# Memory Sequence Length of Data Sampling Impacts the Adaptation of Meta-Reinforcement Learning Agents

 $\begin{array}{c} \text{Menglong Zhang}^{1[0009-0000-0304-3178]}, \, \text{Fuyuan Qian}^{1[0009-0002-2560-7872]}, \, \text{and} \\ \text{Quanying Liu}^{1[0000-0002-2501-7656]} \end{array}$ 

<sup>1</sup> Department of Biomedical Engineering, Southern University of Science and Technology, Shenzhen, 518055, China Corresponding to liuqy@sustech.edu.cn (Q.Liu)

**Abstract.** Fast adaptation to new tasks is extremely important for embodied agents in the real world. Meta-reinforcement learning (meta-RL) has emerged as an effective method to enable fast adaptation in unknown environments. Compared to on-policy meta-RL algorithms, off-policy algorithms rely heavily on efficient data sampling strategies to extract and represent the historical trajectories. However, little is known about how different data sampling methods impact the ability of meta-RL agents to represent unknown environments. Here, we investigate the impact of data sampling strategies on the exploration and adaptability of meta-RL agents. Specifically, we conducted experiments with two types of off-policy meta-RL algorithms based on Thompson sampling and Bayesoptimality theories in continuous control tasks within the MuJoCo environment and sparse reward navigation tasks. Our analysis revealed the long-memory and short-memory sequence sampling strategies affect the representation and adaptive capabilities of meta-RL agents. We found that the algorithm based on Bayes-optimality theory exhibited more robust and better adaptability than the algorithm based on Thompson sampling, highlighting the importance of appropriate data sampling strategies for the agent's representation of an unknown environment, especially in the case of sparse rewards.

**Keywords:** Meta-Reinforcement Learning  $\cdot$  Embodied Agent  $\cdot$  Data Sampling Strategy  $\cdot$  Task Adaptation  $\cdot$  Task Representation

## 1 Introduction

The realization of embodied intelligence relies on an agent's ability to fast adapt and generalize to unfamiliar environments. The core of adaptation and generalization is to transfer the knowledge learned during training to new task scenarios. Meta-RL is considered as one of the most effective approaches to facilitate fast adaptation to new tasks for embodied intelligence. The goal of meta-RL is to learn a policy within a given task distribution, which can efficiently adapt to a new task distribution with minimal data acquisition [5,3,23]. This goal of meta-RL, i.e., learning for fast adaptation, is well aligned with the learning strategy

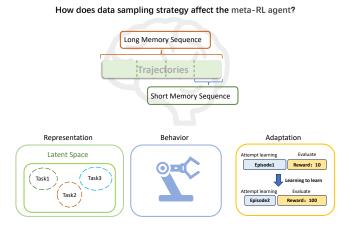


Fig. 1. Motivations of our work.

of humans and embodied AI. Therefore, understanding the mechanism of fast adaptation in meta-RL agents will shed light on the human and embodied AI.

Task representation, as a critical component of meta-RL, impacts the agent's generalization capabilities, particularly in complex control and navigation tasks. Agents are influenced by different task representations during training [31,30,22,28] and need to effectively learn representations to abstract the shared structures in the task distribution. The ability of task representation of an agent mainly depends on data sampling, such as utilizing data relevant to task generalization. Therefore, how to effectively sample task-relevant data during training is crucial. especially for training off-policy meta-RL agents within an off-policy framework. Existing meta-RL algorithms can be categorized into the policy gradient methods [5,33,18] and context-based methods [27,3,20,34]. These meta-RL methods can be categorized into on-policy and off-policy based on the relationship between the policy employed during the learning process and the policy used to generate the data. Compared to on-policy meta-RL algorithms, the off-policy algorithms are more sample-efficient and also rely more on suitable data sampling strategies [20,21,2]. However, how the data sampling strategy affects the off-policy meta-RL agents is still unclear.

