

What about Inputting Policy in Value Function: Policy Representation and Policy-Extended Value Function Approximator

Hongyao Tang,¹ Zhaopeng Meng,¹ Jianye Hao,^{1*} Chen Chen,² Daniel Graves,² Dong Li,² Changmin Yu,³ Hangyu Mao,² Wulong Liu,² Yaodong Yang,¹ Wenyuan Tao,¹ Li Wang¹

¹College of Intelligence and Computing, Tianjin University

²Noah's Ark Lab, Huawei

³Gatsby Computational Neuroscience Unit, University College London

Abstract

We study Policy-extended Value Function Approximator (PeVFA) in Reinforcement Learning (RL), which extends conventional value function approximator (VFA) to take as input not only the state (and action) but also an explicit policy representation. Such an extension enables PeVFA to preserve values of multiple policies at the same time and brings an appealing characteristic, i.e., *value generalization among policies*. We formally analyze the value generalization under Generalized Policy Iteration (GPI). From theoretical and empirical lens, we show that generalized value estimates offered by PeVFA may have lower initial approximation error to true values of successive policies, which is expected to improve consecutive value approximation during GPI. Based on above clues, we introduce a new form of GPI with PeVFA which leverages the value generalization along policy improvement path. Moreover, we propose a representation learning framework for RL policy, providing several approaches to learn effective policy embeddings from policy network parameters or state-action pairs. In our experiments, we evaluate the efficacy of value generalization offered by PeVFA and policy representation learning in several OpenAI Gym continuous control tasks. For a representative instance of algorithm implementation, Proximal Policy Optimization (PPO) re-implemented under the paradigm of GPI with PeVFA achieves about 40% performance improvement on its vanilla counterpart in most environments.

1 Introduction

Reinforcement Learning (RL) has been widely considered as a promising way to learn optimal policies in many decision-making problems (Mnih et al. 2015; Lillicrap et al. 2015; Silver et al. 2016; You et al. 2018; Schreck, Coley, and Bishop 2019; Vinyals et al. 2019; Hafner et al. 2020). One fundamental element of RL is value function which defines the long-term evaluation of a policy. With function approximation (e.g., deep neural networks), a value function approximator (VFA) is able to approximate the values of a policy under large and continuous state spaces. As commonly recognized, most RL algorithms can be described as Generalized Policy Iteration (GPI) (Sutton and Barto 1998). As illustrated on the left of Fig.1, at each iteration the VFA is trained to

approximate the true values of current policy (i.e., policy evaluation), regarding which the policy is further improved (i.e., policy improvement). The value function approximation error hinders the effectiveness of policy improvement and then the overall optimality of GPI (Bertsekas and Tsitsiklis 1996; Scherrer et al. 2015). Unfortunately, such errors are inevitable under function approximation. A large number of samples are usually required to ensure high-quality value estimates, resulting in the sample-inefficiency of deep RL algorithms. Therefore, this raises an urgent need for more efficient value approximation methods (v. Hasselt 2010; Bellemare, Dabney, and Munos 2017; Fujimoto, v. Hoof, and Meger 2018; Kuznetsov et al. 2020).

An intuitive idea to improve the efficiency value approximation is to leverage the knowledge on the values of previous encountered policies. However, a conventional VFA usually approximates the values of one policy and values learned from old policies are over-written gradually during the learning process. This means that the previously learned knowledge cannot be preserved and utilized with one conventional VFA. Thus, such limitations prevent the potentials to leverage the previous knowledge for future learning. In this paper, we study Policy-extended Value Function Approximator (PeVFA), which additionally takes an explicit policy representation as input in contrast to conventional VFA. Thanks to the policy representation input, PeVFA is able to approximate values for multiple policies and induces value generalization among policies. We formally analyze the generalization of approximate values among policies in a general form. From both theoretical and empirical lens, we show that the generalized value estimates can be closer to the true values of the successive policy, which can be beneficial to consecutive value approximation along the policy improvement path, called *local generalization*. Based on above clues, we introduce a new form of GPI with PeVFA (the right of Fig.1) that leverages the local generalization to improve the efficiency of consecutive value approximation along the policy improvement path.

One key point of GPI with PeVFA is the representation of policy since it determines how PeVFA generalizes the values. For this, we propose a framework to learn effective low-dimensional embedding of RL policy. We use network parameters or state-action pairs as policy data and encode them into low-dimensional embeddings; then the embeddings

*Corresponding author: Jianye Hao (jianye.hao@tju.edu.cn).
Copyright © 2022, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

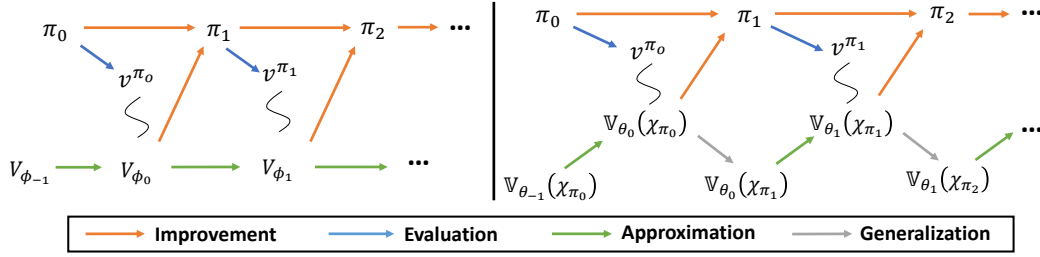


Figure 1: Generalized Policy Iteration (GPI) with function approximation. *Left*: GPI with conventional value function approximator V_ϕ . *Right*: GPI with PeVFA $\mathbb{V}_\theta(\chi_\pi)$ (Sec. 3) where extra generalization steps exist. The subscripts of policy π and value function parameters ϕ, θ denote the iteration number. The squiggle lines represent non-perfect approximation of true values.

are trained to capture the effective information through contrastive learning and policy recovery. Finally, we evaluate the efficacy of GPI with PeVFA and our policy representations. In principle, GPI with PeVFA is general and can be implemented in different ways. As a practical instance, we re-implement Proximal Policy Optimization (PPO) with PeVFA and propose PPO-PeVFA algorithm. Our experimental results on several OpenAI Gym continuous control tasks demonstrate the effectiveness of both value generalization offered by PeVFA and learned policy representations, with an about 40% improvement in average returns achieved by our best variants on standard PPO in most tasks.

We summarize our main contributions below. 1) We study the value generalization among policies induced by PeVFA. From both theoretical and empirical aspects, we shed the light on the situations where the generalization can be beneficial to the learning along policy improvement path. 2) We propose a framework for policy representation learning. To our knowledge, we make the first attempt to learn a low-dimensional embedding of over 10k network parameters for an RL policy. 3) We introduce GPI with PeVFA that leverages the value generalization in a general form. Our experimental results demonstrate the potential of PeVFA in deriving practical and more effective RL algorithms.

2 Related Work

Extensions of Conventional Value Function Sutton et al. (2011) propose General Value Functions (GVFs) as a general form of knowledge representation of rewards and arbitrary cumulants. Later, conventional value functions are extended to take extra inputs for different purposes of generalization. One notable work is Universal Value Function Approximator (UVFA) (Schaul et al. 2015), which is proposed to generalize values among different goals for goal-conditioned RL. UVFA is further developed and various extensions are studied in (Andrychowicz et al. 2017; Rakelly et al. 2019; Lee et al. 2020; Wang et al. 2020; He and Boyd-Graber 2016; Grover et al. 2018). Most of the above works study how to generalize the policy or value function among extrinsic factors, i.e., environments, tasks and opponents (we provide a unified view in Appendix E.2); while we mainly study the value generalization among policies along policy improvement path, an intrinsic learning process of the agent itself.

