

Collaborative Deep Reinforcement Learning

Kaixiang Lin
Computer Science and Engineering
Michigan State University
428 S Shaw Ln.
East Lansing, MI 48824
linkaixi@msu.edu

Shu Wang
Computer Science
Rutgers University
57 US Highway 1
New Brunswick, NJ 088901
sw498@cs.rutgers.edu

Jiayu Zhou
Computer Science and Engineering
Michigan State University
428 S Shaw Ln.
East Lansing, MI 48824
jiayuz@msu.edu

ABSTRACT

Besides independent learning, human learning process is highly improved by summarizing what has been learned, communicating it with peers, and subsequently fusing knowledge from different sources to assist the current learning goal. This *collaborative learning* procedure ensures that the knowledge is shared, continuously refined, and concluded from different perspectives to construct a more profound understanding. The idea of knowledge transfer has led to many advances in machine learning and data mining, but significant challenges remain, especially when it comes to reinforcement learning, heterogeneous model structures, and different learning tasks. Motivated by human collaborative learning, in this paper we propose a collaborative deep reinforcement learning (CDRL) framework that performs adaptive knowledge transfer among heterogeneous learning agents. Specifically, the proposed CDRL conducts a novel deep knowledge distillation method to address the heterogeneity among different learning tasks with a deep alignment network. Furthermore, we present an efficient collaborative Asynchronous Advantage Actor-Critic (cA3C) algorithm to incorporate deep knowledge distillation into the online training of agents, and demonstrate the effectiveness of the CDRL framework using extensive empirical evaluation on OpenAI gym.

CCS CONCEPTS

•Computing methodologies →Machine learning; Reinforcement learning; Transfer learning;

KEYWORDS

Knowledge distillation; Transfer learning; Deep reinforcement learning

ACM Reference format:

Kaixiang Lin, Shu Wang, and Jiayu Zhou. 1997. Collaborative Deep Reinforcement Learning. In *Proceedings of ACM Woodstock conference, El Paso, Texas USA, July 1997 (WOODSTOCK'97)*, 9 pages.
DOI: 10.475/123.4

1 INTRODUCTION

It is the development of cognitive abilities including learning, remembering, communicating that enables human to conduct social

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

WOODSTOCK'97, El Paso, Texas USA

© 2016 Copyright held by the owner/author(s). 123-4567-24-567/08/06...\$15.00
DOI: 10.475/123.4

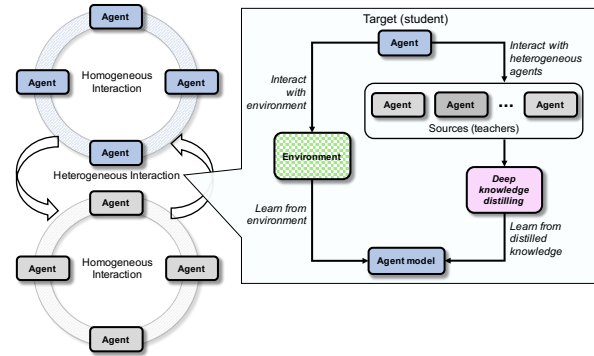


Figure 1: The illustration of Collaborative Deep Reinforcement Learning Framework.

cooperation, which is the key to the rise of humankind. As a social animal, the ability to collaborate awoke the cognitive revolution and reveals the prosperous history of human [14]. In disciplines of cognitive science, education and psychology, *collaborative learning*, a situation in which a group of people learn to achieve a set of tasks together, has been advocated throughout previous studies [9]. It is intuitive to illustrate the concept of collaborative learning in the example of group study. A group of students are studying together to master some challenging course materials. As each student may understand the materials from a distinctive perspective, effective communication would greatly help the entire group achieve a better understanding than those from independent study, and could significantly improve the efficiency and effectiveness of learning process, as well [12].

On the other hand, the study of human learning has largely advanced the design of machine learning and data mining algorithms, especially in reinforcement learning and transfer learning. The recent success of deep reinforcement learning (DRL) has attracted increasing attention from the community, as DRL can discover very competitive strategies by having learning agents interacting with a given environment and using rewards from the environment as the supervision (e.g., [16, 18, 20, 28]). Even though most of current research on DRL has focused on learning from games, it possesses great transformative power to impact many industries with data mining and machine learning techniques such as clinical decision support [32], marketing [2], finance [1], visual navigation [37], and autonomous driving [8]. Although there are many existing efforts towards effective algorithms for DRL [19, 21], the computational cost still imposes significant challenges as training DRL for even a

simple game such as PONG [5] remains very expensive. The underlying reasons for the obstacle of efficient training mainly lie in two aspects: First, the supervision (rewards) from the environment is very sparse and implicit during training. It may take an agent hundreds or even thousands actions to get a single reward, and which actions that actually lead to this reward are ambiguous. Besides the insufficient supervision, training deep neural network itself takes lots of computational resources.

Due to the aforementioned difficulties, performing knowledge transfer from other related tasks or well-trained deep models to facilitate training has drawn lots of attention in the community [16, 24–26, 31]. Existing transfer learning can be categorized into two classes according to the means that knowledge is transferred: *data transfer* [15, 24, 26] and *model transfer* [10, 24, 34, 35]. Model transfer methods implement knowledge transfer from introducing inductive bias during the learning, and has been extensively studied in both transfer learning/multi-task learning (MTL) community and deep learning community. For example, in the regularized MTL models such as [11, 36], tasks with the same feature space are related through some structured regularization. Another example is the multi-task deep neural network, where different tasks share parts of the network structures [35]. One obvious disadvantage of model transfer is the lack of flexibility: usually the feasibility of inductive transfer has largely restricted the model structure of learning task, which makes it not practical in DRL because for different tasks the optimal model structures may be radically different. On the other hand, the recently developed data transfer (also known as knowledge distillation or mimic learning) [15, 24, 26] embeds the source model knowledge into data points. Then they are used as knowledge bridge to train target models, which can have different structures as compared to the source model [6, 15]. Because of the structural flexibility, the data transfer is especially suitable to deal with structure variant models.

