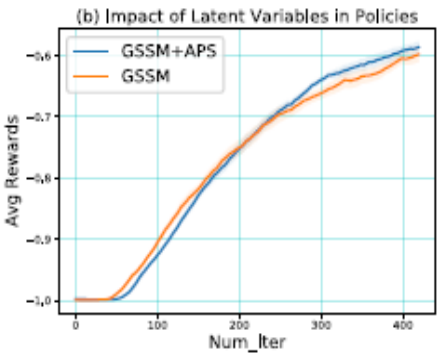
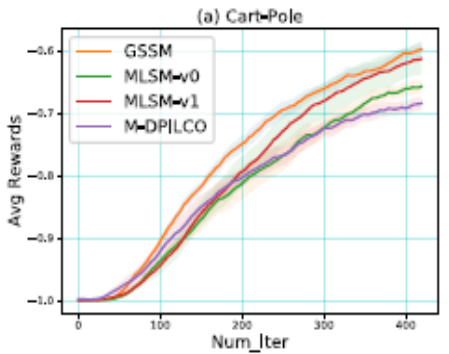
(c) Half-Cheetah



(d) Sllm-Humanoid

Some Results and Analysis

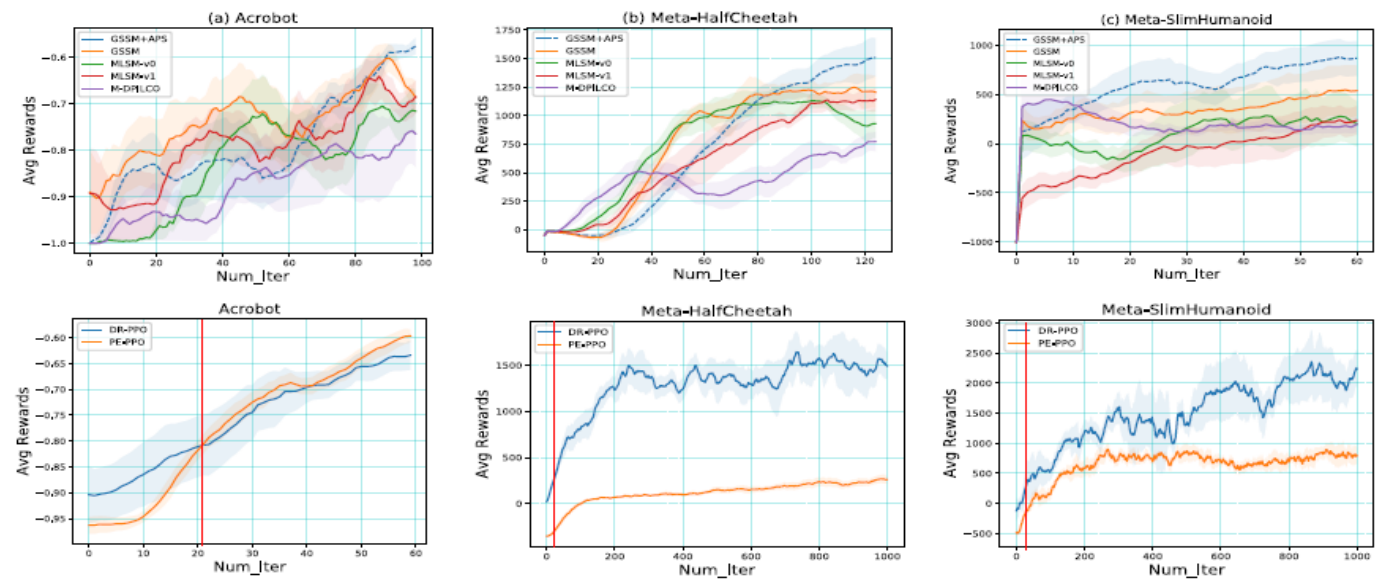
- Cart-Pole Meta-Training (Left)/Testing Results (Right)



- (1) GSSM/MLSM-v1 exhibit best performance in meta-training.
- (2) Latent variables are meaningful by varying physics parameters.
- (3) GSSM + Amortized Policy Search works best in meta-testing.

Some Results and Analysis

- Acrobot/Mujoco Learning Curves in Meta-Training Processes



- (1) Our method (first column blue ones) is also effective in Acrobot/Mujoco environments.
- (2) Most MBMRL baselines outperform model-free ones with same sample complexity.

Some Results and Analysis

- **Acrobot/Mujoco Performance in Meta-Testing Results**

Table 2. Average Rewards in Meta-testing Tasks using Learned Policy Networks. (For each testing task, 50 episodes are sampled and averaged in rewards. Figures in brackets are standard deviations across testing tasks, with bold ones the best.)

ENV	GSSM+APS	GSSM	M-DPILCO	MLSM-v0	MLSM-v1	L2A
ACROBOT	-0.478(±0.049)	-0.506(±0.068)	-0.645(±0.06)	-0.560(±0.064)	-0.524(±0.052)	-0.7775(±0.054)
H-CHEETAH	1597.4(±200)	1306.6(±140)	862.0(±280)	827.3(±190)	1226.8(±64)	-17.9(±130)
S-HUMANOID	1641.8(±170)	717.1(±130)	596.0(±340)	285.9(±360)	745.6(±150)	124.9(±570)

- (1) With non-amortized policy search, GSSM and MLSM-v1 are comparable in cumulative returns.
- (2) GSSM + Amortized Policy Search significantly surpasses others.

Limitations & Future Work

- Dyna-style training is sensitive to hyper-parameters, e.g. number of loops/optimization steps in each iteration and so forth.
- The model-based framework is data efficient but computationally intensive in meta training processes.
- Scalability in vision-based control has not been investigated yet.

Thanks for your Listening~