Policy Embedding and Representation. Although not well studied, representation (or embedding) learning for RL policies is involved in a few works (Hausman et al. 2018; Grover et al. 2018; Arnekjvst, Kragic, and Stork 2019). The most common way to learn a policy representation is to extract from interaction experiences. As a representative, (Grover et al. 2018) propose learning the representation of opponent policy from interaction trajectories with a generative policy recovery loss and a discriminative triplet loss. These losses are later adopted in (Wang et al. 2020; Raileanu et al. 2020). Another straightforward idea is to represent policy parameters. Network Fingerprint (Harb et al. 2020) is such a differentiable representation that uses the concatenation of the vectors of action distribution outputted by policy network on a set of probing states. The probing state set is co-optimized along with the primary learning objective, which can be non-trivial especially when the dimensionality of the set is high. See Appendix for a detailed review. Our work propose a learning framework of policy representation including both above two perspectives.

PVN and PVFs Recently, several works study the generalization among policy space. Harb et al. (2020) propose Policy Evaluation Network (PVN) to directly approximate the distribution of policy π 's objective function $J(\pi) = \mathbb{E}_{\rho_0}[v^\pi(s_0)]$ with initial state $s_0 \sim \rho_0$. PVN takes as input Network Fingerprint (mentioned above) of policy network. After training on a pre-collected set of policies, a random initialized policy can be optimized in a zero-shot manner with the policy gradients of PVN by backpropagating through the differentiable policy input. We call such gradients *GTPI* for short below. Similar ideas are later integrated with task-specific context learning in multi-task RL (Raileanu et al. 2020). In PVN (Harb et al. 2020), as an early attempt, the generalization among policies is studied with small policy network and simple tasks; besides, the most regular online learning setting is not studied. Concurrent to our work, Faccio et al. (2021) propose a class of Parameter-based Value Functions (PVFs) that take vectorized policy parameters as inputs. Based on PVFs, new policy gradient algorithms are introduced in the form of a combination of conventional policy gradients and GTPI. In addition to the zero-shot policy optimization as conducted in PVN, PVFs are also evaluated for online policy learning. Due to directly taking parameters as input, PVFs

suffer from the curse of dimensionality when the number of parameters is high. Besides, GTPI can be non-trivial to rein with usually large policy parameter space and finite policies (much fewer than state-action samples). We provide more discussions on GTPI in Appendix E.4.

Our work differs with PVFs from several aspects. First, we make use of learned policy representation rather than policy network parameters. Second, we do not resort to GTPI for the policy update in our algorithms but focus on utilizing value generalization for more efficient value estimation in GPI. Furthermore, we shed the light on two important problems — how value generalization among policies can happen formally and whether it is beneficial to learning or not — which are neglected in previous works from both theoretical and empirical lens. We refer one to read the original papers and our Appendix E.3 for better understandings of the differences.

3 Policy-extended Value Function Approximator

In this section, we propose Policy-extended Value Function Approximator (PeVFA), an extension of conventional VFA that explicitly takes as input a policy representation.

3.1 Formulation

Consider a Markov Decision Process (MDP) defined as $\langle \mathcal{S}, \mathcal{A}, r, \mathcal{P}, \gamma \rangle$ where \mathcal{S} is the state space, \mathcal{A} is the action space, r is the (bounded) reward function, \mathcal{P} is the transition function and $\gamma \in [0, 1)$ is the discount factor. A policy $\pi \in P(\mathcal{A})^{|\mathcal{S}|}$ defines the distribution over all actions for each state. The goal of an RL agent is to find an optimal policy π^* that maximizes the expected long-term discounted return. The state-value function $v^\pi(s)$ is defined as the expected discounted return by following the policy π from a state s : $v^\pi(s) = \mathbb{E}_\pi [\sum_{t=0}^{\infty} \gamma^t r_{t+1} | s_0 = s]$ where $r_{t+1} = r(s_t, a_t)$. We use V^π to denote the vectorized form of value function.

In a general form, we define *policy-extended value function* $\mathbb{V} : \mathcal{S} \times \Pi \rightarrow \mathbb{R}$ over state and policy space: $\mathbb{V}(s, \pi) = v^\pi(s)$ for all $s \in \mathcal{S}$ and $\pi \in \Pi$. In this paper, we focus on $\mathbb{V}(s, \pi)$ and policy-extended action-value function $\mathbb{Q}(s, a, \pi)$ can be obtained similarly. We use $\mathbb{V}(\pi)$ to denote the value vector for all states in the following. The key point is that \mathbb{V} is able to preserve the values of multiple policies. With function approximation, a PeVFA is expected to approximate the values of policies among policy space, i.e., $\{V^\pi\}_{\pi \in \Pi}$ then enable value generalization among policies.

Formally, given a function $g : \Pi \rightarrow \mathcal{X} \subseteq \mathbb{R}^n$ that maps any policy π to an n -dimensional representation $\chi_\pi = g(\pi) \in \mathcal{X}$, a PeVFA \mathbb{V}_θ with parameter $\theta \in \Theta$ is to minimize the approximation error over all possible states and policies generally: $F_{\mu, \rho, \mathbb{V}_\theta}(\theta, g, \Pi) = \sum_{\pi \in \Pi} \mu(\pi) \|\mathbb{V}_\theta(\chi_\pi) - V^\pi\|_{p, \rho}$, where μ, ρ are distributions over policies and states respectively, $\|f\|_{p, \rho} = (\int_{\mathcal{S}} \rho(ds) |f(s)|^p)^{1/p}$ is ρ -weighted L_p -norm (Lagoudakis and Parr 2003; Scherrer et al. 2015) for any $f : \mathcal{S} \rightarrow \mathbb{R}$. The policy distribution μ of interest depends on the scenario where value generalization is considered. As illustrated in Fig.2, we provide two value generalization scenarios. In the global generalization scenario, a uniform distribution over known policy set may be considered with a

general purpose of value generalization for unknown policies. For the specific local generalization scenario along policy improvement path during GPI, a sophisticated distribution that adaptively weights recent policies more during the learning process may be more suitable in this case. In the following, we care more about the local generalization scenario and use uniform state distribution ρ and L_2 -norm for demonstration. The subscripts are omitted and we use $\|\cdot\|$ for clarity.

3.2 Insights on Generalization among Policies

In this part, we provide preliminary theoretical analysis on value generalization among policies induced by PeVFA, to shed some light on whether the generalization can be beneficial to conventional RL under GPI paradigm. We start from a two-policy case and study whether the value approximation learned for one policy can be generalized to the other one. Later, we study the local generalization scenario (Fig.2(b)) and shed the light on the superiority of PeVFA for GPI. All the proofs are provided in Appendix A.