In this study, we aim to understand how the data sampling strategy affects the online off-policy meta-RL algorithms, in terms of task representation, agent behaviors and adaptation ability (Figure 1). To this end, we analyze the exploration capabilities of two types of off-policy meta-RL algorithms based on different data sampling strategies, specifically those based on Thompson sampling [24] in PEARL [20] and Bayes-optimal policy [4] in VariBAD [34]. By conducting experiments with two meta-RL algorithms based on Bayes-optimal policy and Thompson sampling using different data sampling strategies, we examine the influence of two data sampling strategies, i.e., long and short memory

sequence, on task representation, agent behavior, and adaptability of off-policy meta-RL agents through experiments in continuous control tasks in MuJoCo [25] and complex navigation tasks.

Our findings are summarized as follows.

- Meta-RL based on Bayes-optimal policy has superior robustness to data sampling distributions compared to Thompson sampling-based Meta-RL method in sparse reward tasks. This robustness originates from the better representation of the unknown environment's dynamics and reward models (Figure 5, 6 and 7).
- Experiments on complex robotic navigation tasks demonstrate that although the short memory sampling strategy enables PEARL to converge faster, it does not improve the agent's adaptability; in contrast, the relatively robust off-policy VariBAD algorithm exhibits stronger adaptability (Figure 8, 9 and 10).
- The robustness of algorithms to short memory sequence or long memory sequence sampling strategy is associated with their adaptability capabilities.

# 2 Background and Related Work

In this section, we primarily introduce the foundational concepts of POMDPs and meta-RL, and related work on task representation in reinforcement learning.

## 2.1 POMDP and Meta-RL

A partially observable Markov decision process (POMDP) [9] framework offers a robust mathematical model for decision-making where agents must act under conditions of uncertainty and partial information. A POMDP is defined as a tuple  $(S, A, O, T, Z, R, \gamma)$ , where S is the state space, containing all possible states of the environment; A is the action space, containing all actions that the agent can perform;  $T: S \times A \to \mathcal{P}(S)$  is the state transition function;  $Z: S \times A \to \mathcal{P}(O)$  is the observation function, defining the probability distribution of generating observations given the next state and action;  $R: S \times A \times S \to \mathbb{R}$  is the reward function, which computes the immediate reward based on the current state, chosen action, and resulting state, and  $\gamma$  is the discount factor. The optimization objective of a POMDP is to find a policy  $\pi: \mathcal{H} \to A$ , where  $\mathcal{H}$  represents all possible sequences of historical information. In off-policy meta-RL methods, sufficient task representation from historical trajectories influences the online performance.

The adaptation process of meta-reinforcement learning agents in unknown environments can be seen as a generalization process within POMDPs with a similar distribution. In conventional reinforcement learning algorithms, policy  $\pi$  aims to maximize the expected discounted cumulative reward, expressed as:

$$\mathcal{J}^{\pi} = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^{t} R\left(s_{t}, a_{t}, s_{t+1}\right) \mid \pi\right]. \tag{1}$$

Meta-RL extends the foundational concepts of reinforcement learning by enabling agents to learn how to learn across a variety of tasks, rather than optimizing for a single task. This approach leverages past experience to rapidly adapt to new environments or tasks with minimal additional data. The framework of meta-RL is built around the idea that the skills acquired in previous tasks can inform the agent's policy on unseen tasks, thus reducing the time and data required for learning new tasks. The objective is to train a learning algorithm that can quickly adapt to new tasks using only a few interactions:

$$\max_{\theta} \mathbb{E}_{\tau \sim p(\tau|\theta,T)} \left[ \sum_{t=0}^{T} \gamma^{t} r_{t} \right], \tag{2}$$

where  $\theta$  represents the meta-parameters of the policy, and T represents a task sampled from a distribution of tasks. In meta-RL, we aim to perform well across a variety of such tasks.

#### 2.2 Task Representation in Reinforcement Learning

Effective task representation involves capturing the essential features of different tasks in a manner that accentuates their commonalities and differences, thus enabling the agent to adapt learned strategies to new, yet similar, scenarios. An expressive representation that captures task variations is vital for reducing the number of interactions needed to adapt to new tasks [8]. Previous work has employed reconstruction loss to train auto-encoders to generate low-dimensional representations of tasks, which are then used to assist in policy learning [11,29,17]. This method is also applicable to the extraction of common features in multi-task settings, often utilized in multi-task reinforcement learning [32,31]. Additionally, some works have used contrastive learning in the latent representation space to obtain robust representations across multiple tasks [12,28,30].

