

# Towered Actor Critic for Handling Multiple Action Types in Reinforcement Learning For Drug Discovery

Sai Krishna Gottipati,<sup>1</sup> Yashaswi Pathak,<sup>2</sup> Boris Sattarov,<sup>1</sup> Sahir,<sup>3</sup> Rohan Nuttall,<sup>3</sup> Mohammad Amini,<sup>1</sup> Matthew Taylor,<sup>3,4</sup> Sarath Chandar,<sup>5,6,7</sup>

<sup>1</sup> 99andBeyond

<sup>2</sup> International Institute of Information Technology, Hyderabad

<sup>3</sup> Department of Computing Science, University of Alberta

<sup>4</sup> Alberta Machine Intelligence institute

<sup>5</sup> Mila - Quebec AI Institute

<sup>6</sup> Canada CIFAR AI Chair

<sup>7</sup> Ecole Polytechnique Montréal

saikrishnagv1996@gmail.com

## Abstract

1 Reinforcement learning (RL) has made significant progress  
2 in both abstract and real-world domains, but the majority of  
3 state-of-the-art algorithms deal only with monotonic actions.  
4 However, some applications require agents to reason over  
5 different types of actions. Our application simulates reaction-  
6 based molecule generation, used as part of the drug discovery  
7 pipeline, and includes both uni-molecular and bi-molecular  
8 reactions. This paper introduces a novel framework, *towered*  
9 *actor critic* (TAC), to handle multiple action types. The TAC  
10 framework is general in that it is designed to be combined with  
11 any existing RL algorithms for continuous action space. We  
12 combine it with TD3 to empirically obtain significantly better  
13 results than existing methods in the drug discovery setting.  
14 TAC is also applied to RL benchmarks in OpenAI Gym and  
15 results show that our framework can improve, or at least does  
16 not hurt, performance relative to standard TD3.

## Introduction

17  
18 Reinforcement Learning (RL) has pushed the frontiers of  
19 various domains such as robotics (OpenAI et al. 2019; Kahn,  
20 Abbeel, and Levine 2020), game playing agents (Silver et al.  
21 2017; Anthony, Tian, and Barber 2017), and economics  
22 (Zheng et al. 2020). RL is a core problem of artificial in-  
23 telligence where an agent is situated in an environment and  
24 must learn to adapt its policy to maximize a reward signal. At  
25 each time step  $t$ , the agent interacts with an MDP described  
26 by the tuple  $M = (\mathcal{S}, \mathcal{A}, \mathcal{P}, r)$ , where  $\mathcal{S}$  is the state space,  
27  $\mathcal{A}$  is the action space,  $\mathcal{P}$  is the state transition function, and  
28  $r$  is the reward function. The agent observes a state  $s \in \mathcal{S}$ ,  
29 chooses an action  $a \in \mathcal{A}$  according to its policy  $\pi : \mathcal{S} \rightarrow$   
30  $\mathcal{A}$ , and then receives a reward and next state  $s' \in \mathcal{S}$  from  
31 the environment. Most of the existing RL algorithms are fo-  
32 cused on choosing “monotonic” actions i.e, the action space  
33 is restricted to  $\mathcal{A} \in \mathcal{R}^m$ , where  $m$  is the dimension of action  
34 space. However, several complex and crucial applications in  
35 the real world, like reaction-based molecule generation (as

part of the drug discovery pipeline) (Gottipati et al. 2020b),  
symbolic regression (Udrescu and Tegmark 2020) need to  
accommodate multiple action types (uni-molecular templates  
and bi-molecular templates in the case of molecule genera-  
tion). For bi-molecular templates, one needs to further choose  
a molecule (i.e, reactant) selection action corresponding to  
the template selection action. Existing RL algorithms do not  
implicitly exploit such multiple action types and their hier-  
archical structure. While hierarchical actor critic (Levy et al.  
2017; Röder et al. 2020) or option critic (Bacon, Harb, and  
Precup 2016) deal with breaking action sequences over long  
time horizons into clear spatial or temporal sub-goals, they do  
not necessarily deal with having to choose between multiple  
action types at every time step. To achieve this, we remodel  
the critic to process multiple action types and combine them  
at an intermediate level to predict a hypothetical next state. It  
is then followed by another neural network to compute the  
value function  $V(s')$  of that hypothetical next state. To the  
best of our knowledge, the proposed approach of a “towered”  
actor-critic (TAC) is the only algorithm that handles multiple  
action types in this manner.