For the convenience of demonstration, we use an identical policy representation function, i.e., $\chi_\pi = \pi$, and define the approximation loss of PeVFA \mathbb{V}_θ for any policy $\pi \in \Pi$ as $f_\theta(\pi) = \|\mathbb{V}_\theta(\pi) - V^\pi\| \geq 0$. We use the following definitions for a formal description of value approximation process with PeVFA and local property of loss function f_θ related to generalization (Novak et al. 2018; Wang et al. 2018):

Definition 1 (π -Value Approximation) We define a value approximation process $\mathcal{P}_\pi : \Theta \rightarrow \Theta$ with PeVFA as a γ -contraction on the approximation loss for policy π , i.e., for $\hat{\theta} = \mathcal{P}_\pi(\theta)$, we have $f_{\hat{\theta}}(\pi) \leq \gamma f_\theta(\pi)$ where $\gamma \in [0, 1)$.

Definition 2 (L -Continuity) We call f_θ is L -continuous at policy π if f_θ is Lipschitz continuous at π with a constant $L \in [0, \infty)$, i.e., $|f_\theta(\pi) - f_\theta(\pi')| \leq L \cdot d(\pi, \pi')$ for $\pi' \in \Pi$ with some distance metric d for policy space Π .

With Definition 1, the consecutive value approximation for the policies along policy improvement path during GPI can

be described as: $\theta_{-1} \xrightarrow{\mathcal{P}_{\pi_0}} \theta_0 \xrightarrow{\mathcal{P}_{\pi_1}} \theta_1 \xrightarrow{\mathcal{P}_{\pi_2}} \dots$, as the green arrows illustrated in Fig.1. One may refer to Appendix A.1 for a discussion on the rationality of the two definitions.

To start our analysis, we first study the generalized value approximation loss in a two-policy case where only the value of policy π_1 is approximated by PeVFA as below:

Lemma 1 For $\theta \xrightarrow{\mathcal{P}_{\pi_1}} \hat{\theta}$, if $f_{\hat{\theta}}$ is \hat{L} -continuous at π_1 and $f_\theta(\pi_1) \leq f_\theta(\pi_2)$, we have: $f_{\hat{\theta}}(\pi_2) \leq \gamma f_\theta(\pi_2) + \mathcal{M}(\pi_1, \pi_2, \hat{L})$, where $\mathcal{M}(\pi_1, \pi_2, \hat{L}) = \hat{L} \cdot d(\pi_1, \pi_2)$.

Corollary 1 \mathcal{P}_{π_1} is γ_g -contraction ($\gamma_g \in [0, 1)$) for π_2 when $f_\theta(\pi_2) > \frac{\hat{L} \cdot d(\pi_1, \pi_2)}{1 - \gamma}$.

Lemma 1 shows that the post- \mathcal{P}_{π_1} approximation loss for π_2 is upper bounded by a generalized contraction of prior loss plus a locality margin term \mathcal{M} which is related to π_1, π_2 and the locality property of $f_{\hat{\theta}}$. In general, the form of \mathcal{M} depends on the local property assumed. For a step further, Corollary 1 reveals the condition where a contraction on value approximation loss for π_2 is achieved when PeVFA is only trained to approximate the values of π_1 .

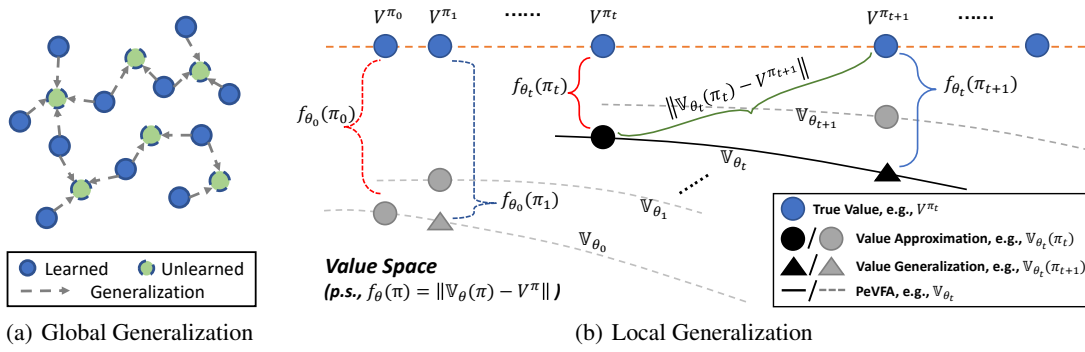


Figure 2: Illustrations of value generalization among policies of PeVFA. Each circle denotes value function (estimate) of a policy. (a) *Global Generalization*: values of known policies can be generalized to unknown policies. (b) *Local Generalization*: values of previous policies (e.g., π_t) can be generalized to successive policies (e.g., π_{t+1}) along policy improvement path.

Then we consider the local generalization scenario as illustrated in Fig.2(b). For any iteration t of GPI, the values of current policy π_t are approximated by PeVFA, followed by a improved policy π_{t+1} whose values are to be approximated in the next iteration. The value generalization from each π_t and π_{t+1} can be similarly considered as the two-policy case. In addition to the former results, we shed the light on the value generalization along policy improvement path below:

Lemma 2 For $\theta_{-1} \xrightarrow{\mathcal{P}_{\pi_0}} \theta_0 \xrightarrow{\mathcal{P}_{\pi_1}} \theta_1 \xrightarrow{\mathcal{P}_{\pi_2}} \dots$ with γ_t for each \mathcal{P}_{π_t} , if f_{θ_t} is L_t -continuous at π_t for any $t \geq 0$, we have $f_{\theta_t}(\pi_{t+1}) \leq \gamma_t f_{\theta_{t-1}}(\pi_t) + \mathcal{M}_t$, where $\mathcal{M}_t = L_t \cdot d(\pi_t, \pi_{t+1})$.

Corollary 2 By induction, we have $f_{\theta_t}(\pi_{t+1}) \leq \prod_{i=0}^t \gamma_i f_{\theta_{-1}}(\pi_0) + \sum_{i=0}^{t-1} \prod_{j=i+1}^t \gamma_j \mathcal{M}_i + \mathcal{M}_t$.

The above results indicate that the value generalization loss can be recursively bounded and has an upper bound formed by a repeated contraction on initial loss plus the accumulation of locality margins induced from each local generalization. An infinity-case discussion is in Appendix A.5.

Although the theory above depicts the value generalization among policies, it is not necessarily useful to conventional RL under GPI paradigm. Therefore, the next question is naturally whether PeVFA with value generalization among policies is preferable to the conventional VFA; if yes, what the case is to be like. To this end, we introduce a desirable condition which reveals the superiority of PeVFA during consecutive value approximation along the policy improvement path:

Theorem 1 During $\theta_{-1} \xrightarrow{\mathcal{P}_{\pi_0}} \theta_0 \xrightarrow{\mathcal{P}_{\pi_1}} \theta_1 \xrightarrow{\mathcal{P}_{\pi_2}} \dots$, for any $t \geq 0$, if $f_{\theta_t}(\pi_t) + f_{\theta_t}(\pi_{t+1}) \leq \|V^{\pi_t} - V^{\pi_{t+1}}\|$, then $f_{\theta_t}(\pi_{t+1}) \leq \|\mathbb{V}_{\theta_t}(\pi_t) - V^{\pi_{t+1}}\|$.

