
Attention-based Partial Decoupling of Policy and Value for Generalization in Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this work, we introduce Attention-based Partially Decoupled Actor-Critic (AP-
2 DAC), an actor-critic architecture for generalization in reinforcement learning,
3 which partially separates the policy and the value function. To learn directly from
4 images, traditional actor-critic architectures use a shared network to represent the
5 policy and value function. While a shared representation for policy and value
6 allows parameter and feature sharing, it can also lead to overfitting that catastrophically
7 hurts generalization performance. On the other hand, two separate networks
8 for policy and value can help to avoid overfitting and reduce the generalization
9 gap, but at the cost of added complexity both in terms of architecture design and
10 hyperparameter tuning. APDAC provides an intermediate tradeoff that combines
11 the strengths of both architectures by sharing the initial part of the network and
12 separating the later parts for policy and value. It also incorporates an attention
13 mechanism to propagate relevant features to the separate policy and value blocks.
14 Our empirical analysis shows that APDAC significantly outperforms the PPO base-
15 line and achieves comparable performance with respect to the recent state-of-the-art
16 method IDAAC on the challenging RL generalization benchmark Procgen.

17 1 Introduction

18 Deep reinforcement learning algorithms have shown human-level performance on a variety of different
19 control tasks Mnih et al. [2015] Mnih et al. [2016] [Harojo et al., 2018]. They can master complex
20 tasks by exploring and specializing over a training task and environment given a large number of
21 samples. Deployment of such intelligent systems in real-world applications requires significant
22 generalization and faster adaptation capabilities with respect to similar but unseen scenarios or
23 environments. However, generalizability of this magnitude has yet to be achieved for standard RL
24 algorithms [Cobbe et al., 2021][Cobbe et al., 2020][Grigsby and Qi, 2020]Justesen et al. [2018].

25 Until recently, deep RL algorithms were trained and tested on the same environment. Thus, the
26 issue of overfitting was not consistently observed and measured, but instead implicitly appreciated.
27 Several recent works reveal the potential side effects of such a limited approach to evaluation in the
28 assessment of generalization. These findings motivate the development of benchmarks that provide a
29 better way to quantify an agent’s ability to generalize. With the emergence of such benchmarks, it
30 has become customary to train and test on different sets of similar scenarios to effectively evaluate
31 generalization[Cobbe et al., 2021][Cobbe et al., 2020]Raileanu and Fergus [2021].

32 Generalization is a fundamental aspect of representation learning in episodic tasks consisting of
33 diverse levels. Compared to hand-designed levels, procedural content generation techniques enable
34 generation of a nearly unlimited number of highly varied levels. In this work, we consider the problem
35 of generalization to unseen scenarios or levels of procedurally generated environments given exposure
36 to a limited number of levels during training. The levels vary in terms of background, dynamics,

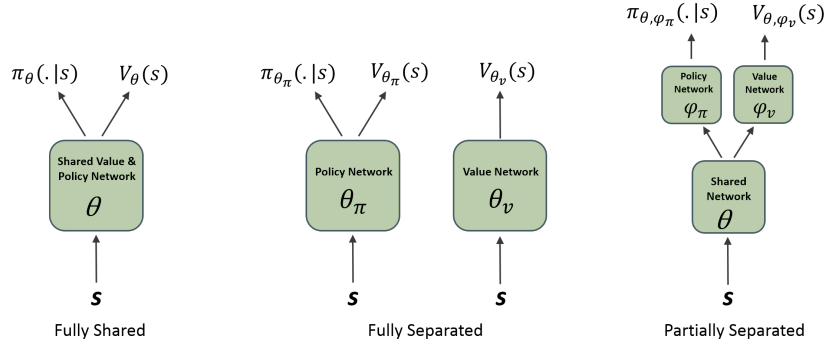


Figure 1: Comparison of architectures: (left) a fully shared network for policy and value (middle), two explicitly separate networks for policy and value (right), and our proposed partially separated network for policy and value function

37 game assets, and the attributes of the entities such as position, spawn time, shape, and color; however,
 38 all the levels share the same end goal. Thus, it significant generalization capability is needed to learn
 39 a robust policy that perform well on levels, episodes, or scenarios that have not yet encountered at
 40 training time.

41 Recently, Raileanu and Fergus [2021] demonstrated a policy-value representation asymmetry, which
 42 suggests that value estimation requires more information compared to that needed to learn an optimal
 43 policy. Thus, shared representation of policy and value function can lead to overfitting. The learned
 44 representation can easily be biased towards the instance specific features responsible for accurate
 45 estimation of the value function. Consequently, the learned policy which generally requires only the
 46 minimal set of task relevant features loses its ability to generalize to unseen variations of the same
 47 task.