In this paper, we utilize two of the most fundamental and effective meta-RL models. PEARL [20] employs an RNN encoder as either the task representation module or the inference module. VariBAD [34,2] uses a Variational Autoencoder (VAE) [10] as both the representation and prediction module, which not only extracts task representations but also predicts the environmental model during training.

## 3 Models and Data Sampling Strategy

In this section we introduce models used and show how to use long-term memory replay and short-term memory reply in two context-based meta-RL methods.

## 3.1 Thompson Sampling and PEARL

Thompson sampling [24] is a Bayesian approach to addressing the exploration-exploitation dilemma and has been effectively applied in the context of meta-RL [16]. In sequential decision-making, meta-learning can be divided into two

phases: the first phase involves abstracting a representation of the distribution of training tasks, while the second phase allows the policy to leverage the prior knowledge of the task distribution acquired in the first phase to predict unknown task distributions, achieving rapid adaptation with minimal interactions with the environment. These models are typically comprised of two components: a task inference module and a policy module [1,20,34,19]. Thompson sampling can be described as the process of sampling actions from the policy based on posterior predictions [15]. According to Bayes' theorem, given the historical trajectory  $\tau$ , we can update the posterior distribution of  $\theta$ :  $P(\theta \mid \tau) \propto P(\tau \mid \theta)P(\theta)$ , where

$$P(\tau \mid \theta) = \prod_{t=0}^{T-1} P(s_{t+1} \mid s_t, a_t, \theta) P(r_t \mid s_t, a_t, \theta).$$
 (3)

We employ PEARL as a meta-RL model based on Thompson sampling on historical information. The inference network (task encoder)  $q_{\phi}$  encodes the agent's historical information to capture task-relevant sufficient statistics, and  $\mathbf{z}$  is the latent embedding of the encoder. During the meta-training phase, the parameters within  $q_{\phi}$  are optimized by modeling the Q-value function  $Q_{\theta}(\mathbf{s}, \mathbf{a}, \mathbf{z})$  within the Soft Actor-Critic (SAC) [7] algorithm and constrained by the variational lower bound to enhance the inference module's ability to learn information pertinent to the current task being performed, the inference module and policy module use different experience.

In the meta-testing phase, as the agent tackles unknown tasks, it updates the posterior of task history in a manner similar to Thompson sampling through the inference network, enhancing its capability to explore unknown tasks. The structure of PEARL is depicted in the upper part of Figure 2.

## 3.2 Bayes-optimality and VariBAD

Bayes-optimality is a principle within the broader Bayesian decision theory that dictates selecting actions based on maximizing expected utility, considering all possible outcomes weighted by their probabilities. In the context of meta-RL, a Bayes-optimal policy aims to maximize the expected reward across a distribution of tasks by leveraging a posterior distribution over tasks [4]. This approach inherently balances exploration and exploitation by considering the uncertainty in the environment's dynamics and the reward function. By integrating prior knowledge with observations gathered during interactions with the environment, the Bayes-optimal approach adapts its strategy to better respond to unseen dynamics. In BAMDP, we aim to maximize the expected reward over T time steps:

$$\mathcal{J}(\pi) = \mathbb{E}_{b_0,\pi} \left[ \sum_{t=0}^{T-1} \mathbb{E}_{R \sim b_t} \left[ R\left(s_t, a_t\right) \right] \right], \tag{4}$$

where  $b_t = p(r, p \mid \tau_{:t})$  represents the belief about the current environment dynamics and reward function based on the historical trajectory.

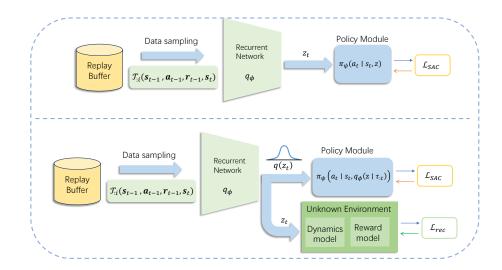


Fig. 2. PEARL and off-policy VariBAD famework.