Theorem 1 shows that the generalized value estimates $\mathbb{V}_{\theta_t}(\pi_{t+1})$ can be closer to the true values of policy π_{t+1} than $\mathbb{V}_{\theta_t}(\pi_t)$. Note that $\mathbb{V}_{\theta_t}(\pi_t)$ is the value approximation for π_t which is equivalent to the counterpart V_{ϕ_t} for a conventional VFA as value generalization among policies does not exist. To consecutive value approximation along policy improvement path, this means that the value generalization of PeVFA has the potential to offer closer start points at each

iteration. If such closer start points can often exist, we expect PeVFA to be preferable to conventional VFA since value approximation can be more efficient with PeVFA and it in turn facilitates the overall GPI process.

However, the condition in Theorem 1 is not necessarily met in practice. It depends on the locality margins that may be related to function family and optimization method of PeVFA, as well as the scale of policy improvement. We conjecture that there are many looser sufficient conditions that lead to the consequence of Theorem 1, and the presented condition is the strictest one among them to achieve. This can be interpreted by considering the geometrical relationship between V^{π_t} , $V^{\pi_{t+1}}$, $\mathbb{V}_{\theta_t}(\pi_t)$ and $\mathbb{V}_{\theta_t}(\pi_{t+1})$. One special case that requires the sufficient condition presented in Theorem 1 is where $\mathbb{V}_{\theta_t}(\pi_t)$ and $\mathbb{V}_{\theta_t}(\pi_{t+1})$ and lie on the line segment between V^{π_t} and $V^{\pi_{t+1}}$. We leave these further theoretical investigations for future work. Instead, we empirically examine the existence of such desirable generalization below.

3.3 Empirical Evidences

We empirically investigate the value generalization of PeVFA with didactic environments. In this section, PeVFA \mathbb{V}_{θ} is parameterized by neural network and we use the concatenation of all weights and biases of the policy network as a straightforward representation χ_{π} for each policy, called *Raw Policy Representation (RPR)*. Details are provided in Appendix B.

First, we demonstrate the global generalization (illustrated in Fig.2(a)) in a continuous 2D Point Walker environment. We build the policy set Π with synthetic policies, each of which is a randomly initialized 2-layer *tanh*-activated neural network with 2 units for each layer. The size of Π is 20k and the behavioral diversity of synthetic policies is verified (see Fig.6(b) in Appendix). We divide Π into the known policies Π_0 for training \mathbb{V}_{θ} and the unseen policies Π_1 for testing. Fig.3(a) shows the value predictions for policies from training and testing set (100 for each). Our results show that a PeVFA trained on Π_0 achieves reasonable generalization performance when evaluating on Π_1 .

Next, we investigate the value generalization along policy improvement path, i.e., local generalization as in Fig.2(b). We use a 2-layer 8-unit policy network trained by standard PPO

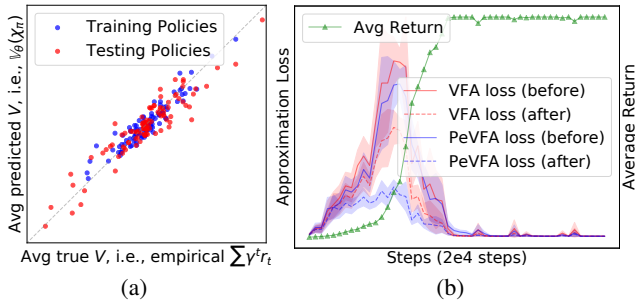


Figure 3: Empirical evidences of two kinds of generalization of PeVFA. (a) *Global generalization*: PeVFA shows comparable value estimation performance on testing policy set (red) after learning on training policy set (blue). (b) *Local generalization*: PeVFA ($\mathbb{V}_\theta(\chi_\pi)$) shows lower losses than conventional VFA (V_ϕ) before and after the value approximation training for successive policies along policy improvement path. In (b), the left axis is for approximation loss (lower is better) and the right axis is for average return as a reference of the policy learning process (green curve).

algorithm (Schulman et al. 2017) in MuJoCo continuous control tasks. Parallel to the conventional value network $V_\phi(s)$ (i.e., VFA) in PPO, we set a PeVFA network $\mathbb{V}_\theta(s, \chi_\pi)$ as a reference for the comparison on value approximation loss. Compared to V_ϕ , PeVFA $\mathbb{V}_\theta(s, \chi_\pi)$ takes RPR as input and approximates the values of all historical policies ($\{\pi_i\}_{i=0}^t$) in addition. We compare the value approximation losses of V_ϕ (red) and \mathbb{V}_θ (blue) before (solid) and after (dashed) updating with on-policy samples collected by the improved policy π_{t+1} at each iteration. Fig.3(b) shows the results for InvertedPendulum-v1. Results for all 7 MuJoCo tasks can be found in Appendix B.2. By comparing approximation losses before updating (red and blue solid curves), we can observe that the approximation loss of $\mathbb{V}_\theta(\chi_{\pi_{t+1}})$ is almost consistently lower than that of V_{ϕ_t} . This means that the generalized value estimates offered by PeVFA are usually closer to the true values of π_{t+1} , demonstrating the consequence arrived in Theorem 1. For the dashed curves, it shows that PeVFA $\mathbb{V}_{\theta_{t+1}}(\chi_{\pi_{t+1}})$ can achieve lower approximation loss for π_{t+1} than conventional VFA $V_{\phi_{t+1}}$ after the same number of training with the same on-policy samples. The empirical evidence above indicates that PeVFA can be preferable to the conventional VFA for consecutive value approximation. The generalized value estimates along policy improvement path have the potential to expedite the process of GPI.

3.4 Reinforcement Learning with PeVFA

Based on the results above, we expect to leverage the value generalization of PeVFA to facilitate RL. In Algorithm 1, we propose a general description of RL algorithm under the paradigm of GPI with PeVFA. For each iteration, the interaction experiences of current policy and the policy representation are stored in a buffer (line 3-4). At an interval of M iterations, PeVFA is trained via value approximation for previous policies with the stored data and the policy representation

Algorithm 1: RL under the paradigm of GPI with PeVFA ($\mathbb{V}(s, \chi_\pi)$ is used for demonstration)

```

1: Initialize policy  $\pi_0$ , policy representation model  $g$ , PeVFA  $\mathbb{V}_{-1}$ 
   and experience buffer  $\mathcal{D}$ 
2: for iteration  $t = 0, 1, \dots$  do
3:   Rollout policy  $\pi_t$  in the environment and obtain  $k$  trajectories
      $\mathcal{T}_t = \{\tau_i\}_{i=0}^k$ 
4:   Get representation  $\chi_{\pi_t} = g(\pi)$  for policy  $\pi_t$  and add experiences
      $(\chi_{\pi_t}, \mathcal{T}_t)$  in buffer  $\mathcal{D}$ 
5:   if  $t \% M = 0$  then
6:     Update PeVFA  $\mathbb{V}_{t-1}(s, \chi_{\pi_i})$  for previous policies with
       data  $\{(\chi_{\pi_i}, \mathcal{T}_i)\}_{i=0}^{t-1}$ 
7:     Update policy representation model  $g$ , e.g., with approaches
       provided in Sec. 4
8:   end if
9:   Update PeVFA  $\mathbb{V}_{t-1}(s, \chi_{\pi_t})$  for current policy  $\chi_{\pi_t}$  and set
      $\mathbb{V}_t \leftarrow \mathbb{V}_{t-1}$ 
10:  Update  $\pi_t$  w.r.t  $\mathbb{V}_t(s, \chi_{\pi_t})$  by policy improvement algorithm
    and set  $\pi_{t+1} \leftarrow \pi_t$ 
11: end for

```

tation model is updated according to the method used (line 5-8). This part is unique to PeVFA for preservation and generalization of knowledge of historical policies. Next, value approximation for current policy is performed with PeVFA (line 9). A key difference here is that the generalized value estimates (i.e., $\mathbb{V}_{t-1}(\chi_{\pi_t})$) are used as start points. Afterwards, a successive policy is obtained from typical policy improvement (line 10). Algorithm 1 can be implemented in different ways and we propose an instance implemented based on PPO (Schulman et al. 2017) in our experiments later. We refer the readers to Appendix C for more discussions on GPI with PeVFA. In the next section, we introduce our methods for policy representation learning.