48 An alternative to the shared representation is to separate the policy and the value networks[Cobbe
 49 et al., 2021]Raileanu and Fergus [2021]. This helps to disentangle the features necessary to properly
 50 estimate the value and policy function. However, the policy function cannot be learned in isolation,
 51 via standalone training: it requires gradients from the value function to learn the optimal policy. Thus,
 52 additional measures need to be taken to improve the policy network, which increases the complexity
 53 of overall training. Moreover, two networks entail increased memory requirements and training time.

54 As shown in Figure 1, we trade-off between these two extreme approaches - fully shared vs. fully
 55 separate networks - by combining the benefits of both while mitigating their disadvantages. We
 56 propose an *Attention-based Partially Decoupled Actor-Critic (APDAC)* that shares some early layer
 57 blocks of the network while separating the later (downstream) ones into policy and value subnetworks
 58 fed by the shared blocks. Additionally, we deploy an attention mechanism in the two separate
 59 branches, which further decouples policy and value function approximation. We attribute the benefits
 60 of our approach to the hierarchical representation of feature and the ability of attention mechanisms
 61 to effectively identify the components of an input pertinent to the optimization task. We also conduct
 62 an ablation study to better understand the significance of each component of our contribution.

63 In summary, the key contributions of this work are as follows: (i) we propose a new approach that
 64 partially shares and partially decouples the value and policy network; (ii) we develop an integrated
 65 attention mechanism to encourage distinct feature learning for policy and value with minimum
 66 overhead; (iii) we demonstrate competitive performance compared to the state-of-the-art methods on
 67 the Procgen benchmark.

68 2 Related Work

69 Many recent works have established that lack of generalization is a systemic issue in the domain deep
 70 reinforcement learning and popular algorithms tend to overfit to the environment, resulting in models
 71 which seem merely to memorize surface-level details of the environment rather than generalizable
 72 skills [Rajeswaran et al., 2017, Justesen et al., 2018, Grigsby and Qi, 2020, Raileanu and Rocktäschel,

73 2020]. Existing solutions to the generalization problem include L2 regularization, dropout, data
74 augmentation, selective noise augmentation, and batch normalization [Igl et al., 2019a, Cobbe et al.,
75 2019, Igl et al., 2019b, Hu et al., 2021]. As established in Cobbe et al. [2019], Procgen is a testing
76 suite which uses procedural content generation to benchmark generalization to greater effect than
77 traditional benchmarks, and became popular with recent works forwarding generalization in deep
78 reinforcement learning [Igl et al., 2020, Wang et al., 2020, Raileanu and Fergus, 2021, Mazouze
79 et al., 2021]. As discussed in Raileanu and Fergus [2021], sharing features between policy and
80 value functions can lead to overfitting, harming the model’s ability to generalize to new, unseen
81 environments. In contrast to the previous methods, Raileanu and Fergus [2021] Cobbe et al. [2021]
82 make use of fully disconnected policy and value functions. This provides greater generalization and
83 sample efficiency than earlier counterparts, as indicated by state-of-the-art performance on nearly
84 all Procgen environments. However, this performance comes at the cost of a greater number of
85 parameters than previous approaches requiring more computing power. In addition, certain Procgen
86 environments requires specific hyperparameters to produce reported performance. Our method
87 provides results consistent with those in [Raileanu and Fergus, 2021] with fewer parameters and
88 reducing the need for hyperparameter tuning.

89 Literatures exploring the potential of attention mechanisms in neural networks have found success
90 across a wide array of domains, including natural language processing and vision, both as part of
91 convolutional layers and as stand-alone layers [Iqbal and Sha, 2019, Hu, 2019, Ramachandran et al.,
92 2019]. Attention has been proven as an useful paradigm in the domain of natural language processing,
93 seeing wide usage in NLP tasks such as sentiment classification, relation classification, and text
94 summarization [Qin et al., 2017, Lei et al., 2018, Hu, 2019]. Attention has also been utilized in vision
95 models to great success, yielding strong performance while requiring less computing power and fewer
96 input parameters, with self-attention and dual attention models being used in pursuits such as image
97 classification and scene segmentation [Fu et al., 2019, Bello et al., 2019, Ramachandran et al., 2019].
98 The use of attention mechanisms in the domain of deep reinforcement learning, however, is less
99 prevalent. The closest similar works involve variations of A2C with a shared attention mechanism
100 [Iqbal and Sha, 2019, Barati and Chen, 2019]. However, our work differs by combining the attention
101 mechanism with a partially split policy and value function model which is designed to prevent
102 overfitting and achieve generalization.

