
Transfer Learning for Deep Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Deep reinforcement learning (RL) has shown the potential to achieve superhuman
2 performance in solving complex decision tasks. Although, unlike humans, it fails to
3 generalise and reuse previously acquired knowledge effectively, which is a crucial
4 ability for a truly intelligent agent. The work proposes an RL-specific modification
5 of CycleGAN, which ensures one-to-one knowledge transfer between different RL
6 tasks. We evaluate the approach on the 2-D Atari game Pong and compare it against
7 two baselines: using GAN and CycleGAN methods. The results demonstrate that
8 our method consistently outperforms the state-of-the-art methods.

9 1 Introduction

10 The inherent ability of reinforcement learning (RL) to dynamically learn complex policies through
11 trial and error has shown great potential in solving diverse decision problems. Deep RL, which
12 combines RL advantages with the power to handle high-dimensional data, recently brought many
13 advances. For instance, model-free methods show significant results in MuJoCo environments,
14 [1], real-world robotic applications, [2] and have demonstrated an ability to achieve super-human
15 performance in Atari games, [3], [4]. Model-based deep RL methods such as AlphaZero, [5], or
16 PlaNET, [6] made significant progress. However, in many real-world tasks, RL remains unsuitable as
17 the errors can be extremely costly. One of the promising ways to address this issue is using *transfer*
18 *learning* (TL), [7] when skills and knowledge collected on similar tasks are applied to the currently
19 solved problem. Besides, learning is essential for developing agents capable of lifelong learning,
20 [8], for simulation-to-real knowledge transfer used in robotics, [9, 10, 11, 12], or for developing the
21 general AI, [13].

22 Despite many advances made, the use of transfer learning in RL and especially TL in deep RL, is
23 limited. For example, 1-pixel perturbations of state observations can lead to useless policies, [14].
24 The RL methods often fail to reuse previously acquired knowledge even in similar tasks when the
25 original image is rotated, or some colours are changed. It has also been shown that learning from
26 scratch can be more efficient than fine-tuning a previously obtained model, [15]. That significantly
27 contrasts with the human ability to generalise and reuse previously acquired knowledge.

28 Main contributions of the paper

- 29 • It introduces a *method for knowledge transfer between two different RL tasks* based on
30 RL-specific modification of CycleGAN. The method is highly applicable in practice. The
31 *method does not rely on paired data* and is independent of the nature of the involved RL
32 tasks. It ensures that the approach can be easily applied to various domains.
- 33 • The work establishes a *correspondence function that reveals similarity of the source and*
34 *target RL tasks*. The proposed formulation suggests learning the correspondence function
35 by minimising the discriminative loss function.

- 36 • The work proposes a *new four-component loss function* which reflects different types of
37 losses. The proposed modification accounts for the actual policy used and takes into account
38 the existing dynamic relationships of the involved RL tasks.
- 39 • We demonstrate how adding two new components *generalises GAN and CycleGAN methods*,
40 i.e. the latter are special cases of the proposed approach.
- 41 • We achieve results indicating a complete reuse of previously acquired knowledge when
42 transferring between the original Pong and rotated Pong.
- 43 • We show that the proposed approach copes with the tasks for which standard approaches
44 (GAN, CycleGAN) fail and learning from scratch remains to be the most efficient method.

45 Experiments with 2-D Atari game Pong, [16], demonstrated that the proposed method notably speeds
46 up the learning and increases the average reward.

47 The paper layout is as follows. Section 2 recalls the necessary background and formulates the TL
48 problem. Section 3.4 constructs the correspondence function and proposes a novel method of its
49 learning. Section 4 describes the experimental evaluation of the proposed approach and compares it
50 with baseline methods. Section 5 provides concluding remarks and outlines future research directions.

51 Related works

52 Survey [17] systematically analyses recent advances in transfer learning for deep RL. The research
53 category to which our approach belongs utilises mapping functions between the source and target
54 task to ensure knowledge transfer. Among them, a line of research that learns common features of
55 RL tasks that can be transferred. It was shown, [18], that the policies learnt on so-called mid-level
56 features can generalise better than those learnt directly on image observations. Work [19] leverages
57 general features of two RL tasks with different dynamics. However, the method is based on paired
58 image observations which are hard or impossible to obtain in practice. Work [20] achieved success in
59 tasks differing in reward function by maintaining successor features and decoupling environment
60 dynamic and reward function. Approach [21] introduces task similarity criterion and builds TL
61 framework based on knowledge shaping, where for similar tasks, efficient transfer is theoretically
62 guaranteed.

63 The pioneering work that used task correspondence was based on unsupervised image-to-image
64 translation models CycleGAN¹, [22], and UNIT², [23]. Approach [15] achieved results on a specific
65 set of tasks by finding correspondence between states of two RL tasks. The application potential
66 of the approach is rather limited as problems like mode-collapse were present. Works [11] and
67 [10] improved the approach proposed in [15] by introducing learnt Q -function or object detection
68 into the learning of the task correspondence. One of the recent approaches, [24], includes an
69 environment model in the learning of the task correspondence. This approach is strongly inspired by
70 the video-to-video translation model, [25].