4 Policy Representation Learning

To derive practical deep RL algorithms, one key point is policy representation, i.e., a low-dimensional embedding of RL policy. Intuitively, the choice of policy representation influences the approximation and generalization of PeVFA. To our knowledge, it remains unclear how effective representation for general RL policies can be obtained in practice. In previous section, we demonstrate the effectiveness of using policy parameters as a naive representation, called RPR, when policy network is small. However, a usual policy network may have large number of parameters, thus making it inefficient and even irrational to use RPR for approximation and generalization. More generally, policy parameters of the policy we wish to represent may not be accessible.

To this end, we propose a general framework of policy representation learning as illustrated in Fig.4. The first thing to consider is data source, i.e., from which we can extract the information for an effective policy representation. Recall that the policy is a distribution over state and action space of high dimensionality. The features of such a distribution is not directly available. Therefore, we consider two kinds of data source below that indirectly contains the information of policies: 1) *Surface Policy Representation (SPR)*: The first data source is state-action pairs (or trajectories (Grover et al.

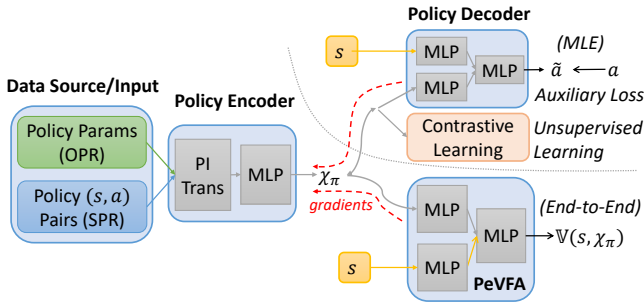


Figure 4: The framework of policy representation training. Policy network parameters used for OPR or policy state-action pairs used for SPR are fed into policy encoder with permutation-invariant (PI) transformations followed by an MLP, producing the representation χ_π . Afterwards, χ_π can be trained by gradients from the value approximation loss of PeVFA (i.e., End-to-End), as well as the auxiliary loss of policy recovery or contrastive learning (i.e., InfoNCE) loss.

2018)), since they reflect the behaviors of policy. This data source is general since no explicit form of policy is assumed. In a geometric view, learning policy representation from state-action pairs can be viewed as extracting the features of policy via scattering sample points on the curved surface of policy distribution. 2) *Origin Policy Representation (OPR)*: The other data source is parameters of policy since they determine the underlying form of policy distribution. Such a data source is often available during the learning process of deep RL algorithms when policy is parameterized by neural networks. Generally, we consider a policy network to be an MLP with well represented state features (e.g., features extracted by CNN for pixels or by LSTM for sequences) as input.

The remaining question is how we extract the policy representation from the data sources mentioned above. As shown in Fig.4, we use permutation-invariant (PI) transformations followed by an MLP to encode the data of policy π into an embedding χ_π for both SPR and OPR. For SPR, each state-action pair of $\{(s_i, a_i)\}_{i=1}^k$ is fed into a common MLP, followed by a Mean-Reduce operation on the outputted features across k . For OPR, we perform PI transformation (similar as done for state-action pairs) inner-layer weights and biases $\{(w_i, b_i)\}_{i=1}^h$ for each layer first, where h denotes the number of nodes in this layer and w_i, b_i is the income weight vector from previous layer and the bias of i th node; then we concatenate encoding of layers and obtain the OPR. An illustration for the encoding of OPR is in Fig.12 of Appendix.

To train the policy embedding χ_π obtained above, the most straightforward way is to backpropagate the value approximation loss of PeVFA in an *End-to-End (E2E)* fashion as illustrated on the lower-right of Fig.4. In addition, we provide two self-supervised training losses for both OPR and SPR, as illustrated on the upper-right of Fig.4. The first one is an *auxiliary loss (AUX)* of policy recovery (Grover et al. 2018), i.e., to recover the action distributions of π from χ_π under different states. To be specific, an auxiliary policy decoder $\tilde{\pi}(\cdot|s, \chi_\pi)$ is trained through behav-

ioral cloning, formally to minimize cross-entropy objective $\mathcal{L}_{\text{AUX}} = -\mathbb{E}_{(s,a)} [\log \tilde{\pi}(a|s, \chi_\pi)]$. For the second one, we propose to train χ_π by *Contrastive Learning (CL)* (Srinivas, Laskin, and Abbeel 2020; Schwarzer et al. 2020): policies are encouraged to be close to similar ones (i.e., positive samples π^+), and to be apart from different ones (i.e., negative samples π^-) in representation space. For each policy, we construct positive samples by data augmentation on policy data, depending on SPR or OPR considered; and different policies along the policy improvement path naturally provide negative samples for each other. Finally, the embedding χ_π is optimized through minimizing the InfoNCE loss (Oord, Li, and Vinyals 2018) below: $\mathcal{L}_{\text{CL}} = -\mathbb{E}_{(\pi^+, \{\pi^-\})} \left[\log \frac{\exp(\chi_{\pi^+}^\top W \chi_{\pi^+})}{\exp(\chi_{\pi^+}^\top W \chi_{\pi^+}) + \sum_{\pi^-} \exp(\chi_{\pi^+}^\top W \chi_{\pi^-})} \right]$.

Now, the training of policy representation in Algorithm 1 can be performed with any combination of data sources and training losses provided above. The pseudo-code of the overall policy representation training framework is shown in Algorithm 4 in Appendix D.4. In addition, complete implementation details and more discussions (e.g., on the scalability, representation criteria) are provided in Appendix D.

5 Experiments

In this section, we focus on the following questions:

- **Question 1:** Can value generalization offered by PeVFA improve a deep RL algorithm in practice?
- **Question 2:** Can our proposed framework to learn effective policy representation?

Our experiments are conducted in several OpenAI Gym continuous control tasks (1 from Box2D and 5 from MuJoCo) (Brockman et al. 2016; Todorov, Erez, and Tassa 2012). All experimental details and curves can be found in Appendix B.

Algorithm Implementation. We use PPO (Schulman et al. 2017) as the basic algorithm and propose a representative implementation of Algorithm 1, called **PPO-PeVFA**. PPO is a policy optimization algorithm that follows the paradigm of GPI (Fig.1, left). A value network $V_\phi(s)$ with parameters ϕ (i.e., conventional VFA) is trained to approximate the value of current policy π ; while π is optimized with respect to a proximal surrogate objective using advantages calculated by V_ϕ and GAE (Schulman et al. 2016). Compared with original PPO, PPO-PeVFA makes use of a PeVFA network $\mathbb{V}_\theta(s, \chi_\pi)$ with parameters θ rather than the conventional VFA $V_\phi(s)$, and follows the training scheme as in Algorithm 1. Note PPO-PeVFA has the same policy optimization as original PPO and only differs at value approximation.