103 **3 Attention-based Partially Decoupled Actor-Critic**

104 In Attention-based Partially Decoupled Actor-Critic (APDAC), we modify the traditional shared
105 representation of the actor-critic model by partially separating the policy and the value function
106 followed by a shared component. Each of the partially separated policy and value sub-networks are
107 enhanced with the inclusion of multiple attention modules.

108 **3.1 Partial Decoupling of Policy and Value function**

109 Decoupling of the policy and the value function is crucial to overcome the problem of overfitting,
110 which is the main drawback of a shared representation [Raileanu and Fergus, 2021]. However, simply
111 employing two explicitly separate network has serious inherent drawbacks due to the dependency of
112 the policy function approximation on the gradient of the value function. Cobbe et al. [2021] shows
113 that, such straightforward method of using two separate networks for policy and value functions
114 brings a performance decrease when compared to the shared network architecture. To address this
115 issue, works which utilize the separate network model to optimize policy and value functions often
116 use an auxiliary value or advantage head in the policy network. This auxiliary head provides a helpful
117 gradient to the separate policy network to learn better task-relevant policy representation, whereas
118 the separate value network optimizes the value function which plays the original role of the critic.
119 Moreover, the two network models bring in additional number of hyperparameters such as the update
120 frequency of the policy network, update frequency of the value network, coefficient for the advantage
121 loss.

122 We leverage the hierarchical representation of image features to design our network. Generally,
123 the low-level features include minor details such as lines, edges, dots, and curves, whereas the
124 high-level features are composed of multiple low-level features. Based on this, we hypothesize
125 that, although the high-level features responsible for accurate estimation of the policy and value

126 function may differ, the low-level features which constitute the high-level features are almost similar
 127 for both. The performance of generalization in RL increases with the number of convolutional
 128 layers. Thus, it has become customary to use very deep networks specially while learning directly
 129 from image observation. Deep convolutional neural networks (CNN) are capable of learning the
 130 feature representations hierarchically using a sequence of convolutional and pooling layers. Initial
 131 convolutional layers in a neural network learn the filters to capture low-level features while the
 132 later layers in the pipeline learn to identify larger objects and shapes. So, we propose to bifurcate
 133 the network at the convolutional layer level instead of merely separating the policy and the value
 134 head as in the case of a fully shared network. This helps to decouple the high-level feature learning
 135 for policy and value on top of the same feature learned by the shared network. Thus, our network
 136 comprises three parts, the part of the network shared between policy and value parameterized by
 137 θ , the part dedicated to policy learning parameterized by ϕ_π , and another part dedicated to value
 138 function approximation which is parameterized by ϕ_v . The overall network is trained all together to
 139 optimize the following objective:

$$J_{APDAC}(\theta, \phi_\pi, \phi_v) = J_\pi(\theta, \phi_\pi) - \alpha_v L_V(\theta, \phi_v) + \alpha_s S_\pi(\theta, \phi_\pi) \quad (1)$$

140 where $J_\pi(\theta, \phi_\pi)$ is the policy gradient objective, $L_V(\theta, \phi_v)$ is the value loss, $S_\pi(\theta, \phi_\pi)$ is an entropy
 141 bonus that enables efficient exploration, and α_v and α_s are the coefficients denoting relative weight
 142 of the corresponding terms. We optimize the same clipped surrogate policy objective as used in
 143 PPO[Schulman et al., 2017]:

$$J_\pi(\theta, \phi_\pi) = \hat{\mathbb{E}}_t \left[\min(r_t(\theta, \phi_\pi) \hat{A}_t, \text{clip}(r_t(\theta, \phi_\pi), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) \right] \quad (2)$$

144 where $r_t(\theta, \phi_\pi) = \frac{\pi(\theta, \phi_\pi)(a_t|s_t)}{\pi(\theta, \phi_\pi)_{old}(a_t|s_t)}$, and \hat{A}_t is the estimation of the advantage function at timestep
 145 t . The only difference is that the parameters ϕ_π are not affected by the value loss L_V while the
 146 parameters θ are. The value loss L_V is a squared error loss and defined as follows:

$$L_V(\theta, \phi_v) = \hat{\mathbb{E}}_t \left[V_{\theta, \phi_v}(s_t) - \hat{V}_t^{targ} \right] \quad (3)$$

147 where \hat{V}_t^{targ} is the value function target.