Many successful and exciting applications of deep neural  
networks and RL have been reported over the last decade in  
the field of *de novo* molecule generation. These generative  
systems were found to be extremely efficient in optimiz-  
ing structures of molecules to fit a desired properties profile  
such as drug-likeness, predicted activity against biological  
targets, lipophilicity, etc. (Gómez-Bombarelli et al. 2018;  
You et al. 2018b; Olivecrona et al. 2017; Segler et al. 2017;  
Popova, Isayev, and Tropsha 2018). However, only a few  
recent publications have identified and proposed solutions to  
the long standing challenge of embedding synthetic accessi-  
bility into the molecule generation framework. Some of them  
reported use of RL to tackle this problem (Gottipati et al.  
2020b). These systems use reaction-based transformations  
of molecules as a set of actions available to the agent. In this  
context, there are two types of reactions / transformations (ac-  
tions): uni-molecular and bi-molecular. These approaches are  
limited because they are built on existing RL algorithms and

75 do not leverage these multiple action types. Later sections  
76 elaborate on this issue and propose a novel architecture and  
77 training paradigm to demonstrate superior performance over  
78 the existing state of the art *de novo* drug design approaches.

79 While this work is focused on applying the TAC formula-  
80 tion to reaction based drug discovery, there are several other  
81 real world use cases where TAC is applicable and can boost  
82 performance. For example, in a simplistic case of symbolic  
83 regression, the mathematical operators can be classified into  
84 “unitary operators” (like abs, negation etc.) that operate on  
85 a single number and “binary operators” (like addition, sub-  
86 traction etc.) that operate on two numbers and thus, another  
87 number to be operated upon has to be chosen by the policy  
88 network. Symbolic regression along with other use cases like  
89 learning-based active localization (Chaplot, Parisotto, and  
90 Salakhutdinov 2018) are described in detail in the Appendix.

91 We further show that this approach can be applied to any  
92 Markov Decision Process (MDP), including the ones without  
93 any inherent requirement for towered structure in the policy  
94 network, and obtain superior performance on a wide range  
95 of tasks. This paper’s contributions are:

- 96 1. The introduction of a novel mechanism, towered actor  
97 critic, to train agents on MDPs with multiple action types.
- 98 2. TAC is applied to the task of reaction-based *de novo* drug  
99 design (where there are different action types at every time  
100 step) and show significant improvement over the existing  
101 state-of-the-art results.
- 102 3. TAC is also applied to standard RL tasks that do not have  
103 an inherent pyramidal/hierarchical structure of actions. Per-  
104 formance is comparable, or improved, relative to TD3 on  
105 several continuous action OpenAI Gym environments.

## 106 Related Work

107 This section summarizes the most relevant background on  
108 RL and drug discovery.

### 109 Reinforcement Learning

110 While most of the RL research is focused on choosing mono-  
111 tonic actions, there are few approaches that attempted to deal  
112 with multiple action types. Hierarchical actor critic (Levy  
113 et al. 2017; Röder et al. 2020; Masson and Konidaris 2015)  
114 learns by creating low-level and high-level sub goals for more  
115 sample efficient training. Option critic (Bacon, Harb, and Pre-  
116 cup 2016), (Kumar and Precup 2017) learns to choose tempo-  
117 rally extended actions. Some approaches (Lakshminarayanan  
118 et al. 2016) use a combination of both spatially and tempo-  
119 rally extended actions, but none of these approaches deal with  
120 having to choose actions from multiple action types at the  
121 same time step. The specific architecture used for the critic  
122 in this work is inspired by the after state MDP (Sutton and  
123 Barto 2018) and the successor representation (Dayan 1993)  
124 (that attempts to solve the reward reevaluation problem by pro-  
125 viding a way to represent the relationship between different  
126 states via storing future state occupancies) and is detailed  
127 further in the later section. In the domain of continuous action  
128 space RL, several improvements have been proposed over the  
129 original deterministic policy gradient algorithm (Silver et al.

2014). Deep deterministic policy gradient (DDPG) (Lillicrap  
et al. 2015) extended the work to include deeper networks  
and demonstrated superior results over a wide range of RL en-  
vironments. Twin delayed deep deterministic networks (TD3)  
(Fujimoto, van Hoof, and Meger 2018) tackled the problem  
of overestimation bias by employing multiple tricks includ-  
ing: using the minimum value of two critics while computing  
the temporal difference target and making delayed updates to  
the policy network.