71 2 Background and notation

72 This section briefly recalls RL formalism and introduces the considered problem.

73 2.1 Notation

74 Throughout the text, sets are denoted by bold capital letters (e.g. \mathbf{X}), \mathbb{N} and \mathbb{R} are sets of natural and
75 real numbers respectively. $\|x\|$ is the L1 norm of x . x_t is the value of x at discrete time $t \in \mathbb{N}$. $E_p[x]$
76 denotes the expected value of x with respect to a probability density p (if provided).

77 We formalise the transfer problem in a general way by considering two RL tasks - the *source task*, S ,
78 and the *target task*, T , characterised by their respective task domains. $\mathbf{S}_S \times \mathbf{A}_S$ and $\mathbf{S}_T \times \mathbf{A}_T$, with \mathbf{S}
79 and \mathbf{A} denoting a set of states and a set of actions respectively.

¹Cycle generative adversarial network

²Unsupervised Image-to-Image Translation Networks

80 2.2 Reinforcement learning

81 Reinforcement learning (RL) considers an *agent* purposefully interacting with an *environment* by
82 selecting actions. RL agent models its environment as a *Markov decision process* (MDP), [26]
83 consisting of discrete sets of observable states \mathbf{S} and actions \mathbf{A} . Set $\mathbf{S} \times \mathbf{A}$ is referred to as the
84 *task domain*. At each time t , the agent observes environment state $s_t \in \mathbf{S}$ and takes action $a_t \in \mathbf{A}$.
85 Executing action a_t at state s_t : i) causes a transition to state s_{t+1} according to *transition function* that
86 describes $p(s_{t+1}|s_t, a_t)$, and ii) provides reward r_t , i. e. the value of reward function $R(s_{t+1}, a_t, s_t)$.
87 The agent’s *goal* is to learn policy $\pi^* : \mathbf{S} \mapsto \mathbf{A}$ that maximises the accumulated reward.

88 Whenever the state space is huge, for instance, when the state is given by a video frame, efficient learn-
89 ing of Q -function calls for numerical approximation. The state-of-the-art in function approximation
90 points to deep neural networks (DNN) as a suitable methodology, [27].

91 Deep Q -networks (DQN), [28], use a standard off-policy Q -learning, [29], and DNN to estimate the
92 Q -function, which then gives the maximizing policy π^* .

93 3 Transfer learning for RL

94 Humans have a remarkable ability to generalise. They do not learn everything from scratch but rather
95 reuse earlier acquired knowledge to a new task or domain³. Generally finding common patterns
96 between different tasks and effectively transferring the concepts learned on one task to another is an
97 essential characteristic of high-level intelligence.

98 In this section, we formalise a problem of transfer learning between two RL tasks, empirically
99 introduce a correspondence function reflecting the similarity of two RL tasks and propose an RL-
100 specific modification of CycleGAN algorithm that realises knowledge transfer between two RL tasks.
101 The proposed transfer i) considers behaviours, which are most useful for the target task; ii) captures
102 and respects common patterns in transition dynamics of the involved RL tasks.

103 3.1 Problem formulation

104 We consider two RL tasks: the *source task*, S , and the *target task*, T with their respective task
105 domains $\mathbf{S}_S \times \mathbf{A}_S$ and $\mathbf{S}_T \times \mathbf{A}_T$. Each of the tasks corresponds to MDP with its own environmental
106 dynamics and reward function, see Section 2.2. Transition functions of the tasks as well as theirs
107 reward functions may be different.

108 This work uses an abstract notion of similarity, inspired by human learning when tackling related
109 problems. Two tasks are similar if they share some common properties, and the knowledge acquired
110 in one task proves to be beneficial in solving the other. This empirical definition can be more formally
111 introduced as follows.

112 **Definition 3.1** (Correspondence function). Consider source S and target T tasks with respective
113 domains $\mathbf{S}_S \times \mathbf{A}_S$ and $\mathbf{S}_T \times \mathbf{A}_T$. A *correspondence function*, $\mathcal{C} : (\mathbf{S}_T \times \mathbf{A}_T) \mapsto (\mathbf{S}_S \times \mathbf{A}_S)$, is a
114 mapping, which reveals the similarity of the involved RL tasks in terms of the dynamics of the tasks’
115 environments and the associated Q -functions.

116 It is clear that function \mathcal{C} establishes the relationship between similar patterns in behaviour of the
117 target and source tasks that are necessary for knowledge transfer. So, if Q_S is the optimal Q -function
118 for the source task, then Q -function

$$Q_S(\mathcal{C}(\cdot, \cdot)) : \mathbf{S}_T \times \mathbf{A}_T \mapsto \mathbb{R} \quad (1)$$

119 gives better performance⁴ on the target task than a random policy.

120 Let us assume (for brevity) that the action spaces of the source and the target RL task are identical,
121 i.e. $\mathbf{A}_S = \mathbf{A}_T$. Let mutually corresponding actions be found using identity mapping regardless
122 of the current state⁵. Thus, we need to learn a mapping indicating corresponding states, i. e. the

³Developmental psychologists have shown that as early as 18 months old, children can infer intentions and imitate the behaviour of adults, [30]. The imitation is complex as children must infer a match between their observations and internal representations, effectively linking the two diverse domains.

⁴Performance is measured by average reward per time.