148 We argue that the additional overhead introduced by the reliance of the policy optimization on value
 149 gradient, increased number of hyperparameters, and higher memory footprint in case of separate
 150 network architecture can be overcome by a single network architecture which partially separates the
 151 policy and the value. At the same time, our experimental results show that this partial separation
 152 prevents the model from being trapped into the common pitfalls of the fully shared network.

153 3.2 Relevant Feature Learning using Attention

154 To ensure maximum separation between the features learned by the partially separated policy and value
 155 sub-networks, we propose to incorporate individual attention mechanisms within the corresponding
 156 blocks of the network. Attention has been shown as an effective means to learn high quality and
 157 meaningful representation. A good number of attention mechanisms have evolved depending on the
 158 type of feature domain they focus on. We propose to use attention modules similar to Squeeze and
 159 Excitation (SE) network that explicitly models the inter-dependencies between the channels of its
 160 convolutional features. The Squeeze and Excitation block leverages global information to put relative
 161 importance to the useful features compared to the less useful ones. The SE block in the later layers of
 162 a deep network enables distinct feature learning in a highly class-specific manner while in the initial
 163 layers it learns in a class-agnostic manner. This characteristic of the SE block has made it a suitable
 164 choice for our task, where we need to learn distinct features relevant to policy and value. Thus, we
 165 propose to deploy the SE block only in the split value and policy section of the network.

166 Our architecture incorporates the SE block in almost the same fashion as SENet for the Residual
 167 Blocks; however, we add extra SE attention block for the convolutional layers outside of the Residual
 168 block (See Section ?? for details). To utilize global information beyond the local receptive field of

169 filters, the squeeze operation in the attention block first encodes a channel descriptor $z \in R^C$ by
 170 global average pooling. Each element of z is defined as:

$$z_c = F_{sq}(u_c) = \frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W u_c(i, j), \quad (4)$$

171 where $u_c \in \mathbb{R}^{H \times W}$ and $U = [u_1, u_2, \dots, u_c]$ denotes the convolved feature output produced by
 172 the previous convolutional layer. In the next phase, the excitation operation attempts to capture
 173 channel-wise nonlinear dependencies based on the channel-descriptor z . The excitation operation is
 174 realized by two fully-connected (FC) layers around the non-linearity:

$$s = F_{ex}(z, W) = \sigma(g(z, W)) = \sigma(W_2 \delta(W_1 z)), \quad (5)$$

175 where σ refers to the sigmoid activation function, δ refers to the ReLU function, $W_1 \in \mathbb{R}^{\frac{C}{r} \times C}$, and
 176 $W_2 \in \mathbb{R}^{C \times \frac{C}{r}}$; r denotes the reduction ratio parameter between the two FC layers. Finally, the input
 177 feature map \mathbf{U} is rescaled as follows using the learned activation s :

$$\bar{x}_c = F_{scale}(u_c, s_c) = s_c u_c, \quad (6)$$

178 where $\bar{X} = [\bar{x}_1, \bar{x}_2, \dots, \bar{x}_c]$. This way, a set of channel weight is learned to recalibrate the channel
 179 response. Thus, the attention block on the policy and value branch enables attended feature learning
 180 specific to the policy and value functions respectively. From our experiment, it is also evident that the
 181 attention mechanism coupled with the contribution of Section 3.1 helps to attain the same level of
 182 decoupling as the method with two separate network.

183 4 Experiments and Results

184 We evaluate our proposed architecture on the complete Procgen benchmark presented in Cobbe et al.
 185 [2020], which consists of 16 distinct environments with procedurally generated levels. The availability
 186 of highly diverse procedurally generated levels across a wide variety of game environments has
 187 motivated us to choose Procgen as our testbed. Procgen provides a difficulty setting to tune the
 188 difficulty of an environment’s level generation, which can be set to either "easy" or "hard". We
 189 experiment with the easy mode of difficulty for 25 million total timesteps as per the recommendation
 190 of Cobbe et al. [2020]. We train the model on 200 levels and test on the full distribution of the levels.

191 4.1 Network Architecture

192 Following previous works involving Procgen, we chose IMPALA’s deeper residual CNN architecture
 193 as our backbone citepcobbe2020procgen[Cobbe et al., 2021][Raileanu and Fergus, 2021]. Although
 194 this is a relatively large model, it strikes a good balance between the performance on the highly
 195 diverse environment and the required computational power [Cobbe et al., 2020]. This particular
 196 IMPALA CNN architecture has 15 convolutional layers divided into three groups [Espeholt et al.,
 197 2018]. Each group has a similar configuration like Conv - Pooling - Residual Block - Residual Block
 198 where in turns each residual block includes two Conv layer. APDAC shares the first two blocks (10
 199 convolutional layers) of the IMPALA CNN, then branches out for the third block. Thus, APDAC
 200 employs five separate convolutional layers each for policy and value function in order to learn features
 201 distinctly. Finally, it incorporates one attention unit per residual block along with one at the very
 202 beginning of the separation and another just following the first convolutional block.