There are two situations that transfer learning methods are essential in DRL:

Certificated heterogeneous transfer. Training a DRL agent is computational expensive. If we have a well-trained model, it will be beneficial to assist the learning of other tasks by transferring knowledge from this model. Therefore we consider following research question: Given one *certificated* task (i.e. the model is well-designed, extensively trained and performs very well), how can we maximize the information that can be used in the training of other related tasks? Some model transfer approaches directly use the weights from the trained model to initialize the new task [24], which can only be done when the model structures are the same. Thus, this strict requirement has largely limited its general applicability on DRL. On the other hand, the initialization may not work well if the tasks are significantly different from each other in nature [24]. This challenge could be partially solved by generating an intermediate dataset (logits) from the existing model to help learning the new task. However, new problems would arise when we are transferring knowledge between *heterogeneous tasks*. Not only the action spaces are different in dimension, the intrinsic action probability distributions and semantic meanings of two tasks could differ a lot. Specifically, one action in PONG may refer to move the paddle upwards while the same action index in RIVERAID [5] would correspond to fire. Therefore, the distilled dataset generated from

the trained source task cannot be directly used to train the heterogeneous target task. In this scenario, the first key challenge we identified in this work is that how to conduct data transfer among heterogeneous tasks so that we can maximally utilize the information from a certificated model while still maintain the flexibility of model design for new tasks. During the transfer, the transferred knowledge from other tasks may contradict to the knowledge that agents learned from its environment. One recently work [25] use an attention network selective eliminate transfer if the contradiction presents, which is not suitable in this setting since we are given a certificated task to transfer. Hence, the second challenge is how to resolve the conflict and perform a meaningful transfer.

Lack of expertise. A more general desired but also more challenging scenario is that DRL agents are trained for multiple heterogeneous tasks without any pre-trained models available. One feasible way to conduct transfer under this scenario is that agents of multiple tasks share part of their network parameters [26, 35]. However, an inevitable drawback is, multiple models lose their task-specific designs since the shared part needs to be the same. Another solution is to learn a domain invariant feature space shared by all tasks [3]. However, some task-specific information is often lost while converting the original state to a new feature subspace. In this case, an intriguing questions is that: can we design a framework that fully utilizes the original environment information and meanwhile leverages the knowledge transferred from other tasks?

This paper investigates the aforementioned problems systematically and proposes a novel Collaborative Deep Reinforcement Learning (CDRL) framework (illustrated in Figure 1) to resolve them. Our major contribution is threefold:

- First, in order to transfer knowledge among heterogeneous tasks while remaining the task-specific design of model structure, a novel deep knowledge distillation is proposed to address the heterogeneity among tasks, with the utilization of deep alignment network designed for the domain adaptation.
- Second, in order to incorporate the transferred knowledge from heterogeneous tasks into the online training of current learning agents, similar to human collaborative learning, an efficient collaborative asynchronously advantage actor-critic learning (cA3C) algorithm is developed under the CDRL framework. In cA3C, the target agents are able to learn from environments and its peers simultaneously, which also ensure the information from original environment is sufficiently utilized. Further, the knowledge conflict among different tasks is resolved by adding an extra distillation layer to the policy network under CDRL framework, as well.
- Last but not least we present extensive empirical studies on OpenAI gym to evaluate the proposed CDRL framework and demonstrate its effectiveness by achieving more than 10% performance improvement compared to the current state-of-the-art.

Notations: In this paper, we use teacher network/source task denotes the network/task contained the knowledge to be transferred to others. Similarly, the student network/target task is referred to those tasks utilizing the knowledge transferred from others to facilitate its own training. The expert network denotes the network that has already reached a relative high averaged reward in its own

environment. In DRL, an agent is represented by a policy network and a value network that share a set of parameters. Homogeneous agents denotes agents that perform and learn under independent copies of same environment. Heterogeneous agents refer to those agents that are trained in different environments.

2 RELATED WORK

Multi-agent learning. One closely related area to our work is multi-agent reinforcement learning. A multi-agent system includes a set of agents interacting in one environment. Meanwhile they could potentially interact with each other [7, 13, 17, 30]. In collaborative multi-agent reinforcement learning, agents work together to maximize a shared reward measurement [13, 17]. There is a clear distinction between the proposed CDRL framework and multi-agent reinforcement learning. In CDRL, each agent interacts with its own environment copy and the goal is to maximize the reward of the target agents. The formal definition of the proposed framework is given in Section 4.1.

Transfer learning. Another relevant research topic is domain adaption in the field of transfer learning [23, 29, 33]. The authors in [29] proposed a two-stage domain adaptation framework that considers the differences among marginal probability distributions of domains, as well as conditional probability distributions of tasks. The method first re-weights the data from the source domain using Maximum Mean Discrepancy and then re-weights the predictive function in the source domain to reduce the difference on conditional probabilities. In [33], the marginal distributions of the source and the target domain are aligned by training a network, which maps inputs into a domain invariant representation. Also, knowledge distillation was directly utilized to align the source and target class distribution. One clear limitation here is that the source domain and the target domain are required to have the same dimensionality (i.e. number of classes) with same semantics meanings, which is not the case in our deep knowledge distillation.

In [3], an invariant feature space is learned to transfer skills between two agents. However, projecting the state into a feature space would lose information contained in the original state. There is a trade-off between learning the common feature space and preserving the maximum information from the original state. In our work, we use data generated by intermediate outputs in the knowledge transfer instead of a shared space. Our approach thus retains complete information from the environment and ensures high quality transfer. The recently proposed A2T approach [25] can avoid negative transfer among different tasks. However, it is possible that some negative transfer cases may be because of the inappropriate design of transfer algorithms. In our work, we show that we can perform successful transfer among tasks that seemingly cause negative transfer.

Knowledge transfer in deep learning. Since the training of each agent in an environment can be considered as a learning task, and the knowledge transfer among multiple tasks belongs to the study of multi-task learning. **The multi-task deep neural network (MTDNN) [35] transfers knowledge among tasks by sharing parameters of several low-level layers.** Since the low-level layers can be considered to perform representation learning, the MTDNN is learning a shared representation for inputs, which is then used

by high-level layers in the network. Different learning tasks are related to each other via this shared feature representation. In the proposed CDRL, we do not use the shared representation due to the inevitable information loss when we project the inputs into a shared representation. We instead perform explicitly knowledge transfer among tasks by distilling knowledge that are independent of model structures. In [15], the authors proposed to compress cumbersome models (teachers) to more simple models (students), where the simple models are trained by a dataset (knowledge) distilled from the teachers. However, this approach cannot handle the transfer among heterogeneous tasks, which is one key challenge we addressed in this paper.