Baselines and Variants. Except for original PPO as a default baseline, we use another two baselines: 1) PPO-PeVFA with randomly generated policy representation for each policy, denoted by **Ran PR**; 2) PPO-PeVFA with Raw Policy Representation (**RPR**), i.e., use the vectorized policy network parameters as in PVFs (Faccio, Kirsch, and Schmidhuber 2021). Our variants of PPO-PeVFA differ at the policy representation used. In total, we consider 6 variants denoted by the combination of the policy data choice (i.e., **OPR**, **SPR**) and representation principle choice (i.e., **E2E**, **CL**, **AUX**).

Environments	Benchmarks			Origin Policy Representation (Ours)			Surface Policy Representation (Ours)		
	PPO	Ran PR	RPR	E2E	CL	AUX	E2E	CL	AUX
HalfCheetah-v1	2621	2470	2325	3171 \pm 427	3725 \pm 348	3175 \pm 517	2774 \pm 233	3349 \pm 341	3216 \pm 506
Hopper-v1	1639	1226	1097	2085 \pm 310	2351 \pm 231	2214 \pm 360	2227 \pm 297	2392 \pm 263	2577 \pm 217
Walker2d-v1	1505	1269	317	1856 \pm 305	2038 \pm 315	2044 \pm 316	1930.57 \pm 456	2203 \pm 381	1980 \pm 325
Ant-v1	2835	2742	2143	3581 \pm 185	4019 \pm 162	3784 \pm 268	3173 \pm 184	3632 \pm 134	3397 \pm 200
InvDouPend-v1	9344	9355	8856	9357 \pm 0.29	9355 \pm 0.64	9355 \pm 0.68	9355 \pm 0.89	9356 \pm 0.96	9355 \pm 1.42
LunarLander-v2	219	226	-22	238 \pm 3.37	239 \pm 3.70	234 \pm 3.47	236 \pm 3.13	234 \pm 3.13	235 \pm 5.70

Table 1: Average returns (\pm half a std) over 10 trials for algorithms. Each result is the maximum average evaluation along the training process. Top two values for each environment are bold.

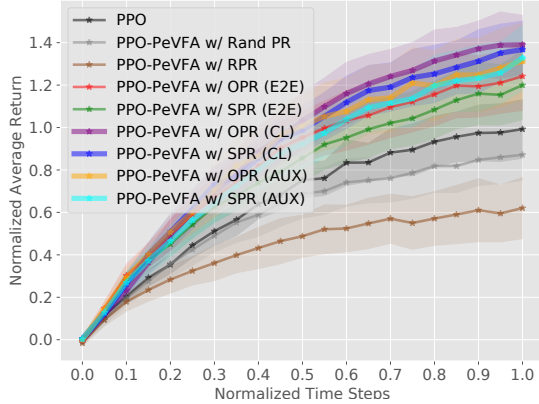


Figure 5: Averaged returns (normalized by the maximum of PPO) aggregated over 4 MuJoCo tasks.

Experimental Details. For all baselines and variants, we use a normal-scale policy network with 2 layers and 64 units for each layer, resulting in over 3k~10k (e.g., Ant-v1) policy parameters depending on the environments. We do not assume pre-collected policies. Thus the size of policy set increases from 1 (i.e., the initial policy) during the learning process, to about 1k~2k for a single trial. The dimensionality of all kinds of policy representation except for RPR is set to 64. The buffer D maintains recent 200k steps of experience and the policy data of corresponding policy. The number of interaction step of each trial is 1M for InvDouPend-v1 and LunarLander-v2, 4M for Ant-v1 and 2M for the others.

Results. The overall experimental results are summarized in Tab.1. In Fig.5, we provide aggregated results across all environments except for InvDouPend-v1 and LunarLander-v2 (since most algorithms achieve near-optimal results), where all returns are normalized by the results of PPO in Tab.1. All learning curves can be found in Appendix F.2.

To Question 1: From Tab.1, we can find that both PPO-PeVFA w/ OPR (E2E) and PPO-PeVFA w/ SPR (E2E) outperforms PPO in all 6 tasks, and achieve over 20% improvement in Fig.5. This demonstrates the effectiveness of PeVFA. Moreover, the improvement is further enlarged (to about 40%) by CL and AUX for both OPR and SPR. This indicates that the superiority of PeVFA can be further utilized with better policy representation that offers a more suitable space for value generalization.

To Question 2: In Tab.1, consistent degeneration is observed for PPO-PeVFA w/ Ran PR due to the negative effects on generalization caused by the randomness and disorder of policy representation. This phenomenon seems to be more severe for PPO-PeVFA w/ RPR due to the complexity of high-dimensional parameter space. In contrast, the improvement achieved by our proposed PPO-PeVFA variants shows that effective policy representation can be learned from policy parameters (OPR) and state-action pairs (SPR) though value approximation loss (i.e., E2E) and further improved when additional self-supervised representation learning is involved as CL and AUX. Overall, OPR slightly outperforms SPR as CL does over AUX. We hypothesize that it is due to the stochasticity of state-action pairs which serve as inputs of SPR and training samples for AUX. This reveals the space for future improvement.

In addition, we visualize the learned representation in Fig.18 and 19 in Appendix F.3. We can observe that policies from different trials are locally continuous and show different modes of embedding trajectories due to random initialization and optimization; while a global evolvment among trials emerges with respect to policy performance.

6 Conclusion and Future Work

In this paper, we propose Policy-extended Value Function Approximator (PeVFA) and study value generalization among policies. We propose a new form of GPI based on PeVFA which is potentially preferable to conventional VFA for value approximation. Moreover, we propose a general framework to learn low-dimensional embedding of RL policy. Our experiments demonstrate the effectiveness of the generalization characteristic of PeVFA and our proposed policy representation learning methods.

Our work opens up some research directions on value generalization among policies and policy representation. A possible future study on the theory of value generalization among policies is to consider the interplay between approximation error, policy improvement and local generalization during GPI with PeVFA. Besides, analysis on influence factors of value generalization among policies (e.g., policy representation, architecture of PeVFA) and other utilization of PeVFA are expected. For better policy representation, inspirations on OPR may be got from studies on Manifold Hypothesis of neural network; the selection of more informative state-action pairs for SPR is also worth research.

Acknowledgments

The work is supported by the National Natural Science Foundation of China (Grant No.: U1836214), and the New Generation of Artificial Intelligence Science and Technology Major Project of Tianjin (Grant No.: 19ZXZNGX00010).