Though the TAC framework is designed to be combined  
with any existing RL algorithm for continuous action spaces,  
this paper combines and benchmarks TAC with TD3.

### 142 *De novo* drug design

143 There has been significant progress in the application of ma-  
144 chine learning methods for *de novo* drug design. Molecule  
145 generation is a well-studied problem and has been handled  
146 by several methods like genetic or evolutionary algorithms  
(Brown et al. (2004); Jensen (2019); Ahn et al. (2020)), gen-  
147 erative models (Simonovsky and Komodakis (2018); Gómez-  
148 Bombarelli et al. (2018); Winter et al. (2019a); Jin, Barzilay,  
149 and Jaakkola (2018); Popova, Isayev, and Tropsha (2018);  
150 Olivecrona et al. (2017); Griffiths and Hernández-Lobato  
(2020)) and RL based approaches (You et al. (2018a); Zhou  
151 et al. (2018)). Although these methods perform very well  
152 on popular benchmarks, such as Guacamol (Brown et al.  
153 (2019)), the molecules proposed can be infeasible to synthe-  
154 size in the real world. This issue has also been systematically  
155 highlighted by Gao and Coley (2020), where a synthesis  
156 planning program is used to quantify how often molecules  
157 proposed by some of the popular generative models cannot  
158 be readily synthesized.

159 Recent research has tried to ameliorate this problem, such  
160 as when Bradshaw et al. (2019) used a method based on vari-  
161 ational auto-encoders to optimize for a target property using  
162 single-step reactions. A similar solution was proposed by  
163 Korovina et al. (2019), where the authors employed random  
164 selection of reactants and reaction conditions in a multi-step  
165 process to generate molecules. These generated molecules  
166 were then subjected to property evaluation. More recently,  
167 Gottipati et al. (2020b) introduced PGFS (policy gradient for  
168 forward synthesis) which attempted to generate molecules via  
169 multi-step chemical synthesis and at the same time optimized  
170 the generation towards maximising an objective function. As  
171 an extension of PGFS, Gottipati et al. (2020a) optimized for  
172 the maximum reward objective instead of the usual cumu-  
173 lative return objective that helped in slightly improving the  
174 performance.

175 PGFS operates in the realm of off-policy continuous action  
176 space. The actor module  $\Pi$  that consists of  $f$  and  $\pi$  networks  
177 predicts a reaction template and a continuous action (which  
178 is in the space defined by the feature representations of all  
179 second reactants). Specifically, the  $f$  network takes in the  
180 current state  $s$  (reactant-1  $R^{(1)}$ ) as input and outputs the best  
181 reaction template  $T$ . The  $\pi$  network then takes in both  $R^{(1)}$   
182 and  $T$  as inputs and outputs the continuous action  $a$ . The  
183 environment then takes in  $a$  and computes  $k$  closest valid  
184 second reactants ( $R^{(2)}$ ). For each of these  $R^{(2)}$ s, it computes  
185  
186

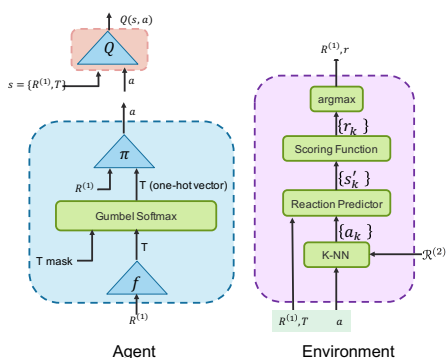


Figure 1: PGFS agent (comprising actor and critic modules) and environment.

in Equation 2, the parameters of the  $\pi$  network are getting updated even though its output is not used by the environment in case of uni-molecular reactions. Thus, every time a uni-molecular template is chosen by the  $f$ -network, the parameters of the critic and  $\pi$  networks are getting updated by arbitrary gradients. Since the percent of uni-molecular templates with respect to all available templates is sufficiently small (about 15 percent), there might not be any noticeable damage if the task under consideration doesn't need to utilize lot of uni-molecular templates. Nonetheless, the overall training becomes slower and we recognize this as a fundamental issue that couldn't be addressed with any of the existing RL algorithms.