VariBAD implements Bayesian optimal policies under the framework of Bayes-Adaptive MDP (BAMDP) [4,6] and it is designed to learn a latent representation of the environment's dynamics, rewards, and task-specific parameters through a variational autoencoder architecture (lower part of Figure 2). During training, the agent learns a universal policy that is conditioned on both the current state and a latent variable that encodes task-specific information. This latent variable is updated as new data is collected, effectively allowing the agent to infer the underlying task dynamics and adapt its policy accordingly. VariBAD thus enables an agent to perform robustly across a variety of tasks by efficiently learning and updating its understanding of the task environment.

## 3.3 Long and Short Memory Sequence Sampling

In this paper, we investigate the impact of data sampling strategies on the prior information extracted by the representation module in off-policy meta-RL, leading to varying exploration outcomes by the algorithm. In our experiments, we have set up long and short memory sequence sampling strategies (Figure 3). The long memory strategy involves storing all historical information from interactions with the environment in the replay buffer, from which the agent and inference network sample during training. In contrast, the short memory sequence sampling setting requires clearing the replay buffer at the start of each training iteration, ensuring that the agent and inference network can only sample the most recent historical data. Typically, Thompson sampling updates the posterior based on the most recent historical information. In our experiments described in Section 4, we observe that the long memory sequence can significantly disrupt the exploration effectiveness of PEARL, while the off-policy VariBAD

remains robust under both settings. The exploration and exploitation trade-off is also reflected in the varying performance of the different representation modules.

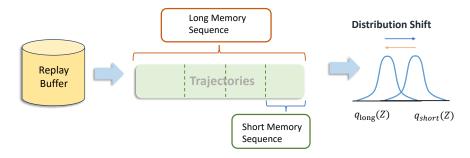


Fig. 3. Long memory sequence sampling and short memory sequence sampling. Different context sampling strategies can lead to shifts in the distribution of task representations, thereby affecting the agent's exploration and adaptation capabilities.

# 4 Experiments

In our study, we conduct comparisons of the performance of PEARL and off-policy VariBAD under various data sampling strategies in the continuous control task Sparse Half-Cheetah-Vel within the MuJoCo environment and two challenging navigation tasks, Ant-Semi-Circle and Sparse-Point-Robot. We analyze the robustness of the algorithms and their task adaptation capabilities based on their performance.

#### 4.1 Task setting

**Sparse Half-Cheetah-Vel.** We have modified the original Half-Cheetah environment in MuJoCo to adopt a sparse reward function, where the cheetah receives a reward of +1 only if it moves a distance greater than a specified threshold in each time step; otherwise, the reward is 0.

$$r_t^{sparse} = \begin{cases} 1, & \|x_t - x_{\text{goal}}\|_2 \le r \\ 0, & \text{otherwise.} \end{cases}$$

Subsequently, we utilize the commonly used Half-Cheetah-Vel task in meta-RL, where the objective for the agent is to reach a target velocity as quickly as possible. The target velocity is randomly sampled from a range of 0.0 to 3.0. Consequently, the reward in this environment is structured as follows:

$$r_t = -\left|v_t - v_{\mathrm{goal}}\right| - 0.05 \cdot \left\|a_t\right\|_2^2 + r_t^{sparse},$$

We set up 100 training tasks and 20 testing tasks.

**Sparse-Point-Robot.** We established the Sparse-Point-Robot environment to evaluate the performance of algorithms in a sparse reward navigation setting. Different goals are set on a semi-circle, with their locations unknown. At the start of each episode, the robot's initial position is randomly placed outside of the semi-circle. The objective is for the robot to locate the target within a single episode. The reward structure is configured as follows:

$$r = \begin{cases} 1, & \text{if } r \ge -\text{goal\_radius} \\ 0, & \text{otherwise.} \end{cases}$$

Ant-Semi-Circle. In the Ant-Semi-Circle task [2], an ant robot is required to navigate toward a goal that is randomly positioned on a semi-circle. Unlike the point robot scenario, this task employs the Ant model from the MuJoCo simulation, which introduces increased control complexity. The reward structure is set as follows:

$$r_t = -\|x_t - x_{\text{goal}}\|_1 - 0.1 \cdot \|a_t\|_2^2$$
.