203 We implement APDAC on top of the implementation of IDAAC Raileanu and Fergus [2021]. We
 204 also used the same PPO code used in Raileanu and Fergus [2021] which is actually built using the
 205 PyTorch implementation of Kostrikov [2018]. When training PPO and IDAAC, whenever applicable,
 206 we followed the same hyperparameter setup from Raileanu and Fergus [2021]. The only difference
 207 is that we reduced the number of mini batch size to minimize computational cost. IDAAC trials
 208 were run using the best hyperparameters for each individual Procgen environment as established in
 209 Raileanu and Fergus [2021].

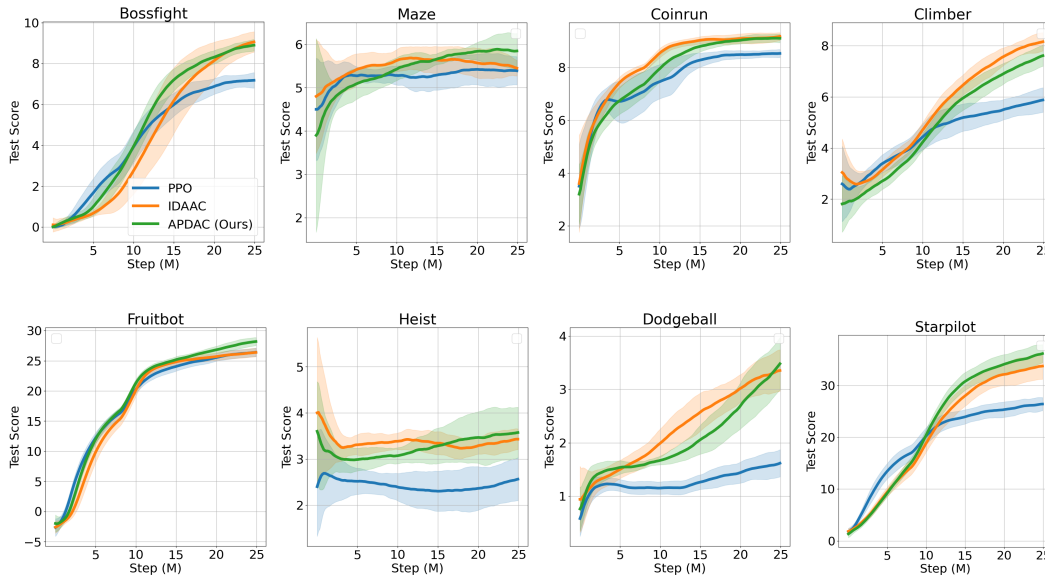


Figure 2: Test performance of PPO, IDAAC, and APDAC over 8 Procgen environments. Mean and standard deviation are calculated over 10 trials, each with a different seed.

210 4.2 Comparative Analysis

211 In our experiments, we compare the performance our approach, APDAC, with the representatives
 212 from two other network topologies. PPO serves as a representative of the models with fully shared
 213 policy and value networks, whereas IDAAC represents those having separate policy and value net-
 214 work [Schulman et al., 2017][Raileanu and Fergus, 2021]. In addition to decoupling the optimization
 215 of the policy and value function using two fully separate networks, IDAAC also uses a discriminator
 216 loss function to decrease the dependency on the instance specific aspects of the environment which
 217 are irrelevant to learning good policy. Figure 2 shows the rolling mean test score averaged over ten
 218 trials for each of the eight environments from the Procgen benchmark. The rolling standard deviation
 219 between these trials is calculated as well, with confidence intervals bounding one standard deviation
 220 above and below each curve. In these results, APDAC sees significant gains in the performance
 221 compared to the shared network approach shown as PPO. Furthermore, it performs similarly, and in
 222 some cases even better, than the existing state-of-the-art, IDAAC, and does so with fewer required
 223 parameters. As such, we conclude that APDAC succeeds as an efficient and performative compromise
 224 between the two existing topologies.