## TAC for *de novo* drug design

The environment (underlying MDP) in reaction based molecule generation assumes that the transition function is deterministic (i.e., the environment computes only the most probable product for a given state and action). Under this assumption, the value function of the next state  $V(s_{i+1})$  is exactly equal to  $Q(s_i, a_i)$  value of the current state  $s_i$ , action  $a_i$  pair. Also, we are interested in the quality (reward) of the product  $P(s_{i+1})$  and not necessarily on the multiple ways in which the same product can be obtained using different chemical reactions (state-action pairs). Thus, drawing further inspiration from after state MDP (Sutton and Barto 2018) enables us to (hypothetically) break the "critic" down into a two step process by introducing two modules in the critic: product predictor module (that predicts the product  $P$  of the chemical reaction) and value function predictor module (that predicts the value function  $V(s_{i+1})$  of the next state  $s_{i+1}$  i.e., of the product  $P$ ). While the after state MDP learns the value functions of next states by computing all of the next states in the environment, we incorporate a (hypothetical) next state prediction module inside the critic and enable gradient propagation through this module to enable it to learn potentially more robust representations of the next state that are useful for the task under consideration. The critic architecture is also inspired from work in successor representations (Dayan 1993). Instead of representing the value function as a product of a reward vector and a successor representation matrix (or their corresponding learned variants (Barreto et al. 2017; Borsa et al. 2018; Hansen et al. 2019; Machado, Bellemare, and Bowling 2018)), we represent the action-value function  $Q(s, a)$  as a composite function of the product prediction module (i.e., next state prediction module) and the value network. Thus, when the task is changed, (i.e., when the reward to be optimized is changed) the product prediction modules need not be retrained (analogous to how the successor representations need not be retrained always).

The product predictor module has two different networks,  $U$ -net (for processing uni-molecular reactions) and  $B$ -net (for processing bi-molecular reactions).  $U$ -net takes in  $R^{(1)}$  and template  $T$  as inputs and computes the (hypothetical) product

$$P_u = U_{\theta_u}(R^{(1)}, T)$$

$B$ -net takes  $R^{(1)}$ , and action  $a$  as inputs and computes the

the corresponding product of the chemical reaction between  $R^{(1)}$  and  $R^{(2)}$ , computes the reward for the obtained product and chooses the product (next state,  $s_{t+1}$ ) that corresponds to the highest reward. All these quantities are stored in the replay buffer. Thus, by formulating forward synthesis as a continuous action space RL problem, PGFS is able to leverage TD3 (Fujimoto, van Hoof, and Meger 2018) to update the actor ( $f, \pi$ ) and critic ( $Q$ ) networks.

Specifically, the authors used the following optimizations:

$$\min L(\theta^Q) = \frac{1}{N} \sum_i |y_i - Q(R_i^{(1)}, \{T_i, a_i\})|^2 \quad (1)$$

where, the parameters  $\theta^Q$  of the critic ( $Q$  network) are updated by minimizing mean squared error between the critic value of the current state ( $R^{(1)}$ ) action ( $\{T_i, a_i\}$ ) pair and the 1-step TD target  $y_i$

$$\min L(\theta^{f,\pi}) = - \sum_i \text{Critic}(R_i^{(1)}, \text{Actor}(R_i^{(1)})) \quad (2)$$

where, the parameters  $\theta^f$  and  $\theta^\pi$  of the actor networks  $f$  and  $\pi$  respectively are updated by minimizing the negative (equivalent to maximizing) of the critic's value of current state-action pair, where the actions  $T_i$  and  $a_i$  are computed by the actor networks.

$$\min L(\theta^f) = - \sum_i (T_i^{(1)}, \log(f(R_i^{(1)}))) \quad (3)$$

where, the parameters  $\theta^f$  of the  $f$ -network were also updated by minimizing the cross entropy loss between the output of the  $f$ -network  $f(R_i^{(1)})$  and the actual valid template  $T_i^{(1)}$  chosen. This is done to encourage the  $f$ -network to predict valid templates during the initial phases of training.

PGFS used two types of reactions (reaction templates): uni-molecular reactions (also called transformation reactions) and bi-molecular reactions. If we consider Equation 1, we notice that for both types of reactions, the same update rule is being used i.e., even for uni-molecular reactions that do not require an action  $a$  (because, uni-molecular reactions do not need a second reactant  $R^{(2)}$ ), the critic is still evaluating such actions and using it to update its parameters. Similarly,