Knowledge transfer in deep reinforcement learning. Knowledge transfer is also studied in deep reinforcement learning. [19] proposed multi-threaded asynchronous variants of several most advanced deep reinforcement learning methods including Sarsa, Q-learning, Q-learning and advantage actor-critic. Among all those methods, asynchronous advantage actor-critic (A3C) achieves the best performance. Instead of using experience replay as in previous work, A3C stabilizes the training procedure by training different agents in parallel using different exploration strategies. This was shown to converge much faster than previous methods and use less computational resources. We show in Section 4.1 that the A3C is subsumed to the proposed CDRL as a special case. In [24], a single multi-task policy network is trained by utilizing a set of expert Deep Q-Network (DQN) of source games. At this stage, the goal is to obtain a policy network that can play source games as close to experts as possible. The second step is to transfer the knowledge from source tasks to a new but related target task. The knowledge is transferred by using the DQN in last step as the initialization of the DQN for the new task. As such, the training time of the new task can be significantly reduced. Different from their approach, the proposed transfer strategy is not to directly mimic experts' actions or initialize by a pre-trained model. In [26], knowledge distillation was adopted to train a multi-task model that outperforms single task models of some tasks. The experts for all tasks are firstly acquired by single task learning. The intermediate outputs from each expert are then distilled to a similar multi-task network with an extra controller layer to coordinate different action sets. One clear limitation is that major components of the model are exactly the same for different tasks, which may lead to degraded performance on some tasks. In our work, transfer can happen even when there are no experts available. Also, our method allow each task to have their own model structures. Furthermore, even the model structures are the same for multiple tasks, the tasks are not trained to improve the performance of other tasks (i.e. it does not mimic experts from other tasks directly). Therefore our model can focus on maximizing its own reward, instead of being distracted by others.

3 BACKGROUND

3.1 Reinforcement Learning

In this work, we consider the standard reinforcement learning setting where each agent interacts with its own environment over a number of discrete time steps. Given the current state $s_t \in \mathcal{S}$ at step t , agent g_i selects an action $a_t \in \mathcal{A}$ according

to its policy $\pi(a_t|s_t)$, and receives a reward r_{t+1} from the environment. The goal of the agent is to choose an action a_t at step t that maximize the sum of future rewards $\{r_t\}$ in a decaying manner: $R_t = \sum_{i=0}^{\infty} \gamma^i r_{t+i}$, where scalar $\gamma \in (0, 1]$ is a discount rate. Based on the policy π of this agent, we can further define a state value function $V(s_t) = E[R_t|s = s_t]$, which estimates the expected discounted return starting from state s_t , taking actions following policy π until the game ends. The goal in reinforcement learning algorithm is to maximize the expected return. Since we are mainly discussing one specific agent's design and behavior throughout the paper, we leave out the notation of the agent index for conciseness.

3.2 Asynchronous Advantage actor-critic algorithm (A3C)

The asynchronous advantage actor-critic (A3C) algorithm [19] launches multiple agents in parallel and asynchronously updates a global shared target policy network $\pi(a|s, \theta_p)$ as well as a value network $V(s, \theta_v)$, parametrized by θ_p and θ_v , respectively. Each agent interacts with the environment, independently. At each step t the agent takes an action based on the probability distribution generated by policy network. After playing a n -step rollout or reaching the terminal state, the rewards are used to compute the advantage with the output of value function. The updates of policy network is conducted by applying the gradient:

$$\nabla_{\theta_p} \log \pi(a_t|s_t; \theta_p) A(s_t, a_t; \theta_v),$$

where the advantage function $A(s_t, a_t; \theta_v)$ is given by:

$$\sum_{i=0}^{T-t-1} \gamma^i r_{t+i} + \gamma^{T-t} V(s_T; \theta_v) - V(s_t; \theta_v).$$

Term T represents the step number for the last step of this rollout, it is either the max number of rollout steps or the number of steps from t to the terminal state. The update of value network is to minimize the squared difference between the environment rewards and value function outputs, i.e.,

$$\min_{\theta_v} \left(\sum_{i=0}^{T-t-1} \gamma^i r_{t+i} + \gamma^{T-t} V(s_T; \theta_v) - V(s_t; \theta_v) \right)^2.$$

The policy network and the value network share the same layers except for the last output layer. An entropy regularization of policy π is added to improve exploration, as well.

3.3 Knowledge distillation

Knowledge distillation [15] is a transfer learning approach that distills the knowledge from a teacher network to a student network using a temperature parameterized "soft targets" (i.e. a probability distribution over a set of classes). It has been shown that it can accelerate the training with less data since the gradient from "soft targets" contains much more information than the gradient obtained from "hard targets" (e.g. 0, 1 supervision).

To be more specific, logits vector $\mathbf{z} \in \mathcal{R}^d$ for d actions can be converted to a probability distribution $\mathbf{h} \in (0, 1)^d$ by a softmax function, raised with temperature τ :

$$\mathbf{h}(i) = \text{softmax}(\mathbf{z}/\tau)_i = \frac{\exp(\mathbf{z}(i)/\tau)}{\sum_j \exp(\mathbf{z}(j)/\tau)}, \quad (1)$$

where $\mathbf{h}(i)$ and $\mathbf{z}(i)$ denotes the i -th entry of \mathbf{h} and \mathbf{z} , respectively.

Then the knowledge distillation can be completed by optimize the following Kullback-Leibler divergence (KL) with temperature τ [15, 26].

$$L_{KL}(D, \theta_p^\beta) = \sum_{t=1} \text{softmax}(\mathbf{z}_t^\alpha / \tau) \ln \frac{\text{softmax}(\mathbf{z}_t^\alpha / \tau)}{\text{softmax}(\mathbf{z}_t^\beta)} \quad (2)$$

where \mathbf{z}_t^α is the logits vector from teacher network (notation α represents teacher) at step t , while \mathbf{z}_t^β is the logits vector from student network (notation β represents student) of this step. θ_p^β denotes the parameters of the student policy network. D is a set of logits from teacher network.

4 COLLABORATIVE DEEP REINFORCEMENT LEARNING FRAMEWORK

In this section, we introduce the proposed collaborative deep reinforcement learning (CDRL) framework. Under this framework, a collaborative Asynchronous Advantage Actor-Critic (cA3C) algorithm is proposed to confirm the effectiveness of the collaborative approach. Before we introduce our method in details, one underlying assumption we used is as follows:

ASSUMPTION 1. *If there is a universe that contains all the tasks $E = \{e_1, e_2, \dots, e_\infty\}$ and k_i represents the corresponding knowledge to master each task e_i , then $\forall i, j, k_i \cap k_j \neq \emptyset$.*

This is a formal description of our common sense that any pair of tasks are not absolutely isolated from each other, which has been implicitly used as a fundamental assumption by most prior transfer learning studies [11, 24, 26]. Therefore, we focus on mining the shared knowledge across multiple tasks instead of providing strategy selecting tasks that share knowledge as much as possible, which remains to be unsolved and may lead to our future work. The goal here is to utilize the existing knowledge as well as possible. For example, we may only have a well-trained expert on playing Pong game, and we want to utilize its expertise to help us perform better on other games. This is one of the situations that can be solved by our collaborative deep reinforcement learning framework.