References

- Andrychowicz, M.; Crow, D.; Ray, A.; Schneider, J.; Fong, R.; Welinder, P.; McGrew, B.; Tobin, J.; Abbeel, P.; and Zaremba, W. 2017. Hindsight Experience Replay. In *NeurIPS*, 5048–5058.
- Andrychowicz, M.; Raichuk, A.; Stanczyk, P.; Orsini, M.; Girgin, S.; Marinier, R.; Hussenot, L.; Geist, M.; Pietquin, O.; Michalski, M.; Gelly, S.; and Bachem, O. 2020. What Matters In On-Policy Reinforcement Learning? A Large-Scale Empirical Study. *CoRR*, abs/2006.05990.
- Arnekjvist, I.; Kragic, D.; and Stork, J. A. 2019. VPE: Variational Policy Embedding for Transfer Reinforcement Learning. In *ICRA*, 36–42.
- Bellemare, M. G.; Dabney, W.; and Munos, R. 2017. A Distributional Perspective on Reinforcement Learning. In *ICML*, volume 70, 449–458.
- Bertsekas, D. P.; and Tsitsiklis, J. N. 1996. *Neuro-dynamic programming*, volume 3 of *Optimization and neural computation series*. Athena Scientific.
- Brockman, G.; Cheung, V.; Pettersson, L.; Schneider, J.; Schulman, J.; Tang, J.; and Zaremba, W. 2016. OpenAI Gym. *CoRR*, abs/1606.01540.
- Chen, T.; Kornblith, S.; Norouzi, M.; and Hinton, G. E. 2020. A Simple Framework for Contrastive Learning of Visual Representations. *CoRR*, abs/2002.05709.
- Cobbe, K.; Hilton, J.; Klimov, O.; and Schulman, J. 2021. Phasic Policy Gradient. In *ICML*, volume 139 of *Proceedings of Machine Learning Research*, 2020–2027.
- Degrís, T.; White, M.; and Sutton, R. S. 2012. Off-Policy Actor-Critic. *CoRR*, abs/1205.4839.
- D’Oro, P.; and Jaskowski, W. 2020. How to Learn a Useful Critic? Model-based Action-Gradient-Estimator Policy Optimization. In *NIPS*.
- Engstrom, L.; Ilyas, A.; Santurkar, S.; Tsipras, D.; Janoos, F.; Rudolph, L.; and Madry, A. 2020. Implementation Matters in Deep RL: A Case Study on PPO and TRPO. In *ICLR*.
- Faccio, F.; Kirsch, L.; and Schmidhuber, J. 2021. Parameter-Based Value Functions. In *ICLR*.
- Fu, H.; Tang, H.; Hao, J.; Chen, C.; Feng, X.; Li, D.; and Liu, W. 2020. Towards Effective Context for Meta-Reinforcement Learning: an Approach based on Contrastive Learning. *CoRR*, abs/2009.13891.
- Fujimoto, S.; v. Hoof, H.; and Meger, D. 2018. Addressing Function Approximation Error in Actor-Critic Methods. In *ICML*.
- Gaier, A.; Asteroth, A.; and Mouret, J. 2020. Discovering representations for black-box optimization. In *GECCO*, 103–111.
- Grosnit, A.; Tutunov, R.; Maraval, A.; Griffiths, R.; Cowen-Rivers, A.; Yang, L.; Zhu, L.; Lyu, W.; Chen, Z.; Wang, J.; Peters, J.; and Bou-Ammar, H. 2021. High-Dimensional Bayesian Optimisation with Variational Autoencoders and Deep Metric Learning. *CoRR*, abs/2106.03609.
- Grover, A.; Al-Shedivat, M.; Gupta, J. K.; Burda, Y.; and Edwards, H. 2018. Learning Policy Representations in Multiagent Systems. In *ICML*, volume 80, 1797–1806.
- Haarnoja, T.; Zhou, A.; Abbeel, P.; and Levine, S. 2018. Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor. In *ICML*, 1856–1865.
- Hafner, D.; Lillicrap, T. P.; Ba, J.; and Norouzi, M. 2020. Dream to Control: Learning Behaviors by Latent Imagination. In *ICLR*.
- Harb, J.; Schaul, T.; Precup, D.; and Bacon, P. 2020. Policy Evaluation Networks. *CoRR*, abs/2002.11833.
- Hausman, K.; Springenberg, J. T.; Wang, Z.; Heess, N.; and Riedmiller, M. A. 2018. Learning an Embedding Space for Transferable Robot Skills. In *ICLR*.
- He, H.; and Boyd-Graber, J. L. 2016. Opponent Modeling in Deep Reinforcement Learning. In *ICML*, volume 48, 1804–1813.
- He, K.; Fan, H.; Wu, Y.; Xie, S.; and Girshick, R. B. 2020. Momentum Contrast for Unsupervised Visual Representation Learning. In *CVPR*, 9726–9735.
- Igl, M.; Farquhar, G.; Luketina, J.; Boehmer, W.; and Whiteson, S. 2020. The Impact of Non-stationarity on Generalisation in Deep Reinforcement Learning. *CoRR*, abs/2006.05826.
- Keskar, N. S.; Mudigere, D.; Nocedal, J.; Smelyanskiy, M.; and Tang, P. T. P. 2017. On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima. In *ICLR*.
- Kingma, D. P.; and Welling, M. 2014. Auto-Encoding Variational Bayes. In *ICLR*.
- Kostrikov, I.; Yarats, D.; and Fergus, R. 2020. Image Augmentation Is All You Need: Regularizing Deep Reinforcement Learning from Pixels. *CoRR*, abs/2004.13649.
- Kuznetsov, A.; Shvechikov, P.; Grishin, A.; and Vetrov, D. P. 2020. Controlling Overestimation Bias with Truncated Mixture of Continuous Distributional Quantile Critics. In *ICML*, volume 119, 5556–5566.
- Lagoudakis, M. G.; and Parr, R. 2003. Least-Squares Policy Iteration. *J. Mach. Learn. Res.*, 4: 1107–1149.
- Lan, Q.; Pan, Y.; Fyshe, A.; and White, M. 2020. Maxmin Q-learning: Controlling the Estimation Bias of Q-learning. In *ICLR*.
- Laskin, M.; Lee, K.; Stooke, A.; Pinto, L.; Abbeel, P.; and Srinivas, A. 2020. Reinforcement Learning with Augmented Data. *CoRR*, abs/2004.14990.
- Lee, K.; Seo, Y.; Lee, S.; Lee, H.; and Shin, J. 2020. Context-aware Dynamics Model for Generalization in Model-Based Reinforcement Learning. *CoRR*, abs/2005.06800.
- Levy, A.; Konidaris, G. D.; Jr., R. P.; and Saenko, K. 2019. Learning Multi-Level Hierarchies with Hindsight. In *ICLR*.
- Lillicrap, T. P.; Hunt, J. J.; Pritzel, A.; Heess, N.; Erez, T.; Tassa, Y.; Silver, D.; and Wierstra, D. 2015. Continuous control with deep reinforcement learning. In *ICLR*.
- Liu, H.; Simonyan, K.; and Yang, Y. 2019. DARTS: Differentiable Architecture Search. In *ICLR*.
- Luo, R.; Tian, F.; Qin, T.; Chen, E.; and Liu, T. 2018. Neural Architecture Optimization. In *NeurIPS*, 7827–7838.
- Maaten, L. V. D.; and Hinton, G. E. 2008. Visualizing Data using t-SNE. *Journal of Machine Learning Research*, 9: 2579–2605.
- Mnih, V.; Badia, A. P.; Mirza, M.; Graves, A.; Lillicrap, T. P.; Harley, T.; Silver, D.; and Kavukcuoglu, K. 2016. Asynchronous Methods for Deep Reinforcement Learning. In *ICML*.
- Mnih, V.; Kavukcuoglu, K.; Silver, D.; Rusu, A. A.; Veness, J.; Bellemare, M. G.; Graves, A.; Riedmiller, M. A.; Fidjeland, A.; Ostrovski, G.; Petersen, S.; Beattie, C.; Sadik, A.; Antonoglou, I.; King, H.; Kumaran, D.; Wierstra, D.; Legg, S.; and Hassabis, D. 2015. Human-level control through deep reinforcement learning. *Nature*, 518(7540): 529–533.