⁵More specifically, all actions of the source and target task have the same labels and meanings (e.g. $a = 1$ stands for "up"). Therefore, no mapping between source and target task action spaces is necessary

123 correspondence function for states. The searched correspondence function \mathcal{C} is then obtained as
 124 follows:

$$\mathcal{C}(s_T, a_T) = (G_T(s_T), I(a_T)), \quad \forall (s_T, a_T) \in \mathbf{S}_T \times \mathbf{A}_T, \quad (2)$$

125 where G_T is the generator from (3) mapping states from the *target task* to states from the *source task*
 126 and $I(\cdot)$ is an identity mapping.

127 The correspondence function is unknown to RL agent and the next section describes how to learn it.

128 3.2 Learning of correspondence function

129 The proposed learning is inspired by CycleGAN, see Section 3.3, where the learning minimises a
 130 discriminative *loss function*, which makes the similarity metric small for similar patterns and large
 131 otherwise. Even direct application of CycleGAN to the states brought some success in policy transfer,
 132 see for instance [15]. However, data records in experience memories comprise richer yet unused
 133 information that may be helpful for the transfer of knowledge. We propose to include additional
 134 components into the loss function minimised in CycleGAN learning making the method entirely
 135 relevant to RL.

136 This work proposes adding two new components to the CycleGAN losses, (4), (5):

- 137 • Q -loss \mathcal{L}_Q - a loss that reflects how the Q -function learned from the source task, Q_S , copes
 138 with impreciseness in learned generators G_T and G_S .
- 139 • Model-loss \mathcal{L}_M - a loss that reflects the influence of the environment model of the source
 140 task.

141 Let us briefly summarize CycleGAN and explain the reasons for introducing the new components
 142 and their forms.

143 3.3 CycleGAN

144 Cycle-consistent Generative Adversarial Network (CycleGAN), [22], is based on GAN⁶, [31], and
 145 was originally proposed for image-to-image translation. The idea behind *cycle consistency* is that
 146 data that has been translated to a new domain and then recovered from it, should not change.

147 CycleGAN operates with two mappings G_S and G_T called *generators*⁷

$$G_S : \mathbf{S}_S \rightarrow \mathbf{S}_T \quad \text{and} \quad G_T : \mathbf{S}_T \rightarrow \mathbf{S}_S. \quad (3)$$

148 They are learnt as two GANs, that is, simultaneously with the corresponding discriminators D_S and
 149 D_T . Generators learn to map states from \mathbf{S}_S to \mathbf{S}_T and vice-versa, while discriminators learn to
 150 *distinguish* a real state from a state mapped by a generator.

151 Learning in CycleGAN minimises a two-component loss. The first is *adversarial loss*, \mathcal{L}_{GAN} comes
 152 from GAN and is given by

$$\begin{aligned} \mathcal{L}_{GAN} = & E_{s_S} [\log D_S(s_S)] + E_{s_T} [\log (1 - D_S(G_S(s_T)))] \\ & + E_{s_T} [\log D_T(s_T)] + E_{s_S} [\log (1 - D_T(G_T(s_S)))] \end{aligned} \quad (4)$$

153 The adversarial training encourages mappings G_S and G_T (3) to produce outputs indistinguishable
 154 from the real ones, i. e. respective sets \mathbf{S}_S and \mathbf{S}_T .

155 The second component is *cycle-consistency loss*, \mathcal{L}_{Cyc} , that has the following form:

$$\mathcal{L}_{Cyc} = E_{s_S} [\|G_T(G_S(s_S)) - s_S\|] + E_{s_T} [\|G_S(G_T(s_T)) - s_T\|]. \quad (5)$$

156 Minimisation of cycle-consistency loss \mathcal{L}_{Cyc} ensures that every state $s_S \in \mathbf{S}_S$ must be recoverable
 157 after mapping it back to \mathbf{S}_T , i.e. $G_T(G_S(s_S)) \approx s_S$. The same requirement applies to every state
 158 $s_T \in \mathbf{S}_T$.

⁶Generative adversarial network

⁷that translate data between source and target domains

159 ***Q-loss***

160 The available Q -function, Q_S , should be incorporated in the learning of correspondence function \mathcal{C}
 161 as it determines the optimal policy. We introduce loss \mathcal{L}_Q in the following form:

$$\mathcal{L}_Q = E_{s_S} [\|Q_S(G_T(G_S(s_S))) - Q_S(s_S)\|] \quad (6)$$

162 The loss (6) will make the learned correspondence more suitable for transferring knowledge between
 163 tasks S and T because the learned correspondence function, \mathcal{C} , will retain the parts of the states that
 164 are the most important for choosing the optimal action.

165 ***Model loss***

166 So far, all considered losses (4) - (6) are associated with state values. However, every RL task is a
 167 dynamic one, and the time sequence of states is essential. Consider states of the target and the source
 168 tasks at times t and $t + 1$. If generator G_T ensures mapping s_{Tt} on $s_{S_{t+1}}$ and generator G_S maps
 169 $s_{S_{t+1}}$ back to s_{Tt} , then losses \mathcal{L}_{GAN} , \mathcal{L}_{Cyc} and \mathcal{L}_Q (4) - (6) are minimal. However, it would not
 170 help to solve the target RL task.