4.1 Collaborative deep reinforcement learning

In deep reinforcement learning, since the training of agents are computational expensive, the well-trained agents should be further utilized as source agents (agents where we transferred knowledge from) to facilitate the training of target agents (agents that are provided with the extra knowledge from source). In order to incorporate this type of collaboration to the training of DRL agents, we formally define the collaborative deep reinforcement learning (CDRL) framework as follows:

Definition 4.1. Given m independent environments $\{\epsilon_1, \epsilon_2, \dots, \epsilon_m\}$ of m tasks $\{e_1, e_2, \dots, e_m\}$, the corresponding m agents $\{g_1, g_2, \dots, g_m\}$ are collaboratively trained in parallel to maximize the rewards (master each task) with respect to target agents.

- **Environments.** There is no restriction on the environments: The m environments can be totally different or with some duplications.

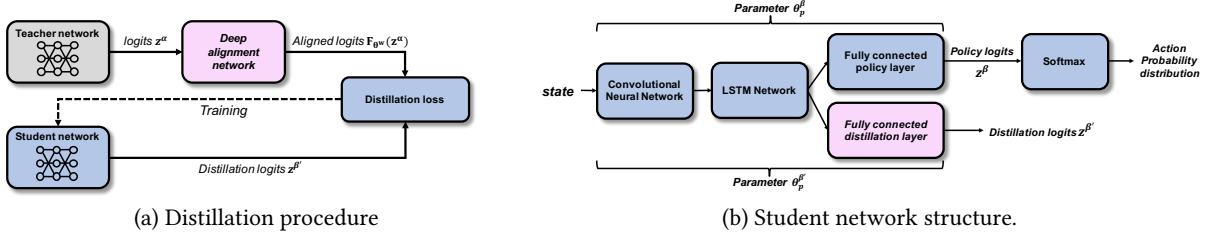


Figure 2: Deep knowledge distillation. In (a), the teacher’s output logits z^α is mapped through a deep alignment network and the aligned logits $\mathcal{F}_{\theta^\omega}(z^\alpha)$ is used as the supervision to train the student. In (b), the extra fully connected layer for distillation is added for learning knowledge from teacher. For simplicity’s sake, time step t is omitted here.

- *In parallel.* Each environment ε_i only interacts with the one corresponding agent g_i , i.e., the action a_t^j from agent g_j at step t has no influence on the state s_{t+1}^i in ε_i , $\forall i \neq j$.
- *Collaboratively.* The training procedure of agent g_i consists of interacting with environment ε_i and interacting with other agents as well. The agent g_i is not necessary to be at same level as “collaborative” defined in cognitive science [9]. E.g., g_1 can be an expert for task e_1 (environment ε_1) while he is helping agent g_2 which is a student agent in task e_2 .
- *Target agents.* The goal of CDRL can be set as maximizing the rewards that agent g_i obtains in environment ε_i with the help of interacting with other agents, similar to inductive transfer learning where g_i is the target agent for target task and others are source tasks. The knowledge is transferred from source to target g_i by interaction. When we set the goal to maximize the rewards of multiple agents jointly, it is similar to multi-task learning where all tasks are source tasks and target tasks at the same time.

Notice that our definition is very different from the previously defined collaborative multiagent Markov Decision Process (collaborative multiagent MDP) [13, 17] where a set of agents select a global joint action to maximize the sum of their individual rewards and the environment is transitted to a new state based on that joint action. First, MDP is not a requirement in CDRL framework. Second, in CDRL, each agent has its own copy of environment and maximizes its own cumulative rewards. The goal of collaboration is to improve the performance of collaborative agents, compared with isolated ones, which is different from maximizing the sum of global rewards in collaborative multiagent MDP. Third, CDRL focuses on how agents collaborate among heterogeneous environments, instead of how joint action affects the rewards. In CDRL, different agents are acting in parallel, the actions taken by other agents won’t directly influence current agent’s rewards. While in collaborative multiagent MDP, the agents must coordinate their action choices since the rewards will be directly affected by the action choices of other agents.

Furthermore, CDRL includes different types of interaction, which makes this a general framework. For example, the current state-of-the-art is A3C [19] can be categorized as one homogeneous CDRL method with advantage actor-critic interaction. Specifically, multiple agents in A3C are trained in parallel with the same environment. All agents first synchronize parameters from a global network, and then update the global network with their individual gradients. This procedure can be seen as each agent maintains its own model (a

different version of global network) and interacts with other agents by sending and receiving gradients.

In this paper, we propose a novel interaction method named deep knowledge distillation under the CDRL framework. It is worth noting that the interaction in A3C only deals with the homogeneous tasks, i.e. all agents have the same environment and the same model structure so that their gradients can be accumulated and interacted. By deep knowledge distillation, the interaction can be conducted among heterogeneous tasks.

4.2 Deep knowledge distillation

As we introduced before, knowledge distillation [15] is trying to train a student network that can behave similarly to the teacher network by utilizing the logits from the teacher as supervision. However, transferring the knowledge among heterogeneous tasks faces several difficulties. First, the action spaces of different tasks may have different dimensions. Second, even if the dimensionality of action space is same among tasks, the action probability distributions for different tasks could vary a lot, as we illustrated in Figure 5 (a) and (b). Thus, the action patterns represented by the logits of different policy networks are usually different from task to task. If we directly force a student network to mimic the action pattern of a teacher network for a different task, it could be trained in a wrong direction, and finally ends up with worse performance than isolated training. In fact, this suspect has been empirically verified in our experiments.

Based on the above observation, we propose deep knowledge distillation to transfer knowledge between heterogeneous tasks. As illustrated in Figure 2 (a), the approach for deep knowledge distillation is straightforward. We use a deep alignment network to map the logits of the teacher network from a heterogeneous source task e^α (environment ε^α), then the logits is used as our supervision to update the student network of target task e^β (environment ε^β). This procedure is performed by minimizing following objective function over student policy network parameters $\theta_p^{\beta'}$:

$$L_{KL}(D, \theta_p^{\beta'}, \tau) = \sum_t l_{KL}(\mathcal{F}_{\theta^\omega}(z_t^\alpha), z_t^{\beta'}, \tau), \quad (3)$$

where

$$l_{KL}(\mathcal{F}_{\theta^\omega}(z_t^\alpha), z_t^{\beta'}, \tau) = \text{softmax}(\mathcal{F}_{\theta^\omega}(z_t^\alpha)/\tau) \ln \frac{\text{softmax}(\mathcal{F}_{\theta^\omega}(z_t^\alpha)/\tau)}{\text{softmax}(z_t^{\beta'})}.$$