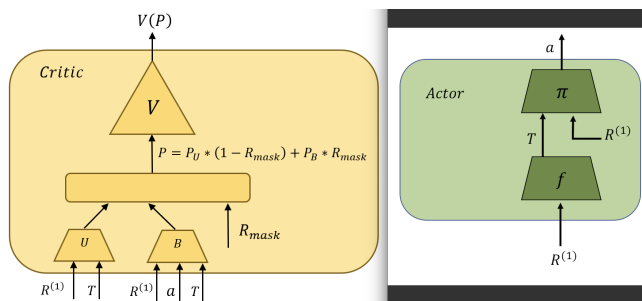


Figure 2: Illustration of TAC-FS in the context of reaction-based molecule generation for handling multiple action types (uni-molecular and bi-molecular templates)

(hypothetical) product:

$$P_b = B_{\theta_B}(R^{(1)}, a)$$

These are then combined by:

$$P = P_u \times (1 - R_{mask}) + P_b \times R_{mask},$$

where  $P$  represents the final (hypothetical) product of the chemical reaction and  $R_{mask}$  is a reaction mask that is equal to 1 when its a bi-molecular reaction and 0 if uni-molecular. The final predicted product obtained from these product predictor modules is passed through value network  $V$  to get the  $Q(s, a)$ .

Like in PGFS, the actor module takes in reactant at the current time step  $R^{(1)}$  and outputs the best template  $T$  and action  $a$ . Additionally, the actor module in the proposed approach - TAC-FS (towered actor critic for forward synthesis) also identifies the type of reaction template (whether uni-molecular or bi-molecular) and outputs the corresponding binary mask  $R_{mask}$ . The  $f$  network computes  $T$  which is a tensor that contains the probability of each template. Invalid templates are masked off using a binary template mask  $T_{mask}$  (as, only a few reaction templates are valid for a given  $R^{(1)}$ ). It is then passed through a gumbel softmax layer to obtain the template in one-hot tensor format, which indicates the final chosen template.

$$\begin{aligned} T &= f(R^{(1)}) \\ T &= T \odot T_{mask} \\ T &= \text{GumbelSoftmax}(T, \tau) \end{aligned} \quad (4)$$

Based on the chosen template, the environment identifies the type of template (either uni-molecular or bi-molecular). We then compute the action  $a$  by passing the reactant  $R^{(1)}$  and chosen template  $T$  as inputs to the  $\pi$  network. Thus, the actor module finally returns  $T$ ,  $a$  and  $R_{mask}$ .

During the ‘‘backward’’ phase, after sampling a random minibatch from the buffer, the networks are updated as follows: First, we compute the action  $a_{i+1}$ , template  $T_{i+1}$  and  $R_{mask_{i+1}}$  for the next time step given state  $(R_{i+1}^{(1)})$  using the actor module.

$$T_{i+1}, a_{i+1}, R_{mask_{i+1}} = \text{Actor-target}(R_{i+1}^{(1)}) \quad (5)$$

We then add clipped noise to the action output  $a_{i+1}$

$$a_{i+1} = a_{i+1} + \epsilon, \quad \epsilon \sim \text{clip}(\mathcal{N}(0, \bar{\sigma}), -c, c) \quad (6)$$

The one-step TD target can then be computed as:

$$y_i = r_i + \min_{j=1,2} \text{Critic-target}_j(R_{i+1}^{(1)}, \{T_{i+1}, R_{mask_{i+1}}, a_{i+1}\}) \quad (7)$$

The value loss computes the mean square error between the one-step TD target computed above and the  $Q(s, a)$  value of current state  $R_i^{(1)}$ , action  $(T_i, a_i)$  pair. Using the critic, we compute the predicted product  $P_i$  and the  $Q(s, a)$  value which are then used to compute the value loss and an auxiliary loss.

$$P_i, Q_i = \text{CRITIC}(R_i^{(1)}, \{T_i, R_{mask_i}, a_i\}) \quad (8)$$

$$\mathcal{L}_{\text{value}} = \sum_i |y_i - Q_i|^2 \quad (9)$$

We also ensure that the predicted product representation is almost the same as the actual product representation by minimizing an auxiliary loss which is mean squared error between these two representations:

$$\mathcal{L}_{\text{auxil}} = |P_i - R_{i+1}^{(1)}|^2, \quad (10)$$

where  $P_i$  is product predicted by the critic, and  $R_{i+1}^{(1)}$  is the actual product. Thus, the overall critic loss is a weighted sum of the value loss and auxiliary loss

$$\mathcal{L}_{\text{critic}} = \mathcal{L}_{\text{value}} + \alpha \mathcal{L}_{\text{auxil}}, \quad (11)$$

where  $\alpha$  is a hyper parameter.