171 To ensure that the correspondence function, \mathcal{C} , grasps all essential dynamic relationships of the source
 172 and target task, the overall loss minimised must consider the learnt environment model, F^8 of the
 173 source task:

$$\mathcal{L}_M = E_{s_{Tt}, a_{Tt}, s_{T_{t+1}}} [\|F(G_T(s_{Tt}), a_{Tt}) - G_T(s_{T_{t+1}})\|] \quad (7)$$

174 ***Total loss***

175 The proposed total loss comprises all the components (4), (5), (6) and (7) and, thus, has the following
 176 form:

$$\mathcal{L} = \mathcal{L}_{GAN} + \lambda_{Cyc} \mathcal{L}_{Cyc} + \lambda_Q \mathcal{L}_Q + \lambda_M \mathcal{L}_M, \quad (8)$$

177 where λ_{Cyc} , λ_Q and λ_M are *loss parameters* that define relative influence (weight) of the respective
 178 components.

179 The proposed approach, which minimises 4-component loss (8), generalises GAN, [31], and Cycle-
 180 GAN, [22], methods often used for transfer learning. It is easy to see that GAN and CycleGAN can
 181 be obtained by setting some of parameters λ_Q , λ_M , λ_{Cyc} in (8) to zeros as follows:

- 182 • $\lambda_Q = \lambda_M = \lambda_{Cyc} = 0$ (for GAN),
- 183 • $\lambda_Q = \lambda_M = 0$ (for CycleGAN).

184 **3.4 Transfer learning: Algorithm**

185 The main steps of the proposed algorithm:

186 **Step 1** The agent first solves task S by the DQN algorithm. The obtained knowledge, $\mathbf{K}_S =$
 187 (Q_S, \mathbf{M}_S) , consists of learned Q -function, Q_S , and collected experience memory $\mathbf{M}_S =$
 188 $((s_t, a_t, s_{t+1}, r_t)_{i=1}^{n_M})$.

189 **Step 2** The agent applies a random decision rule to task T , collects experience memory \mathbf{M}_T .

190 **Step 3** The assumed similarity of the tasks S and T guarantees the existence of correspondence
 191 function \mathcal{C} (see Definition 3.1). The agent uses knowledge $\mathbf{K}_S = (Q_S, \mathbf{M}_S)$ and memory
 192 \mathbf{M}_T to learn correspondence function \mathcal{C} .

193 **Step 4** Existence of a correspondence function \mathcal{C} , allows to express Q -function of the target task,
 194 Q_T , via Q -function of the source task, Q_S , and learnt correspondence function \mathcal{C} as follows:

$$Q_T(s_T, a_T) = Q_S(\mathcal{C}(s_T, a_T)), \quad \forall (s_T, a_T) \in (\mathbf{S}_T \times \mathbf{A}_T). \quad (9)$$

195 Then the agent can use Q -function Q_S of the source task to choose the optimal actions in the target
 196 task.

⁸Environment model $F: \mathbf{S} \times \mathbf{A} \mapsto \mathbf{S}$ is a mapping taking current state s_t and action a_t and giving the next state s_{t+1} . It is learned using the experience memory, \mathbf{M}_S



Figure 1: **Standard Pong**, [16]

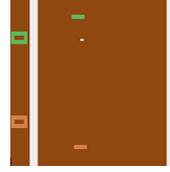


Figure 2: **Pong rotated by 90 degrees**, [16]

197 4 Experimental part

198 To test the efficiency of the proposed approach, two experiments on the Atari game Pong, [16], were
 199 conducted. The performance of the approach was evaluated based on an average accumulated reward
 200 per game. GAN and CycleGAN were used as baseline methods.

201 4.1 Experiment description and setup

202 The proposed TL method was tested in two experiments.

203 **Experiment 1:** The source and target tasks were the same, i.e. game Pong (screenshot is shown in
 204 Fig. 1). The main aim of this experiment was to verify the ability of the proposed approach to find
 205 the identity transformation.

206 **Experiment 2:** The source task was the original Pong while the target task was rotated Pong (see
 207 screenshot in Fig. 2). The game remained the same, but all image frames were rotated by 90 degrees.

208 Each experiment consists of the following steps:

- 209 1) The agent played the *source task* (standard Pong), learned the optimal policy by DQN and
 210 obtained the optimal Q -function Q_S , environment model F and experience memory \mathbf{M}_S
 211 containing 10000 data entries collected at the end of the game.
- 212 2) The agent played the *target task* (standard Pong in Experiment 1 or rotated Pong in Experi-
 213 ment 2) using random policy and obtained data for experience memory \mathbf{M}_T containing 10000
 214 data entries.
- 215 3) The agent started learning the correspondence function \mathcal{C} using the method from Section 3.4,
- 216 4) For every 1000 learning steps, the agent:
 - 217 • suspends learning of correspondence function \mathcal{C} ,
 - 218 • uses learnt \mathcal{C} and the Q -function transformed from the *source task*, see (9), to play five
 219 games of the *target task*, and
 - 220 • computes the average accumulated reward per game.
- 221 5) The agent played the *target task* while using the learned correspondence⁹ and Q -function
 222 Q_S transferred from the *source task*. At the same time the agent uses DQN and fixed \mathcal{C} to
 223 continuously fine-tune Q -function Q_T of the target task.

224 The *key metric* to evaluate the success of the knowledge transfer was the average accumulated reward
 225 per game.

226 **Baselines:** The results are compared with two baselines - using GAN and CycleGAN methods, [31],
 227 [22], which have been recently used for knowledge transfer in similar settings, [15].