Here θ^ω denotes the parameters of the deep alignment network, which transfers the logits z_t^α from the teacher policy network for

knowledge distillation by function $\mathcal{F}_{\theta^\omega}(\mathbf{z}_t^\alpha)$ at step t . As we show in Figure 2 (b), θ_p^β is the student policy network parameters (including parameters of CNN, LSTM and policy layer) for task e^β , while $\theta_p^{\beta'}$ denotes student network parameters of CNN, LSTM and distillation layer. It is clear that the distillation logits $\mathbf{z}_t^{\beta'}$ from the student network does not determine the action probability distribution directly, which is established by the policy logits \mathbf{z}_t^β , as illustrated in Figure 2 (b). We add another fully connected distillation layer to deal with the mismatch of action space dimensionality and the contradiction of the transferred knowledge from source domain and the learned knowledge from target domain. The input to both of the teacher and the student network is the state of environment e^β of target task e^β . It means that we want to transfer the expertise from an expert of task e^α towards the current state. Symbol D is a set of logits from the teacher network in one batch and τ is the temperature same as described in Eq (1). In a trivial case that the teacher network and the student network are trained for same task (e^α equals e^β), then the deep alignment network $\mathcal{F}_{\theta^\omega}$ would reduce to an identity mapping, and the problem is also reduced to a single task policy distillation, which has been proved to be effective in [26]. Before we can apply the deep knowledge distillation, we need to first train a good deep alignment network. In this work, we provide two types of training protocols for different situations:

Offline training: This protocol first trains two teacher networks in both environment e^α and e^β . Then we use the logits of both two teacher networks to train a deep alignment network $\mathcal{F}_{\theta^\omega}$. After acquiring a pre-trained $\mathcal{F}_{\theta^\omega}$, we train a student network of task e^β from scratch, in the meanwhile the teacher network of task e^α and $\mathcal{F}_{\theta^\omega}$ are used for deep knowledge distillation.

Online training: Suppose we only have a teacher network of task e^α , and we want to use the knowledge from task e^α to train the student network for task e^β to get higher performance from scratch. The pipeline of this method is that, we firstly train the student network by interacting with the environment e^β for a certain amount of steps T_1 , and then start to train the alignment network $\mathcal{F}_{\theta^\omega}$, using the logits from the teacher network and the student network. Afterwards, at step T_2 , we start performing deep knowledge distillation. Obviously T_2 is larger than T_1 , and the value of them are task-specific, which is decided empirically in this work.

The offline training could be useful if we have already had a reasonably good model for task e^β , while we want to further improve the performance using the knowledge from task e^α . The online training method is used when we need to learn the student network from scratch. Both types of training protocol can be extended to multiple heterogeneous tasks.

4.3 Collaborative Asynchronous Advantage Actor-Critic

In this section, we introduce the proposed collaborative asynchronous advantage actor-critic (cA3C) algorithm. As we described in section 4.1, the agents are running in parallel. Each agent goes through the same training procedure as described in Algorithm 1. As it shows, the training of agent g_1 can be separated into two parts: The first part is to interact with the environment, get the reward and compute the gradients to minimize the value loss and policy

loss based on Generalized Advantage Estimation (GAE) [27]. The second part is to interact with source agent g_2 so that the logits distilled from agent g_2 can be transferred by the deep alignment network and used as supervision to bias the training of agent g_1 .

To be more concrete, the pseudo code in Algorithm 1 is an involved version of A3C based on online training of deep knowledge distillation. At T -th iteration, the agent interacts with the environment for t_{max} steps or until the terminal state is reached (Line 6 to Line 15). Then the updating of value network and policy network is conducted by GAE. This variation of A3C is firstly implemented in OpenAI universe starter agent [22]. Since the main asynchronous framework is the same as A3C, we still use the A3C to denote this algorithm although the updating is the not the same as advantage actor-critic algorithm used in original A3C paper [19].

The online training of deep knowledge distillation is mainly completed from Line 25 to Line 32 in Algorithm 1. The training of the deep alignment network starts from T_1 steps (Line 25 - 28). After T_1 steps, the student network is able to generate a representative action probability distribution, and we have suitable supervision to train the deep alignment network as well, parameterized by θ^ω . After T_2 steps, θ^ω will gradually converge to a local optimal, and we start the deep knowledge distillation. As illustrated in Figure 2 (b), we use symbol $\theta_p^{\beta'}$ to represent the parameters of CNN, LSTM and the fully connected distillation layer, since we don't want the logits from heterogeneous directly affect the action pattern of target task. To simplify the discussion, the above algorithm is described based on interacting with a single agent from a heterogeneous task. In algorithm 1, logits \mathbf{z}_t^α can be acquired from multiple teacher networks of different tasks, each task will train its own deep alignment network θ^ω and distill the aligned logits to the student network.

As we described in previous section 4.1, there are two types of interactions in this algorithm: 1). GAE interaction uses the gradients shared by all homogeneous agents. 2) Distillation interaction is the deep knowledge distillation from teacher network. The GAE interaction is performed only among homogeneous tasks. By synchronizing the parameters from a global student network in Algorithm 1 (line 3), the current agent receives the GAE updates from all the other agents who interacts with the same environment. In line 21 and 22, the current agent sends his gradients to the global student network, which will be synchronized with other homogeneous agents. The distillation interaction is then conducted in line 31, where we have the aligned logits $\mathcal{F}_{\theta^\omega}(\mathbf{z}_t^\alpha)$ and the distillation logits $\mathbf{z}_t^{\beta'}$ to compute the gradients for minimizing the distillation loss. The gradients of distillation are also sent to the global student network. The role of global student network can be regarded as a parameter server that helps sending interactions among the homogeneous agents. From a different angle, each homogeneous agent maintains an instinct version of global student network. Therefore, both two types of interactions affect all homogeneous agents, which means that the distillation interactions from agent g_2 and agent g_1 would affect all homogeneous agents of agent g_1 ¹.