225 4.3 Ablation

226 To evaluate the contribution of each proposed component we further experiment with an ablated
 227 version of our proposed approach, which eliminates all attention blocks. Thus, this network includes
 228 only the partially separated policy and value representation and do not incorporate the contribution
 229 mentioned in Section 3.2. We denote this model as Partially Decoupled Actor-Critic (PDAC). To
 230 determine the difference in performance brought by this ablation, we compare the results achieved
 231 by PPO, APDAC, and the ablation, PDAC, across the entire Procgen benchmark, via the same
 232 experimental setup as before. Figure 3 shows the ablation results from three example environments,
 233 in which we see APDAC performing better, to an extent, than PDAC. Indeed, it is clear from the
 234 comparison with PPO that APDAC’s main performance gain comes from the partial separation of
 235 policy and value network. As such, sharing the initial part of the network doesn’t harm performance,
 236 and in fact reduces the number of parameters as compared to the network architecture with two
 237 separate networks.

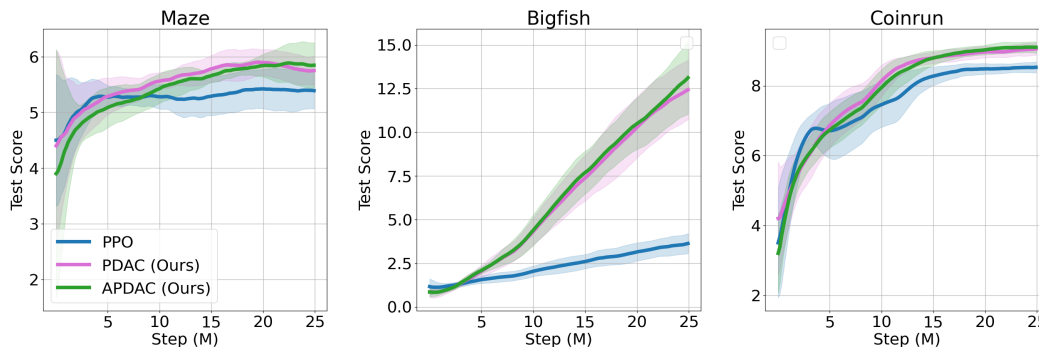


Figure 3: Ablation study for PPO, PDAC, and APDAC on the test distribution for three Procgen environments. Mean and standard deviation are calculated over 10 trials, each with a different seed.

238 5 Conclusion

239 In this work, we discuss a current systemic issue in deep RL regarding lack of performant generaliza-
 240 tion. There is an asymmetry in the information required to train the policy and value functions. Thus,
 241 fully shared policy and value networks are prone to overfitting, harming generalization. However,
 242 the policy approximation requires gradients from the value function to learn the optimal policy, so a
 243 performant model cannot completely isolate the two. Our solution, APDAC, addresses these issues in
 244 an efficient way by partially decoupling the policy and value networks while also adding attention
 245 mechanisms to each sub-network in order to efficiently identify relevant important features. Our
 246 approach contrasts with current methods in the partial separation of its networks and keeping the num-
 247 ber of convolutional layers lower than its fully separate counterpart. Our results demonstrate similar
 248 benefits of a fully decoupled approach while reducing the overall parameters and computational cost.
 249 We argue that such a compromise is a promising way forward in the pursuit of generalization in deep
 250 RL on the grounds of both performance and efficiency. The low performance gains between APDAC
 251 and the ablation can be considered as a limitation to this work, however, attention as it relates to the
 252 field of deep RL is still a growing field of study, and our work proves that there is an opportunity for
 253 the inclusion of attention to achieve generalization. A hopeful future direction is to investigate more
 254 beneficial structures for the attention mechanism.

255 References

- 256 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare,
 257 Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control
 258 through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- 259 Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim
 260 Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement
 261 learning. In *International conference on machine learning*, pages 1928–1937. PMLR, 2016.
- 262 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
 263 maximum entropy deep reinforcement learning with a stochastic actor. In *International conference
 264 on machine learning*, pages 1861–1870. PMLR, 2018.
- 265 Karl W Cobbe, Jacob Hilton, Oleg Klimov, and John Schulman. Phasic policy gradient. In *Interna-
 266 tional Conference on Machine Learning*, pages 2020–2027. PMLR, 2021.
- 267 Karl Cobbe, Chris Hesse, Jacob Hilton, and John Schulman. Leveraging procedural generation
 268 to benchmark reinforcement learning. In *International conference on machine learning*, pages
 269 2048–2056. PMLR, 2020.
- 270 Jake Grigsby and Yanjun Qi. Measuring visual generalization in continuous control from pixels.
 271 *CoRR*, abs/2010.06740, 2020. URL <https://arxiv.org/abs/2010.06740>.