228 The following sections provide the key details of the experiments performed and their results.

229 4.2 Experiment 1

230 This experiment aimed to test transfer learning when *source* and *target* tasks are identical.

⁹the correspondence function that achieved the highest average accumulated reward per game in the previous step was used here

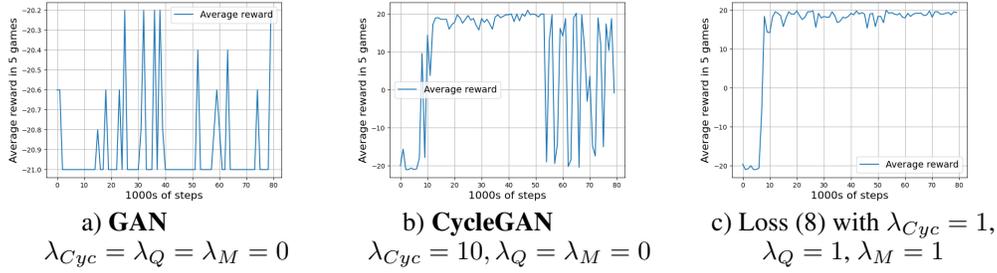


Figure 3: **Experiment 1:** Average accumulated reward per game when playing five games with the transformed Q -function (9) every 1000 learning steps. The performance is shown for different values of loss parameters λ_{Cyc} , λ_Q and λ_M . Fig. 3a and 3b show the **baselines** using *GAN* and *CycleGAN* methods.

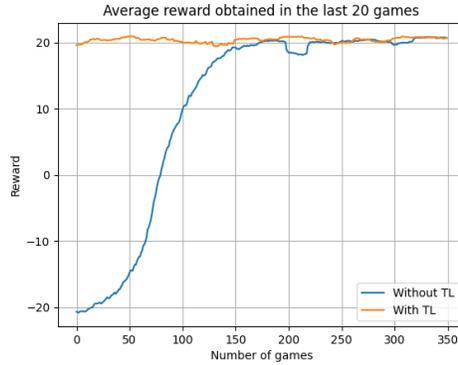


Figure 4: Moving average of reward per game computed from the last 20 games depending on the number of *Pong* games played. The blue line denotes learning from scratch, i. e. without TL, the orange line denotes the case with TL. The Q -function Q_T is continuously learned during the game in both cases.

231 G_S and G_T generators (see Section 3.2) were constructed as neural networks with convolutional
 232 layers. Their specific architecture was taken from [32]. The discriminators D_S and D_T were also
 233 constructed as neural networks with convolutional layers with the architecture as in [33].

234 The transfer learning with the loss (8) was tested for all the combinations of the parameters: $\lambda_{Cyc} \in$
 235 $\{0, 1, 10\}$, $\lambda_Q \in \{0, 1\}$ and $\lambda_M \in \{0, 1, 10\}$.

236 The results are presented in Fig. 3 and Fig. 4. The main findings are:

- 237 • The best results are obtained for the proposed loss function, (8), that contains the proposed
 238 components \mathcal{L}_Q and \mathcal{L}_M , i. e. $\lambda_{Cyc} = \lambda_Q = \lambda_M = 1$, Fig. 3e.
- 239 • *GAN* baseline, Fig. 3a, does not produce good results. Performance of *CycleGAN* baseline,
 240 Fig. 3b, soon became unstable though it provides good rewards at the beginning.
- 241 • The agent successfully learned the correspondence function and completely reused previ-
 242 ously acquired knowledge, Fig. 4.

243 4.3 Experiment 2

244 In Experiment 2, the *target task* is the original *Pong* with image frames rotated by 90 degrees (see
 245 Fig. 2).

246 Generators G_S and G_T , (see (3) and Section 3.2) are constructed as neural networks with two
 247 different architectures (for both of them). The architecture of the first one, referred to here as the