¹Code is publicly available at <https://github.com/illidanlab/cdr1>

Algorithm 1 online cA3C

Require: Global shared parameter vectors Θ_p and Θ_v and global shared counter $T = 0$; Agent-specific parameter vectors Θ'_p and Θ'_v , GAE [27] parameters γ and λ . Time step to start training deep alignment network and deep knowledge distillation T_1, T_2 .

```

1: while  $T < T_{max}$  do
2:   Reset gradients:  $d\theta_p = 0$  and  $d\theta_v = 0$ 
3:   Synchronize agent-specific parameters  $\theta'_p = \theta_p$  and  $\theta'_v = \theta_v$ 
4:    $t_{start} = t$ , Get state  $s_t$ 
5:   Receive reward  $r_t$  and new state  $s_{t+1}$ 
6:   repeat
7:     Perform  $a_t$  according to policy
8:     Receive reward  $r_t$  and new state  $s_{t+1}$ 
9:     Compute value of state  $v_t = V(s_t; \theta'_v)$ 
10:    if  $T \geq T_1$  then
11:      Compute the logits  $z_t^\alpha$  from teacher network.
12:      Compute the policy logits  $z_t^\beta$  and distillation logits  $z_t^{\beta'}$  from student network.
13:    end if
14:     $t = t + 1, T = T + 1$ 
15:  until terminal  $s_t$  or  $t - t_{start} \geq t_{max}$ 
16:
17:     $R = v_t = \begin{cases} 0 & \text{for terminal } s_t \\ V(s_t, \theta'_v) & \text{for non-terminal } s_t \end{cases}$ 
18:  for  $i \in \{t-1, \dots, t_{start}\}$  do
19:     $\delta_i = r_i + \gamma v_{i+1} - v_i$ 
20:     $A = \delta_i + (\gamma \lambda) A$ 
21:     $R = r_i + \gamma R$ 
22:     $d\theta_p \leftarrow d\theta_p + \nabla \log \pi(a_i | s_i; \theta'_p) A$ 
23:     $d\theta_v \leftarrow d\theta_v + \partial(R - v_i)^2 / \partial \theta'_v$ 
24:  end for
25:  Perform asynchronous update of  $\theta_p$  using  $d\theta_p$  and of  $\theta_v$  using  $d\theta_v$ .
26:  if  $T \geq T_1$  then
27:    // Training deep alignment network.
28:     $\min_{\theta^\omega} \sum_t l_{KL}(z_t^\beta, z_t^\alpha, \tau)$ ,  $l_{KL}$  is defined in Eq (3).
29:  end if
30:  if  $T \geq T_2$  then
31:    // online deep knowledge distillation.
32:     $\min_{\theta^{\beta'}} \sum_t l_{KL}(\mathcal{F}_{\theta^\omega}(z_t^\alpha), z_t^{\beta'})$ 
33:  end if
34: end while

```

5 EXPERIMENTS

5.1 Training and Evaluation

In this work, training and evaluation are conducted in OpenAI Gym [5], a toolkit that includes a collection of benchmark problems such as classic Atari games using Arcade Learning Environment (ALE) [4], classic control games, etc. Same as the standard RL setting, an agent is stimulated in an environment, taking an action and receiving rewards and observations at each time step. The training of the agent is divided into episodes, and the goal is to maximize the expectation of the total reward per episode or to reach higher performance using as few episodes as possible.

5.2 Certificated Homogeneous transfer

In this subsection, we verify the effectiveness of knowledge distillation as a type of interaction in collaborative deep reinforcement

learning for homogeneous tasks. This is also to verify the effectiveness of the simplest case for deep knowledge distillation. Although the effectiveness of policy distillation in deep reinforcement learning has been verified in [26] based on DQN, there is no prior studies on asynchronous online distillation. Therefore, our first experiment is to demonstrate that the knowledge distilled from a certificated task can be used to train a decent student network for a homogeneous task. Otherwise, the even more challenging task of transferring among heterogeneous sources may not work. We note that in this case, the Assumption 1 is fully satisfied given $k_1 = k_2$, where k_1 and k_2 are the knowledge needed to master task e_1 and e_2 , respectively. In this experiment, we conduct experiments in a gym environment named PONG. It is a classic Atari game that an agent controls a paddle to bounce a ball pass another player agent. The maximum reward that each episode can reach is 21.

First, we train a teacher network that learns from its own environment by asynchronously performing GAE updates. We then train a student network using only online knowledge distillation from the teacher network. For fair comparisons, we use 8 agents for all environments in the experiments. Specifically, both the student and the teacher are training in PONG with 8 agents. The 8 agents of the teacher network are trained using the A3C algorithm (equivalent to CDRL with GAE updates in one task). The 8 agents of student network are trained using normal policy distillation, which uses the logits generated from the teacher network as supervision to train the policy network directly. From the results in Figure 3 (a) we see that the student network can achieve a very competitive performance that is almost same as the state-of-arts, using online knowledge distillation from a homogeneous task. It also suggests that the teacher doesn't necessarily need to be an expert, before it can guide the training of a student in the homogeneous case. Before 2 million steps, the teacher itself is still learning from the environment, while the knowledge distilled from teacher can already be used to train a reasonable student network. Moreover, we see that the hybrid of two types of interactions in CDRL has a positive effect on the training, instead of causing performance deterioration.

In the second experiment, the student network is learning from both the online knowledge distillation and the GAE updates from the environment. We find that the convergence is much faster than the state-of-art, as shown in Figure 3 (b). In this experiment, the knowledge is distilled from the teacher to student in the first one million steps and the distillation is stopped after that. We note that in homogeneous CDRL, knowledge distillation is used directly with policy logits other than distillation logits. The knowledge transfer setting in this experiment is not a practical one because we already have a well-trained model of PONG, but it shows that when knowledge is correctly transferred, the combination of online knowledge distillation and the GAE updates is an effective training procedure.

5.3 Certificated Heterogeneous Transfer

In this subsection, we design experiments to illustrate the effectiveness of CDRL in certificated heterogeneous transfer, with the proposed deep knowledge distillation. Given a certificated task PONG, we want to utilize the existing expertise and apply it to facilitate the training of a new task BOWLING. In the following experiments, we do not tune any model-specific parameters such

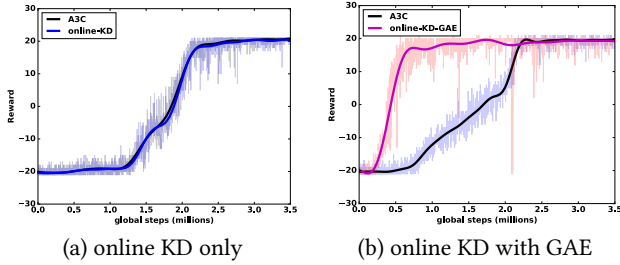


Figure 3: Performance of online homogeneous knowledge distillation. The results show that the combination of knowledge distillation and GAE is an effective training strategy for homogeneous tasks.

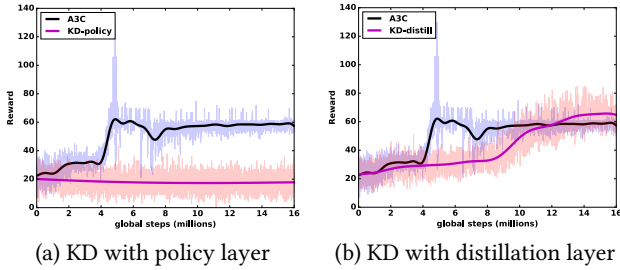


Figure 4: Performance of online knowledge distillation from a heterogeneous task. (a) distillation from a PONG expert using the policy layer to train a BOWLING student (KD-policy). (b) distillation from a PONG expert to a BOWLING student using an extra distillation layer (KD-distill).