- 272 Niels Justesen, Ruben Rodriguez Torrado, Philip Bontrager, Ahmed Khalifa, Julian Togelius, and
273 Sebastian Risi. Illuminating generalization in deep reinforcement learning through procedural
274 level generation. *arXiv preprint arXiv:1806.10729*, 2018.
- 275 Roberta Raileanu and Rob Fergus. Decoupling value and policy for generalization in reinforcement
276 learning. *arXiv preprint arXiv:2102.10330*, 2021.
- 277 Aravind Rajeswaran, Kendall Lowrey, Emanuel Todorov, and Sham Kakade. Towards generalization
278 and simplicity in continuous control. *arXiv preprint arXiv:1703.02660*, 2017.
- 279 Roberta Raileanu and Tim Rocktäschel. Ride: Rewarding impact-driven exploration for procedurally-
280 generated environments. *arXiv preprint arXiv:2002.12292*, 2020.
- 281 Maximilian Igl, Kamil Ciosek, Yingzhen Li, Sebastian Tschitschek, Cheng Zhang, Sam Devlin,
282 and Katja Hofmann. Generalization in reinforcement learning with selective noise injection and
283 information bottleneck. *arXiv preprint arXiv:1910.12911*, 2019a.
- 284 Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. Quantifying generalization
285 in reinforcement learning. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings
286 of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine
287 Learning Research*, pages 1282–1289. PMLR, 09–15 Jun 2019. URL [https://proceedings.
288 mlr.press/v97/cobbe19a.html](https://proceedings.mlr.press/v97/cobbe19a.html).
- 289 Maximilian Igl, Kamil Ciosek, Yingzhen Li, Sebastian Tschitschek, Cheng Zhang, Sam Devlin,
290 and Katja Hofmann. Generalization in reinforcement learning with selective noise injection and
291 information bottleneck. *arXiv preprint arXiv:1910.12911*, 2019b.
- 292 Tianyang Hu, Wenjia Wang, Cong Lin, and Guang Cheng. Regularization matters: A nonparamet-
293 ric perspective on overparametrized neural network. In *International Conference on Artificial
294 Intelligence and Statistics*, pages 829–837. PMLR, 2021.
- 295 Maximilian Igl, Gregory Farquhar, Jelena Luketina, Wendelin Boehmer, and Shimon Whiteson. The
296 impact of non-stationarity on generalisation in deep reinforcement learning. *arXiv e-prints*, pages
297 arXiv–2006, 2020.
- 298 Kaixin Wang, Bingyi Kang, Jie Shao, and Jiashi Feng. Improving generalization in reinforcement
299 learning with mixture regularization. *arXiv preprint arXiv:2010.10814*, 2020.
- 300 Bogdan Mazouze, Ahmed M Ahmed, Patrick MacAlpine, R Devon Hjelm, and Andrey Kolobov.
301 Cross-trajectory representation learning for zero-shot generalization in rl. *arXiv preprint
302 arXiv:2106.02193*, 2021.
- 303 Shariq Iqbal and Fei Sha. Actor-attention-critic for multi-agent reinforcement learning. In *Interna-
304 tional Conference on Machine Learning*, pages 2961–2970. PMLR, 2019.
- 305 Dichao Hu. An introductory survey on attention mechanisms in nlp problems. In *Proceedings of SAI
306 Intelligent Systems Conference*, pages 432–448. Springer, 2019.
- 307 Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jonathon
308 Shlens. Stand-alone self-attention in vision models. *arXiv preprint arXiv:1906.05909*, 2019.
- 309 Pengda Qin, Weiran Xu, and Jun Guo. Designing an adaptive attention mechanism for relation
310 classification. In *2017 International Joint Conference on Neural Networks (IJCNN)*, pages 4356–
311 4362. IEEE, 2017.
- 312 Zeyang Lei, Yujiu Yang, Min Yang, and Yi Liu. A multi-sentiment-resource enhanced attention
313 network for sentiment classification. *arXiv preprint arXiv:1807.04990*, 2018.
- 314 Jun Fu, Jing Liu, Haijie Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. Dual attention
315 network for scene segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision
316 and Pattern Recognition*, pages 3146–3154, 2019.
- 317 Irwan Bello, Barret Zoph, Ashish Vaswani, Jonathon Shlens, and Quoc V Le. Attention augmented
318 convolutional networks. In *Proceedings of the IEEE/CVF international conference on computer
319 vision*, pages 3286–3295, 2019.