## 4.2 Convergence of Algorithms

We conducted experiments in the aforementioned three environments, initially testing the convergence of the algorithms on different tasks under default parameters (PEARL using short sequence, off-policy VariBAD using long sequence).

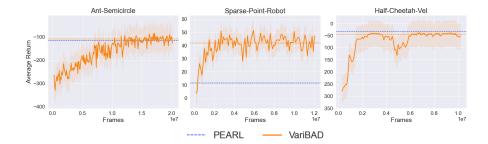


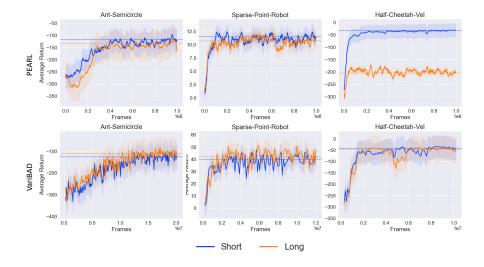
Fig. 4. Tasks training. Dashed lines correspond to the maximum return achieved by PEARL after 1e6 steps. Solid lines correspond to average return achieved by VariBAD. In the Ant-Semicircle and Half-Cheetah-Vel tasks, PEARL and VariBAD converge to similar average returns. However, in the Sparse-Point-Robot task, VariBAD significantly outperforms PEARL.

#### 4.3 Robustness of Algorithms

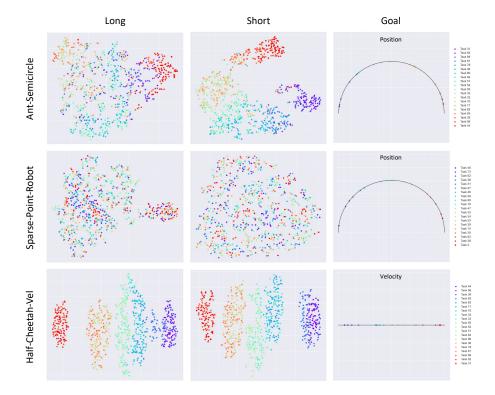
For algorithm robustness, the analysis mainly focuses on the task representations during the meta-training and meta-testing phases. During the meta-training

phase, for tasks involving the control of simulated robots such as Ant-SemiCircle and Half-Cheetah-Vel, PEARL is more susceptible to the influence of memory sequence sampling. Here, short memory sequence sampling proves more advantageous for PEARL's adaptation to new environments, while off-policy VariBAD remains relatively stable in comparison (Figure 5).

We randomly generate 20 goals from each environment, and for each goal, the agent performs 40 runs, each containing 5 episodes. We utilize t-SNE [14] to visualize the latent embeddings of the task inference module at the last time step of the fifth episode during the meta-testing phase for PEARL (Figure 6) and off-policy VariBAD (Figure 7). For navigation tasks, especially the Sparse-Point-Robot task, PEARL does not sufficiently learn the task-relevant information of the environment, whereas off-policy VariBAD exhibits better performance, forming clusters for similar positions on the semicircle. In the case of Ant-Semi-Circle, off-policy VariBAD accurately represents each target's position on the semicircle in the latent space, whereas the clusters formed by PEARL are more dispersed. In the Half-Cheetah-Vel environment, the representation of off-policy VariBAD is more sensitive to the memory sequence sampling strategy because, inherently, based on Bayes-optimality, it requires sampling of the whole history, and the representation in continuous control tasks is significantly influenced by the history.



**Fig. 5.** The average return during the meta-training phase of PEARL and off-policy VariBAD after using short and long memory sampling strategies.