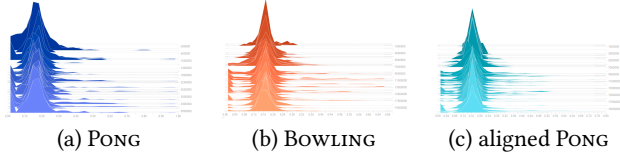


Figure 5: The action probability distributions of a PONG expert, a BOWLING expert and an aligned PONG expert.

as number of layers, size of filter or network structure for BOWLING. We first directly perform transfer learning from PONG to BOWLING by knowledge distillation. Since the two tasks has different action patterns and action probability distributions, directly knowledge distillation with a policy layer is not successful, as shown in Figure 4 (a). In fact, the knowledge distilled from PONG contradicts to the knowledge learned from BOWLING, which leads to the much worse performance than the baseline. We show in Figure 5 (a) and (b) that the action distributions of PONG and BOWLING are very different. To resolve this, we distill the knowledge through an extra distillation layer as illustrated in Figure 2 (b). As such, the knowledge distilled from the certificated heterogeneous task can be successfully transferred to the student network with improved performance after the learning is complete. However, this leads to a much slower convergence than the baseline as shown in Figure 4 (b), because that it takes time to learn a good distillation layer to align the knowledge distilled from PONG to the current learning task. An interesting question is that, is it possible to have both improved performance and faster convergence?

Deep knowledge distillation – Offline training. To handle the heterogeneity between PONG and BOWLING, we first verify the effectiveness of deep knowledge distillation with an offline training procedure. The offline training is split into two stages. In the first stage, we train a deep alignment network with four fully connected layers using the Relu activation function. The training data are logits generated from an expert PONG network and BOWLING network. The rewards of the networks at convergence are 20 and 60 respectively. In stage 2, with the PONG teacher network and trained deep alignment network, we train a BOWLING student network from scratch. The student network is trained with both GAE interactions with its environment, and the distillation interactions from the teacher network and the deep alignment network. The results in Figure 6 (a) show that deep knowledge distillation can transfer knowledge from PONG to BOWLING both efficiently and effectively.

Deep knowledge distillation – Online training. A more practical setting of CDRL is the online training, where we simultaneously train deep alignment network and conduct the online deep knowledge distillation. We use two online training strategies: 1) The training of deep alignment network starts after 4 million steps, when the student BOWLING network can perform reasonably well, and the knowledge distillation starts after 6 million steps. 2) The training of deep alignment network starts after 0.1 million steps, and the knowledge distillation starts after 1 million steps. Results are shown in Figure 6 (b) and (c) respectively. The results show that both strategies reach higher performance than the baseline. Moreover, the results suggest that we do not have to wait until the student network reaches a reasonable performance before we start to train the deep alignment network. This is because the deep alignment network is train to align two distributions of PONG and BOWLING, instead of transferring the actual knowledge. Recall that the action probability distribution of PONG and BOWLING are quite different as shown in Figure 5 (a) and (b). After we projecting the logits of PONG using the deep alignment network, the distribution is very similar to BOWLING, as shown in Figure 5 (c).

5.4 Collaborative Deep Reinforcement Learning

In previous experiments, we assume that there is a well-trained PONG expert, and we transfer knowledge from the PONG expert to the BOWLING student via deep knowledge distillation. A more challenging settings that both of BOWLING and PONG are trained from scratch. In this experiment, we show that the CDRL framework can still be effective in this setting. In this experiment, we train a BOWLING network and a PONG network from scratch using the proposed cA3C algorithm. The PONG agents are trained with GAE interactions only, and the target BOWLING receive supervision from both GAE interactions and distilled knowledge from PONG via a deep alignment network. We start to train the deep alignment network after 3 million steps, and perform deep knowledge distillation after 4 million steps, where the PONG agents are still updating from the environment. We note that in this setting, the teacher network is constantly being updated, as knowledge is distilled from the teacher until 15 million steps. Results in Figure 6 (d) show that the proposed cA3C is able to converge to a higher performance than the current state-of-art. The reward of last one hundred episodes

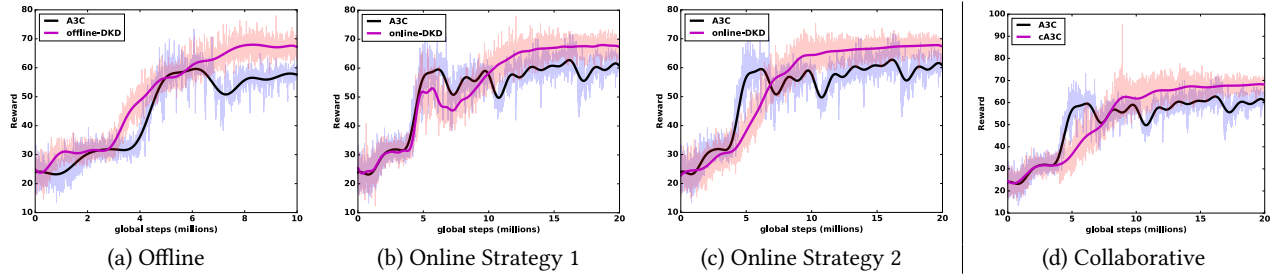


Figure 6: Performance of offline, online deep knowledge distillation, and collaborative learning. Results averaged over 3 runs.

of A3C is 61.48 ± 1.48 , while cA3C achieves 68.35 ± 1.32 , with a significant reward improvement of 11.2%.

6 CONCLUSION

In conclusion, we propose a collaborative deep reinforcement learning framework that can address the knowledge transfer among heterogeneous tasks. Under this framework, we propose deep knowledge distillation to adaptively align the domain of different tasks with the utilization of deep alignment network. Furthermore, we develop an efficient cA3C algorithm and demonstrate its effectiveness by extensive evaluation on OpenAI gym.

ACKNOWLEDGMENTS

This research is supported in part by the Office of Naval Research (ONR) under grant number N00014-14-1-0631 and National Science Foundation under Grant IIS-1565596, IIS-1615597.