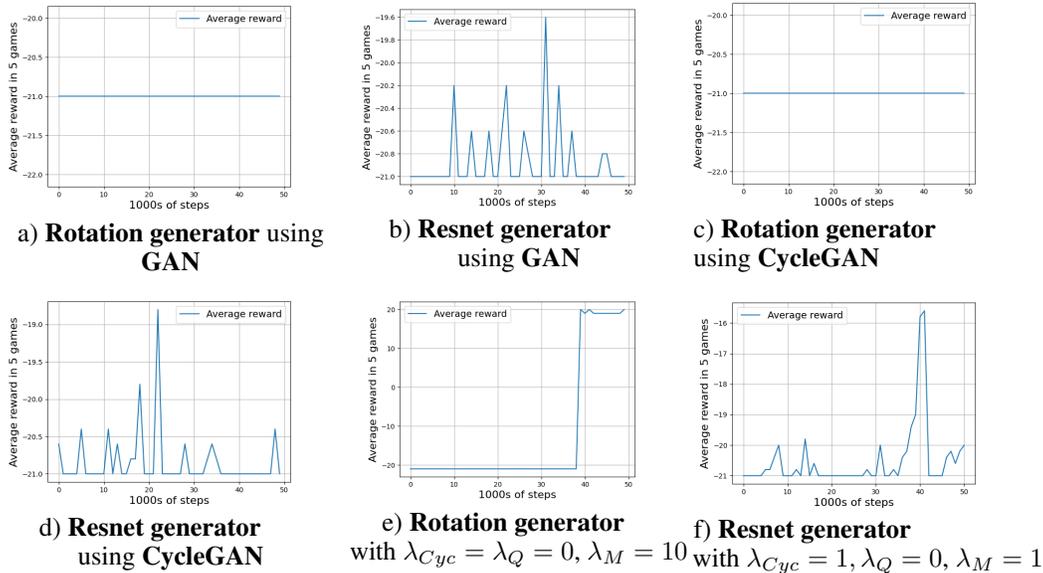


Figure 5: Average accumulated reward in five games when playing Rotated Pong with the transformed Q -function (9) every 1000 learning steps. The results are shown for the **rotation** and the **resnet** generator with the best settings of the loss parameters in each case (e, f) as well as with using **GAN** and **CycleGAN** baselines (a-d).

248 **resnet generator**, was taken from [32] and then followed by a rotation layer, see [34]. The second
 249 type, referred to as the **rotation generator**, was composed of the mentioned rotation layer only.

250 The proposed approach was tested for different values of loss parameters λ_{Cyc} , λ_Q and λ_M , (8).
 251 Fig. 5 - Fig. 7 present the best-reached performance compared with the performance of the baseline
 252 methods.

253 The main findings are the following:

- 254 • The GAN and CycleGAN baselines did not produce a correspondence function suitable for
 255 knowledge transfer.
- 256 • The rotation generator yields the perfect correspondence function. Learning the corre-
 257 spondence function with the **resnet generator** provided significantly better results than learning
 258 the Q -function from scratch.
- 259 • Fine-tuning the Q -function from the source task gives worse performance on the target task
 260 than learning a new Q -function from scratch.

261 5 Conclusion

262 We propose a method for knowledge transfer between two different reinforcement learning tasks.
 263 Our approach establishes the correspondence function that reveals the similarity between the *source*
 264 and *target task*. The neural network approximates the correspondence function and learns it from
 265 unpaired data using dynamic cycle consistency. To ensure that the essential dynamic relationships
 266 between the involved RL tasks are exploited, we have introduced a four-component loss (8) with two
 267 novel components: model loss and Q -loss.

268 We show the efficacy of our approach on simulated experiments involving the 2-D Atari game Pong
 269 and compare it against two baselines using GAN and CycleGAN methods.

270 The results show that the proposed approach outperforms baseline methods. The introduced corre-
 271 spondence function respects Q -function and environment model of the source task and establishes
 272 them into learning. This allows the agent to gradually build, adapt, and use a set of skills while