- 320 Elaheh Barati and Xuewen Chen. An actor-critic-attention mechanism for deep reinforcement
321 learning in multi-view environments. *arXiv preprint arXiv:1907.09466*, 2019.
- 322 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
323 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 324 Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Vlad Mnih, Tom Ward, Yotam Doron,
325 Vlad Firoiu, Tim Harley, Iain Dunning, et al. Impala: Scalable distributed deep-rl with importance
326 weighted actor-learner architectures. In *International Conference on Machine Learning*, pages
327 1407–1416. PMLR, 2018.
- 328 Ilya Kostrikov. Pytorch implementations of reinforcement learning algorithms. *GitHub repository*,
329 2018.

330 **A Appendix**

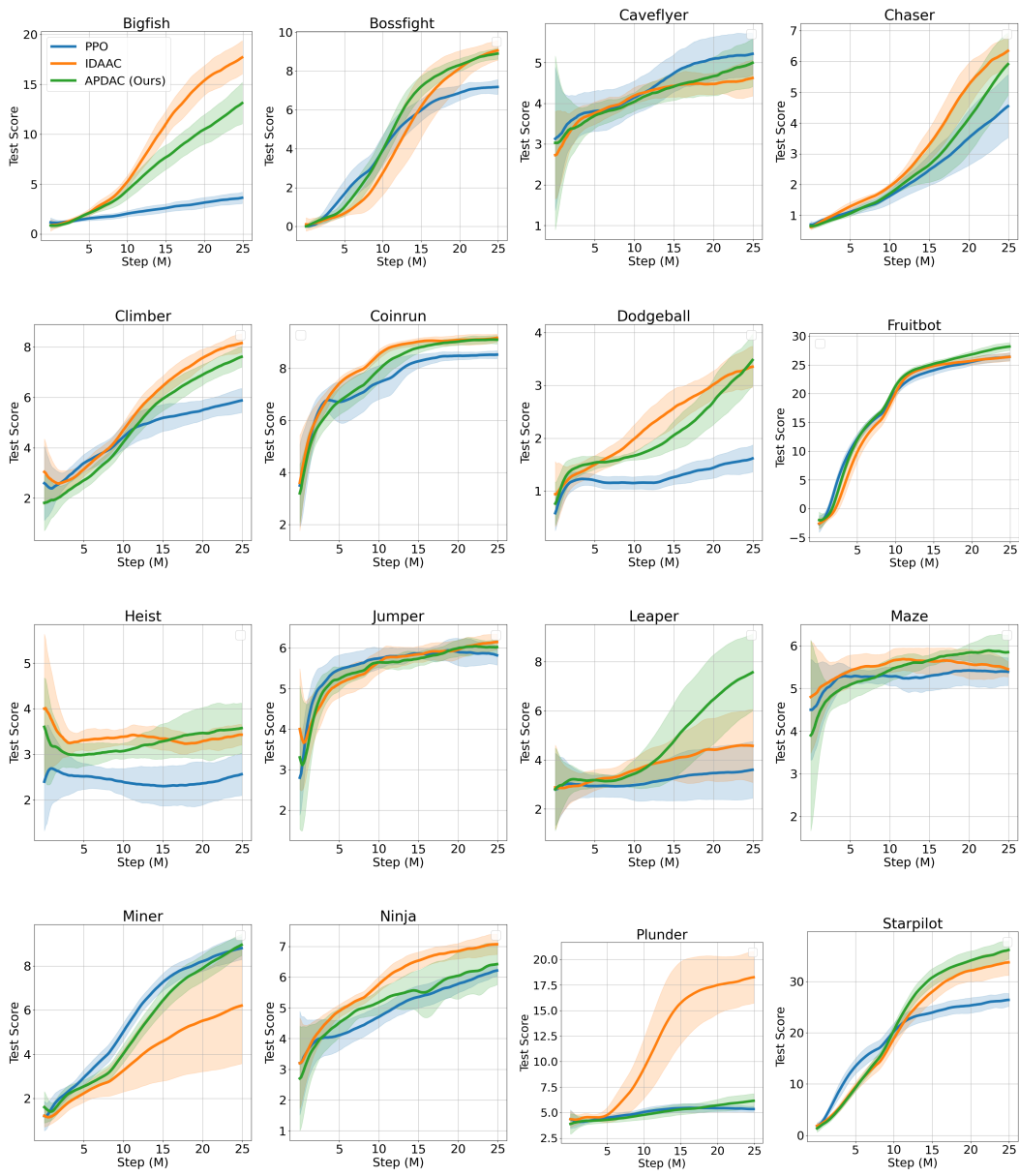


Figure 4: Testing performance for PPO, IDAAC, and APDAC across the entire Progen benchmark. Mean and standard deviation are calculated over 10 trials, each with a different seed.