REFERENCES

- [1] Naoki Abe, Prem Melville, Cezar Pendus, Chandan K Reddy, David L Jensen, Vince P Thomas, James J Bennett, Gary F Anderson, Brent R Cooley, Melissa Kowalczyk, and others. 2010. Optimizing debt collections using constrained reinforcement learning. In *SIGKDD*. ACM, 75–84.
- [2] Naoki Abe, Naval Verma, Chid Apte, and Robert Schroko. 2004. Cross channel optimized marketing by reinforcement learning. In *SIGKDD*. ACM, 767–772.
- [3] YuXuan Liu Pieter Abbeel Sergey Levine Abhishek Gupta, Coline Devinfi. 2017. Learning Invariant Feature Spaces to Transfer Skills with Reinforcement Learning. In *Under review as a conference paper at ICLR 2017*.
- [4] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. 2013. The Arcade Learning Environment: An evaluation platform for general agents. *J. Artif. Intell. Res. (JAIR)* 47 (2013), 253–279.
- [5] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. 2016. OpenAI gym. *arXiv preprint arXiv:1606.01540* (2016).
- [6] Cristian Bucilu, Rich Caruana, and Alexandru Niculescu-Mizil. 2006. Model compression. In *SIGKDD*. ACM, 535–541.
- [7] Lucian Buşoniu, Robert Babuška, and Bart De Schutter. 2010. Multi-agent reinforcement learning: An overview. In *Innovations in multi-agent systems and applications-1*. Springer, 183–221.
- [8] Marc Carreras, Junku Yuh, Joan Batlle, and Pere Ridao. 2005. A behavior-based scheme using reinforcement learning for autonomous underwater vehicles. *IEEE Journal of Oceanic Engineering* 30, 2 (2005), 416–427.
- [9] Pierre Dillenbourg. 1999. *Collaborative Learning: Cognitive and Computational Approaches*. Advances in Learning and Instruction Series. ERIC.
- [10] A Evgeniou and Massimiliano Pontil. 2007. Multi-task feature learning. *NIPS* 19 (2007), 41.
- [11] Theodoros Evgeniou and Massimiliano Pontil. 2004. Regularized multi-task learning. In *SIGKDD*. ACM, 109–117.
- [12] Anuradha A Gokhale. 1995. Collaborative learning enhances critical thinking. (1995).
- [13] Carlos Guestrin, Michail Lagoudakis, and Ronald Parr. 2002. Coordinated reinforcement learning. In *ICML*, Vol. 2. 227–234.
- [14] Yuval N. Harari. 2015. *Sapiens: a brief history of humankind*.
- [15] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531* (2015).
- [16] Max Jaderberg, Volodymyr Mnih, Wojciech Marian Czarnecki, Tom Schaul, Joel Z Leibo, David Silver, and Koray Kavukcuoglu. 2016. Reinforcement learning with unsupervised auxiliary tasks. *arXiv preprint arXiv:1611.05397* (2016).
- [17] Jelle R Kok and Nikos Vlassis. 2006. Collaborative multiagent reinforcement learning by payoff propagation. *JMLR* 7, Sep (2006), 1789–1828.
- [18] Guillaume Lample and Devendra Singh Chaplot. 2016. Playing FPS games with deep reinforcement learning. *arXiv preprint arXiv:1609.05521* (2016).
- [19] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy P Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. 2016. Asynchronous methods for deep reinforcement learning. *arXiv preprint arXiv:1602.01783* (2016).
- [20] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, and others. 2015. Human-level control through deep reinforcement learning. *Nature* 518, 7540 (2015), 529–533.
- [21] Arun Nair, Praveen Srinivasan, Sam Blackwell, Cagdas Alceick, Rory Fearon, Alessandro De Maria, Vedavyas Panneershelvam, Mustafa Suleyman, Charles Beattie, Stig Petersen, and others. 2015. Massively parallel methods for deep reinforcement learning. *arXiv preprint arXiv:1507.04296* (2015).
- [22] OpenAI. 2017. OpenAI universe-starter-agent. <https://github.com/openai/universe-starter-agent>. (2017). Accessed: 2017-02-01.
- [23] Sinno Jialin Pan and Qiang Yang. 2010. A survey on transfer learning. *TKDE* 22, 10 (2010), 1345–1359.
- [24] Emilio Parisotto, Jimmy Lei Ba, and Ruslan Salakhutdinov. 2015. Actor-mimic: Deep multitask and transfer reinforcement learning. *arXiv preprint arXiv:1511.06342* (2015).
- [25] Janarthanan Rajendran, Aravind Lakshminarayanan, Mitesh M Khapra, Balaram Ravindran, and others. 2015. A2T: Attend, Adapt and Transfer: Attentive Deep Architecture for Adaptive Transfer from multiple sources. *arXiv preprint arXiv:1510.02879* (2015).
- [26] Andrei A Rusu, Sergio Gomez Colmenarejo, Caglar Gulcehre, Guillaume Desjardins, James Kirkpatrick, Razvan Pascanu, Volodymyr Mnih, Koray Kavukcuoglu, and Raia Hadsell. 2015. Policy distillation. *arXiv preprint arXiv:1511.06295* (2015).
- [27] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. 2015. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438* (2015).
- [28] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, and others. 2016. Mastering the game of Go with deep neural networks and tree search. *Nature* 529, 7587 (2016), 484–489.
- [29] Qian Sun, Rita Chattopadhyay, Sethuraman Panchanathan, and Jieping Ye. 2011. A two-stage weighting framework for multi-source domain adaptation. In *NIPS*. 505–513.
- [30] Ming Tan. 1993. Multi-agent reinforcement learning: Independent vs. cooperative agents. In *ICML*. 330–337.
- [31] Matthew E Taylor and Peter Stone. 2009. Transfer learning for reinforcement learning domains: A survey. *JMLR* 10, Jul (2009), 1633–1685.
- [32] Devinder Thapa, In-Sung Jung, and Gi-Nam Wang. 2005. Agent based decision support system using reinforcement learning under emergency circumstances. In *International Conference on Natural Computation*. Springer, 888–892.
- [33] Eric Tzeng, Judy Hoffman, Trevor Darrell, and Kate Saenko. 2015. Simultaneous deep transfer across domains and tasks. In *ICCV*. 4068–4076.
- [34] Yu Zhang and Dit-Yan Yeung. 2012. A convex formulation for learning task relationships in multi-task learning. *arXiv preprint arXiv:1203.3536* (2012).
- [35] Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou Tang. 2014. Facial landmark detection by deep multi-task learning. In *ECCV*. Springer, 94–108.
- [36] Jiayu Zhou, Jianhui Chen, and Jieping Ye. 2011. MALSAR: Multi-task learning via structural regularization. *Arizona State University* (2011).
- [37] Yuke Zhu, Roozbeh Mottaghi, Eric Kolve, Joseph J Lim, Abhinav Gupta, Li Fei-Fei, and Ali Farhadi. 2016. Target-driven visual navigation in indoor scenes using deep reinforcement learning. *arXiv preprint arXiv:1609.05143* (2016).