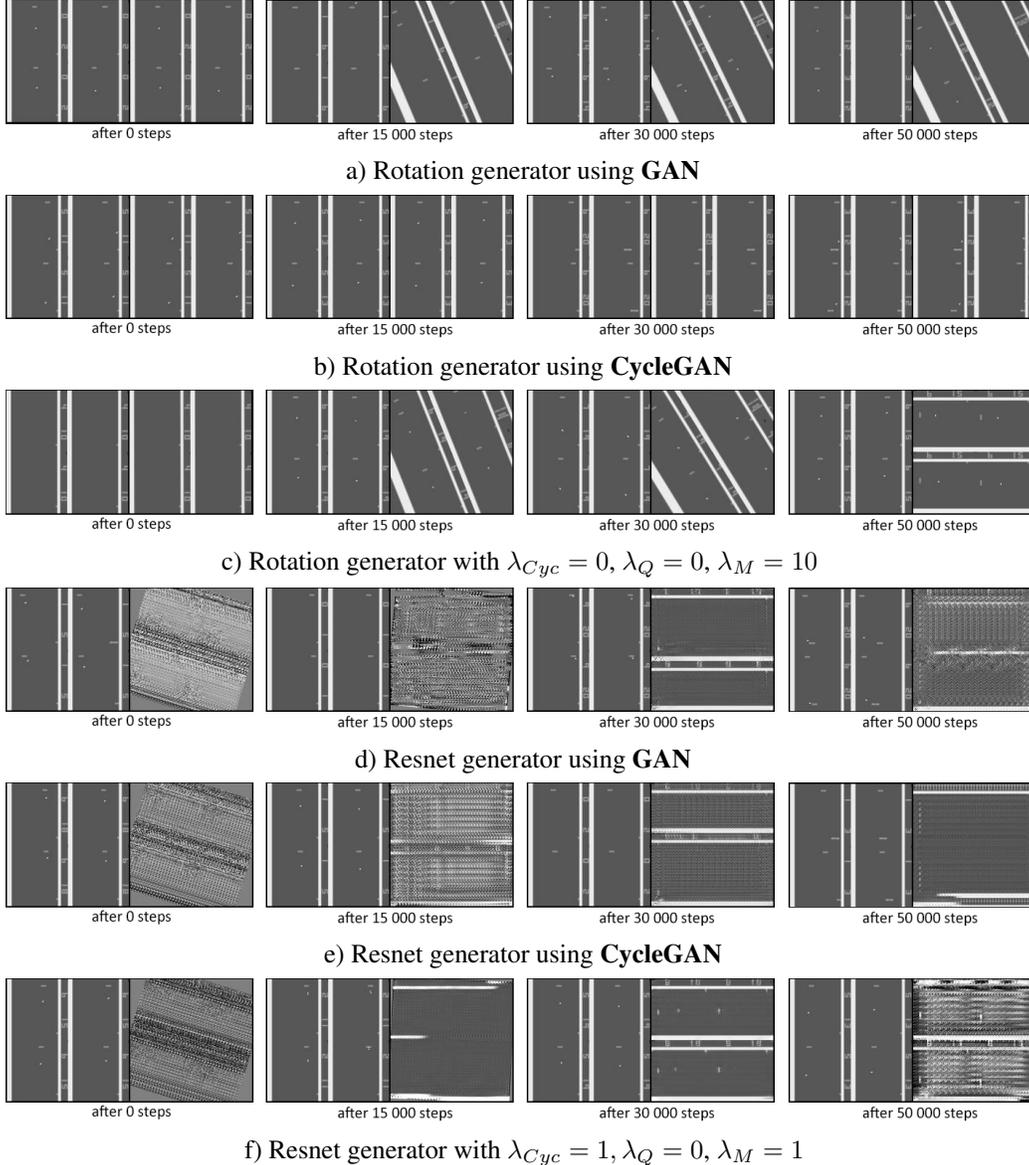


Figure 6: **Experiment 2:** Screenshots of the game depicting progress in learning the correspondence function \mathcal{C} (2) after 0, 15000, 30000 and 50000 steps. The results are shown for the rotation and the resnet generators with the best settings of parameters λ_{Cyc} , λ_Q and λ_M (8) as well as with using **GAN** and **CycleGAN** baselines. The left parts of the pictures are game frames of the *target* task representing the states, and the right parts are the same states transformed by the correspondence function \mathcal{C} .

273 interacting with the dynamically changing environment, which is generally different from the source
 274 task.

275 The most significant advantage of the proposed method is its practical applicability. The solution
 276 does not rely on paired data and is independent of the nature of the involved RL tasks. It ensures that
 277 the approach can be easily applied to various domains.

278 The foreseen research should focus on the open problems: i) how to perform transfer learning between
 279 tasks having low similarity, ii) how to identify and transfer relevant knowledge from several source
 280 tasks, iii) how to choose the only relevant source tasks similarly to [35], iv) what is a better network
 281 architecture for the correspondence function learning.

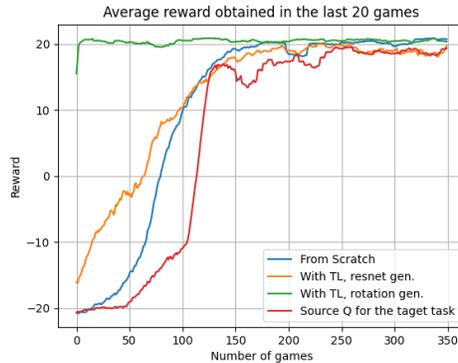


Figure 7: Moving average of reward per game computed from the last 20 games depending on the number of played games for the game rotated Pong for four different agents - an agent learning the game from scratch (blue line), an agent using the correspondence function learned with the resnet generator (orange line), an agent using the correspondence function learned with the rotation generator (green line) and an agent reusing only the Q -function without any correspondence function (red line). The agents were continuously learning the Q -function.

282 **Method implementation:** The method implementation in Python is available at
 283 https://github.com/***/*** (anonymized)

284 References

- 285 [1] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa,
 286 David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv*
 287 *preprint arXiv:1509.02971*, 2015.
- 288 [2] A Rupam Mahmood, Dmytro Korenkevych, Gautham Vasan, William Ma, and James Bergstra.
 289 Benchmarking reinforcement learning algorithms on real-world robots. In *Conference on Robot*
 290 *Learning*, pages 561–591. PMLR, 2018.
- 291 [3] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan
 292 Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint*
 293 *arXiv:1312.5602*, 2013.
- 294 [4] Yueyue Hu, Shiliang Sun, Xin Xu, and Jing Zhao. Attentive multi-view reinforcement learning.
 295 *International Journal of Machine Learning and Cybernetics*, 11(11):2461–2474, 2020.
- 296 [5] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur
 297 Guez, Marc Lanctot, Laurent Sifre, Dharrshan Kumaran, Thore Graepel, et al. A general
 298 reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*,
 299 362(6419):1140–1144, 2018.
- 300 [6] Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee,
 301 and James Davidson. Learning latent dynamics for planning from pixels. In *International*
 302 *Conference on Machine Learning*, pages 2555–2565. PMLR, 2019.
- 303 [7] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on*
 304 *Knowledge and Data Engineering*, 22(10):1345–1359, 2009.
- 305 [8] Haitham Bou Ammar, Eric Eaton, José Marcio Luna, and Paul Ruvolo. Autonomous cross-
 306 domain knowledge transfer in lifelong policy gradient reinforcement learning. In *Twenty-fourth*
 307 *International Joint Conference on Artificial Intelligence*, 2015.
- 308 [9] Stephen James, Paul Wohlhart, Mrinal Kalakrishnan, Dmitry Kalashnikov, Alex Irpan, Julian
 309 Ibarz, Sergey Levine, Raia Hadsell, and Konstantinos Bousmalis. Sim-to-real via sim-to-
 310 sim: Data-efficient robotic grasping via randomized-to-canonical adaptation networks. In