 ${\bf Fig.\,6.}$  The t-SNE visualization of latent embedding during PEARL adaptation to the environment.

## 4.4 Tasks Adaptation Performance

To compare the adaptability of PEARL and VariBAD using different sampling strategies in unknown environments, we conducted rollouts of 5 episodes in each environment (Figure 8). The experiments revealed that the long memory sequence sampling strategy prevented PEARL from adapting effectively to navigation tasks and from achieving satisfactory results in the Sparse-Point-Robot task, consistent with its performance during the meta-training phase. In contrast, off-policy VariBAD demonstrated stable adaptability, successfully adjusting to tasks and achieving high average returns from the first episode across all three environments. Moreover, VariBAD showed less sensitivity to memory sequence sampling strategy, with agents trained using short memory sequence in the Sparse-Point-Robot task exhibiting enhanced exploratory capabilities.

Furthermore, to intuitively assess the impact of sampling strategies on the exploratory capabilities of agents during the meta-testing process, we visualized the exploration trajectories of agents in the Ant-Semi-Circle environment. The length of the memory sequence affects PEARL's exploratory capabilities, primarily because Thompson sampling focuses on updating the encoder's pos-

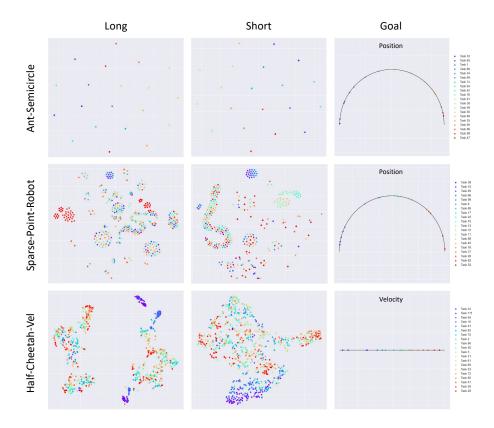


Fig. 7. The t-SNE visualization of latent embedding during off-policy VariBAD adaptation to the environment.

terior based on recent context, hence performing better with a short memory sequence length. Figure 9 shows that PEARL, when using short memory, can successfully reach the target in all five episodes; however, when switched to long memory, the agent fails to accurately locate the target for the same test task. Figure 10 demonstrates that off-policy VariBAD can successfully find the target in two episodes and is unaffected by the memory sequence length. This indicates that applying Bayes-optimality to off-policy meta-reinforcement learning results in stronger online adaptability to complex environments and more robust task representation compared to algorithms based on Thompson sampling.

# 5 Conclusion

**Discussion.** In this work, we comprehensively compared the impact of memory sequence length sampling on task representation, agent behavior, and the explo-

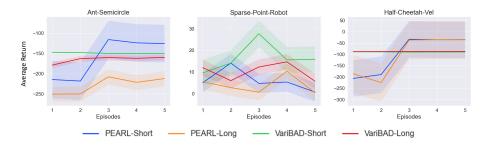
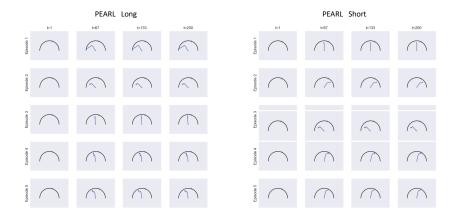


Fig. 8. Adaptation performance of PEARL and off-policy VariBAD using short and long memory sequence sampling strategy.



 ${\bf Fig.\,9.} \ {\bf Behavior} \ {\bf visualization} \ {\bf of} \ {\bf PEARL} \ {\bf during} \ {\bf adaptation} \ {\bf in} \ {\bf Ant-Semi-Circle}.$ 

ration and adaptation capabilities of two types of context-based off-policy meta-RL algorithms. The off-policy VariBAD algorithm, based on Bayes-optimality, demonstrated stronger robustness in sparse reward environments and its adaptability to unknown environments and task representations were less influenced by the training-time memory sequence. Although the PEARL algorithm, based on Thompson sampling, achieved similar average returns during the training phase, different memory lengths significantly affected its exploratory capabilities during the adaptation phase, fundamentally because the memory sequence length can cause shifts in the distribution of task representations.