- Nachum, O.; Gu, S.; Lee, H.; and Levine, S. 2018. Data-Efficient Hierarchical Reinforcement Learning. *CoRR*, abs/1805.08296.
- Nachum, O.; Gu, S.; Lee, H.; and Levine, S. 2019. Near-Optimal Representation Learning for Hierarchical Reinforcement Learning. In *ICLR*.
- Nesterov, Y. E.; and Polyak, B. T. 2006. Cubic regularization of Newton method and its global performance. *Math. Program.*, 108(1): 177–205.
- Notin, P.; Hernández-Lobato, J.; and Gal, Y. 2021. Improving black-box optimization in VAE latent space using decoder uncertainty. *CoRR*, abs/2107.00096.
- Novak, R.; Bahri, Y.; Abolafia, D. A.; Pennington, J.; and Sohl-Dickstein, J. 2018. Sensitivity and Generalization in Neural Networks: an Empirical Study. In *ICLR*.
- Oord, A.; Li, Y.; and Vinyals, O. 2018. Representation Learning with Contrastive Predictive Coding. *CoRR*, abs/1807.03748.
- Raileanu, R.; Goldstein, M.; Szlam, A.; and Fergus, R. 2020. Fast Adaptation via Policy-Dynamics Value Functions. *CoRR*, abs/2007.02879.
- Rakelly, K.; Zhou, A.; Finn, C.; Levine, S.; and Quillen, D. 2019. Efficient Off-Policy Meta-Reinforcement Learning via Probabilistic Context Variables. In *ICML*, volume 97, 5331–5340.
- Rakicevic, N.; Cully, A.; and Kormushev, P. 2021. Policy Manifold Search: Exploring the Manifold Hypothesis for Diversity-based Neuroevolution. *CoRR*, abs/2104.13424.
- Schaul, T.; Horgan, D.; Gregor, K.; and Silver, D. 2015. Universal Value Function Approximators. In *ICML*, volume 37, 1312–1320.
- Scherrer, B.; Ghavamzadeh, M.; Gabillon, V.; Lesner, B.; and Geist, M. 2015. Approximate modified policy iteration and its application to the game of Tetris. *J. Mach. Learn. Res.*, 16: 1629–1676.
- Schreck, J. S.; Coley, C. W.; and Bishop, K. J. 2019. Learning retrosynthetic planning through simulated experience. *ACS central science*, 5(6): 970–981.
- Schulman, J.; Levine, S.; Abbeel, P.; Jordan, M. I.; and Moritz, P. 2015. Trust Region Policy Optimization. In *ICML*, 1889–1897.
- Schulman, J.; Moritz, P.; Levine, S.; Jordan, M. I.; and Abbeel, P. 2016. High-Dimensional Continuous Control Using Generalized Advantage Estimation. In *ICLR*.
- Schulman, J.; Wolski, F.; Dhariwal, P.; Radford, A.; and Klimov, O. 2017. Proximal Policy Optimization Algorithms. *CoRR*, abs/1707.06347.
- Schwarzer, M.; Anand, A.; Goel, R.; Hjelm, R. D.; Courville, A. C.; and Bachman, P. 2020. Data-Efficient Reinforcement Learning with Momentum Predictive Representations. *CoRR*, abs/2007.05929.
- Silver, D.; Huang, A.; Maddison, C. J.; Guez, A.; Sifre, L.; van den Driessche, G.; Schrittwieser, J.; Antonoglou, I.; Panneershelvam, V.; Lanctot, M.; Dieleman, S.; Grewe, D.; Nham, J.; Kalchbrenner, N.; Sutskever, I.; Lillicrap, T. P.; Leach, M.; Kavukcuoglu, K.; Graepel, T.; and Hassabis, D. 2016. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587): 484–489.
- Silver, D.; Lever, G.; Heess, N.; Degris, T.; Wierstra, D.; and Riedmiller, M. A. 2014. Deterministic Policy Gradient Algorithms. In *ICML*, 387–395.
- Srinivas, A.; Laskin, M.; and Abbeel, P. 2020. CURL: Contrastive Unsupervised Representations for Reinforcement Learning. *CoRR*, abs/2004.04136.
- Sutton, R. S.; and Barto, A. G. 1998. *Reinforcement learning - an introduction*. Adaptive computation and machine learning. MIT Press.
- Sutton, R. S.; Modayil, J.; Delp, M.; Degris, T.; Pilarski, P. M.; White, A.; and Precup, D. 2011. Horde: a scalable real-time architecture for learning knowledge from unsupervised sensorimotor interaction. In *AAMAS*, 761–768.
- Tacchetti, A.; Song, H. F.; Mediano, P. A. M.; Zambaldi, V. F.; Kramár, J.; Rabinowitz, N. C.; Graepel, T.; Botvinick, M.; and Battaglia, P. W. 2019. Relational Forward Models for Multi-Agent Learning. In *ICLR*.
- Todorov, E.; Erez, T.; and Tassa, Y. 2012. MuJoCo: A physics engine for model-based control. *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 5026–5033.
- Tsai, Y. H.; Wu, Y.; Salakhutdinov, R.; and Morency, L. 2020. Demystifying Self-Supervised Learning: An Information-Theoretical Framework. *CoRR*, abs/2006.05576.
- Unterthiner, T.; Keyser, D.; Gelly, S.; Bousquet, O.; and Tolstikhin, I. O. 2020. Predicting Neural Network Accuracy from Weights. *CoRR*, abs/2002.11448.
- v. Hasselt, H. 2010. Double Q-learning. In Lafferty, J. D.; Williams, C. K. I.; Shawe-Taylor, J.; Zemel, R. S.; and Culotta, A., eds., *NeurIPS*, 2613–2621.
- Vinyals, O.; Babuschkin, I.; Czarnecki, W. M.; Mathieu, M.; Dudzik, A.; Chung, J.; Choi, D. H.; Powell, R.; Ewalds, T.; Georgiev, P.; Oh, J.; Horgan, D.; Kroiss, M.; Danihelka, I.; Huang, A.; Sifre, L.; Cai, T.; Agapiou, J. P.; Jaderberg, M.; Vezhnevets, A. S.; Leblond, R.; Pohlen, T.; Dalibard, V.; Budden, D.; Sulsky, Y.; Molloy, J.; Paine, T. L.; Gulcehre, C.; Wang, Z.; Pfaff, T.; Wu, Y.; Ring, R.; Yogatama, D.; Wünsch, D.; McKinney, K.; Smith, O.; Schaul, T.; Lillicrap, T.; Kavukcuoglu, K.; Hassabis, D.; Apps, C.; and Silver, D. 2019. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature*, 575(7782): 350–354.
- Wang, H.; Kesar, N. S.; Xiong, C.; and Socher, R. 2018. Identifying Generalization Properties in Neural Networks. *CoRR*, abs/1809.07402.
- Wang, R.; Yu, R.; An, B.; and Rabinovich, Z. 2020. I²HRL: Interactive Influence-based Hierarchical Reinforcement Learning. In *IJCAI*, 3131–3138.
- You, J.; Liu, B.; Ying, Z.; Pande, V. S.; and Leskovec, J. 2018. Graph Convolutional Policy Network for Goal-Directed Molecular Graph Generation. In *NeurIPS 2018*, 6412–6422.