- 311 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
312 12627–12637, 2019.
- 313 [10] Daniel Ho, Kanishka Rao, Zhuo Xu, Eric Jang, Mohi Khansari, and Yunfei Bai. Retinagan:
314 An object-aware approach to sim-to-real transfer. In *2021 IEEE International Conference on*
315 *Robotics and Automation (ICRA)*, pages 10920–10926. IEEE, 2021.
- 316 [11] Kanishka Rao, Chris Harris, Alex Irpan, Sergey Levine, Julian Ibarz, and Mohi Khansari. Rl-
317 cyclegan: Reinforcement learning aware simulation-to-real. In *Proceedings of the IEEE/CVF*
318 *Conference on Computer Vision and Pattern Recognition*, pages 11157–11166, 2020.
- 319 [12] Wei Zhu, Xian Guo, Dai Owaki, Kyo Kutsuzawa, and Mitsuhiro Hayashibe. A survey of
320 sim-to-real transfer techniques applied to reinforcement learning for bioinspired robots. *IEEE*
321 *Transactions on Neural Networks and Learning Systems*, 2021.
- 322 [13] Jeff Clune. AI-GAs: AI-generating algorithms, an alternate paradigm for producing general
323 artificial intelligence. *arXiv preprint arXiv:1905.10985*, 2019.
- 324 [14] Xinghua Qu, Zhu Sun, Yew-Soon Ong, Abhishek Gupta, and Pengfei Wei. Minimalistic attacks:
325 How little it takes to fool deep reinforcement learning policies. *IEEE Transactions on Cognitive*
326 *and Developmental Systems*, 13(4):806–817, 2020.
- 327 [15] Shani Gamrian and Yoav Goldberg. Transfer learning for related reinforcement learning tasks
328 via image-to-image translation. In *International Conference on Machine Learning*, pages
329 2063–2072. PMLR, 2019.
- 330 [16] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning
331 environment: An evaluation platform for general agents. *Journal of Artificial Intelligence*
332 *Research*, 47:253–279, 2013.
- 333 [17] Zhuangdi Zhu, Kaixiang Lin, Anil K. Jain, and Jiayu Zhou. Transfer learning in deep reinforce-
334 ment learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,
335 45(11):13344–13362, 2023.
- 336 [18] Bryan Chen, Alexander Sax, Gene Lewis, Iro Armeni, Silvio Savarese, Amir Zamir, Jitendra
337 Malik, and Lerrel Pinto. Robust policies via mid-level visual representations: An experimental
338 study in manipulation and navigation. *arXiv preprint arXiv:2011.06698*, 2020.
- 339 [19] Abhishek Gupta, Coline Devin, YuXuan Liu, Pieter Abbeel, and Sergey Levine. Learn-
340 ing invariant feature spaces to transfer skills with reinforcement learning. *arXiv preprint*
341 *arXiv:1703.02949*, 2017.
- 342 [20] André Barreto, Will Dabney, Rémi Munos, Jonathan J Hunt, Tom Schaul, Hado P van Hasselt,
343 and David Silver. Successor features for transfer in reinforcement learning. *Advances in Neural*
344 *Information Processing Systems*, 30, 2017.
- 345 [21] Xiang Gao, Jennie Si, and He Huang. Reinforcement learning control with knowledge shaping.
346 *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–12, 2023.
- 347 [22] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image
348 translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE International*
349 *Conference on Computer Vision*, pages 2223–2232, 2017.
- 350 [23] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation
351 networks. *Advances in Neural Information Processing Systems*, 30, 2017.
- 352 [24] Qiang Zhang, Tete Xiao, Alexei A Efros, Lerrel Pinto, and Xiaolong Wang. Learning
353 cross-domain correspondence for control with dynamics cycle-consistency. *arXiv preprint*
354 *arXiv:2012.09811*, 2020.
- 355 [25] Aayush Bansal, Shugao Ma, Deva Ramanan, and Yaser Sheikh. Recycle-gan: Unsupervised
356 video retargeting. In *Proceedings of the European Conference on Computer Vision (ECCV)*,
357 pages 119–135, 2018.

- 358 [26] Martin L Puterman. Markov decision processes. *Handbooks in operations Research and*
359 *Management Science*, 2:331–434, 1990.
- 360 [27] Balázs Csanád Csáji et al. Approximation with artificial neural networks. *Faculty of Sciences,*
361 *Etvös Loránd University, Hungary*, 24(48):7, 2001.
- 362 [28] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G
363 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.
364 Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- 365 [29] Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8(3):279–292, 1992.
- 366 [30] Andrew N Meltzoff. Understanding the intentions of others: re-enactment of intended acts by
367 18-month-old children. *Developmental Psychology*, 31(5):838, 1995.
- 368 [31] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil
369 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in Neural*
370 *Information Processing Systems*, 27, 2014.
- 371 [32] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer
372 and super-resolution. In *European Conference on Computer Vision*, pages 694–711. Springer,
373 2016.
- 374 [33] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with
375 conditional adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision*
376 *and Pattern Recognition*, pages 1125–1134, 2017.
- 377 [34] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks.
378 *Advances in Neural Information Processing Systems*, 28, 2015.
- 379 [35] Marzieh Davoodabadi Farahani and Nasser Mozayani. Evaluating skills in hierarchical reinforce-
380 ment learning. *International Journal of Machine Learning and Cybernetics*, 11(10):2407–2420,
381 2020.