**Future work.** Task representation extraction, especially in multi-task scenarios, is crucial for an agent's ability to adapt to unknown tasks [22,26,31,13]. This paper explores the impact of memory sequence length on shifts in task feature distributions, providing insights into how to maximize the extraction of information relevant to the current task from limited historical data and generate

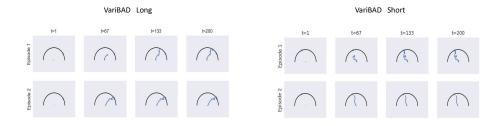


Fig. 10. Behavior visualization of off-policy VariBAD during adaptation in Ant-Semi-Circle.

a robust representation of the task distribution. This prompts us to employ effective representation learning methods in our subsequent work to enhance the performance of the task inference module.

# Acknowledgments

This work was funded in part by the National Key R&D Program of China (2021YFF1200804), Shenzhen Science and Technology Innovation Committee (2022410129, KCXFZ20201221173400001, KJZD20230923115221044).

#### References

- 1. Beck, J., Vuorio, R., Liu, E.Z., Xiong, Z., Zintgraf, L., Finn, C., Whiteson, S.: A survey of meta-reinforcement learning, arXiv preprint arXiv:2301.08028 (2023)
- Dorfman, R., Shenfeld, I., Tamar, A.: Offline meta reinforcement learningidentifiability challenges and effective data collection strategies. Advances in Neural Information Processing Systems 34, 4607–4618 (2021)
- 3. Duan, Y., Schulman, J., Chen, X., Bartlett, P.L., Sutskever, I., Abbeel, P.: RL2: Fast reinforcement learning via slow reinforcement learning. arXiv preprint arXiv:1611.02779 (2016)
- 4. Duff, M.O.: Optimal Learning: Computational procedures for Bayes-adaptive Markov decision processes. University of Massachusetts Amherst (2002)
- Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep networks. In: International conference on machine learning. pp. 1126–1135. PMLR (2017)
- 6. Ghavamzadeh, M., Mannor, S., Pineau, J., Tamar, A., et al.: Bayesian reinforcement learning: A survey. Foundations and Trends® in Machine Learning 8(5-6), 359–483 (2015)
- 7. Haarnoja, T., Zhou, A., Abbeel, P., Levine, S.: Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In: International conference on machine learning. pp. 1861–1870. PMLR (2018)
- 8. Humplik, J., Galashov, A., Hasenclever, L., Ortega, P.A., Teh, Y.W., Heess, N.: Meta reinforcement learning as task inference. arXiv preprint arXiv:1905.06424 (2019)

- Kaelbling, L.P., Littman, M.L., Cassandra, A.R.: Planning and acting in partially observable stochastic domains. Artificial intelligence 101(1-2), 99–134 (1998)
- Kingma, D.P., Welling, M.: Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114 (2013)
- 11. Lange, S., Riedmiller, M.: Deep auto-encoder neural networks in reinforcement learning. In: The 2010 international joint conference on neural networks (IJCNN). pp. 1–8. IEEE (2010)
- 12. Laskin, M., Srinivas, A., Abbeel, P.: Curl: Contrastive unsupervised representations for reinforcement learning. In: International conference on machine learning. pp. 5639–5650. PMLR (2020)
- 13. Lee, S., Chung, S.Y.: Improving generalization in meta-rl with imaginary tasks from latent dynamics mixture. Advances in Neural Information Processing Systems **34**, 27222–27235 (2021)
- 14. Van der Maaten, L., Hinton, G.: Visualizing data using t-sne. Journal of machine learning research 9(11) (2008)
- 15. Ortega, P.A., Braun, D.A.: A minimum relative entropy principle for learning and acting. Journal of Artificial Intelligence Research **38**, 475–511 (2010)
- Ortega, P.A., Wang, J.X., Rowland, M., Genewein, T., Kurth-Nelson, Z., Pascanu, R., Heess, N., Veness, J., Pritzel, A., Sprechmann, P., et al.: Meta-learning of sequential strategies. arXiv preprint arXiv:1905.03030 (2019)
- Péré, A., Forestier, S., Sigaud, O., Oudeyer, P.Y.: Unsupervised learning of goal spaces for intrinsically motivated goal exploration. arXiv preprint arXiv:1803.00781 (2018)
- Raghu, A., Raghu, M., Bengio, S., Vinyals, O.: Rapid learning or feature reuse? towards understanding the effectiveness of maml. arXiv preprint arXiv:1909.09157 (2019)
- Raileanu, R., Goldstein, M., Szlam, A., Fergus, R.: Fast adaptation to new environments via policy-dynamics value functions. In: Proceedings of the 37th International Conference on Machine Learning. pp. 7920–7931 (2020)
- 20. Rakelly, K., Zhou, A., Finn, C., Levine, S., Quillen, D.: Efficient off-policy metareinforcement learning via probabilistic context variables. In: International conference on machine learning. pp. 5331–5340. PMLR (2019)
- Rusu, A.A., Rao, D., Sygnowski, J., Vinyals, O., Pascanu, R., Osindero, S., Hadsell, R.: Meta-learning with latent embedding optimization. arXiv preprint arXiv:1807.05960 (2018)
- Sodhani, S., Zhang, A., Pineau, J.: Multi-task reinforcement learning with context-based representations. In: International Conference on Machine Learning. pp. 9767–9779. PMLR (2021)
- 23. Stadie, B.C., Yang, G., Houthooft, R., Chen, X., Duan, Y., Wu, Y., Abbeel, P., Sutskever, I.: Some considerations on learning to explore via meta-reinforcement learning. arXiv preprint arXiv:1803.01118 (2018)
- 24. Thompson, W.R.: On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. Biometrika **25**(3-4), 285–294 (1933)
- Todorov, E., Erez, T., Tassa, Y.: Mujoco: A physics engine for model-based control. In: 2012 IEEE/RSJ international conference on intelligent robots and systems. pp. 5026–5033. IEEE (2012)
- 26. Wang, B., Xu, S., Keutzer, K., Gao, Y., Wu, B.: Improving context-based meta-reinforcement learning with self-supervised trajectory contrastive learning. arXiv preprint arXiv:2103.06386 (2021)

- 27. Wang, J.X., Kurth-Nelson, Z., Tirumala, D., Soyer, H., Leibo, J.Z., Munos, R., Blundell, C., Kumaran, D., Botvinick, M.: Learning to reinforcement learn. arXiv preprint arXiv:1611.05763 (2016)
- Wang, M., Bing, Z., Yao, X., Wang, S., Kai, H., Su, H., Yang, C., Knoll, A.: Meta-reinforcement learning based on self-supervised task representation learning. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 37, pp. 10157–10165 (2023)
- Watter, M., Springenberg, J., Boedecker, J., Riedmiller, M.: Embed to control: A locally linear latent dynamics model for control from raw images. Advances in neural information processing systems 28 (2015)
- 30. Yuan, H., Lu, Z.: Robust task representations for offline meta-reinforcement learning via contrastive learning. In: International Conference on Machine Learning. pp. 25747–25759. PMLR (2022)
- 31. Zhang, A., Sodhani, S., Khetarpal, K., Pineau, J.: Learning robust state abstractions for hidden-parameter block mdps. arXiv preprint arXiv:2007.07206 (2020)
- 32. Zhang, Y., Yang, Q.: A survey on multi-task learning. IEEE Transactions on Knowledge and Data Engineering **34**(12), 5586–5609 (2021)
- Zintgraf, L., Shiarli, K., Kurin, V., Hofmann, K., Whiteson, S.: Fast context adaptation via meta-learning. In: International Conference on Machine Learning. pp. 7693–7702. PMLR (2019)
- 34. Zintgraf, L., Shiarlis, K., Igl, M., Schulze, S., Gal, Y., Hofmann, K., Whiteson, S.: Varibad: A very good method for bayes-adaptive deep rl via meta-learning. arXiv preprint arXiv:1910.08348 (2019)