The parameters of the critic networks ( $\theta_U, \theta_B, \theta_V$ ) are updated by minimizing the overall critic loss  $\mathcal{L}_{\text{critic}}$ . Note that we actually use two critics (to prevent overestimation bias, as elaborated by Fujimoto, van Hoof, and Meger (2018)), and follow the same steps in Equations 6-10 for both the critics to update their parameters.

The policy loss is defined as negative of the critic’s value for the current state-action pair, where the actions are computed by the current version of the actor networks.

$$\mathcal{L}_{\text{policy}} = - \sum_i \text{CRITIC}(R_i^{(1)}, \text{ACTOR}(R_i^{(1)})) \quad (12)$$

Similar to (Gottipati et al. 2020b), to enable faster training, we also aim to minimize an auxiliary loss, which is the cross entropy loss between the template tensor predicted by the  $f$ -network and the corresponding template  $T$  obtained for the reactant  $R^{(1)}$ .

$$\mathcal{L}_{\text{auxil-actor}} = - \sum_i (T_i^{(1)}, \log(f(R_i^{(1)}))) \quad (13)$$

The overall actor loss is a weighted sum of the policy loss  $\mathcal{L}_{\text{policy}}$  and the actor’s auxiliary loss  $\mathcal{L}_{\text{auxil-actor}}$ :

$$\mathcal{L}_{\text{actor}} = \mathcal{L}_{\text{policy}} + \beta \mathcal{L}_{\text{auxil-actor}}, \quad (14)$$

where  $\beta$  is a hyper parameter. The parameters of the  $f$  and  $\pi$  networks are updated by minimizing the actor loss  $\mathcal{L}_{\text{actor}}$ . The pseudo-code for the entire algorithm is provided in Algorithm 1.

---

**Algorithm 1** TAC-FS
 

---

```

1: procedure ACTOR( $R^{(1)}$ )
2:    $T \leftarrow f(R^{(1)})$ 
3:    $T \leftarrow T \odot T_{mask}$ 
4:    $T \leftarrow \text{GumbelSoftmax}(T, \tau)$ 
5:    $R_{mask} \leftarrow \text{Env.reaction.type}(T)$ 
6:    $a \leftarrow \pi(R^{(1)}, T)$ 
7:   return  $T, a, R_{mask}$ 
8: procedure CRITIC( $R^{(1)}, T, R_{mask}, a$ )
9:    $P_u \leftarrow \text{Unet}(R^{(1)}, T)$ 
10:   $P_b \leftarrow \text{Bnet}(R^{(1)}, a, T)$ 
11:   $P \leftarrow P_u \times (1 - R_{mask}) + P_b \times R_{mask}$ 
12:  return  $P, V(P)$ 
13: procedure ENV.STEP( $R^{(1)}, T, a$ )
14:   $\mathcal{R}^{(2)} \leftarrow \text{GetValidReactants}(T)$ 
15:   $\mathcal{A} \leftarrow \text{kNN}(a, \mathcal{R}^{(2)})$ 
16:   $\mathcal{R}_{t+1}^{(1)} \leftarrow \text{ForwardReaction}(R^{(1)}, T, \mathcal{A})$ 
17:   $\text{Rewards} \leftarrow \text{ScoringFunction}(\mathcal{R}_{t+1}^{(1)})$ 
18:   $r_t, R_{t+1}^{(1)}, done \leftarrow \arg \max \text{Rewards}$ 
19:  return  $R_{t+1}^{(1)}, r_t, done$ 
20: procedure BACKWARD(buffer minibatch)
21:   $T_{i+1}, a_{i+1}, R_{mask_{i+1}} \leftarrow \text{Actor-target}(R_{i+1}^{(1)})$ 
22:   $a_{i+1} \leftarrow a_{i+1} + \epsilon, \quad \epsilon \sim \text{clip}(\mathcal{N}(0, \tilde{\sigma}), -c, c)$ 
23:   $A_{i+1} \leftarrow \{T_{i+1}, R_{mask_{i+1}}, a_{i+1}\}$ 
24:   $P_{i+1}, Q_{i+1} \leftarrow \min_{j=1,2} \text{Critic-target}_j(R_{i+1}^{(1)}, A_{i+1})$ 
25:   $y_i \leftarrow r_i + \gamma Q_{i+1}$ 
26:   $P_i, Q_i \leftarrow \text{CRITIC}(R_i^{(1)}, \{T_i, R_{mask_i}, a_i\})$ 
27:   $\mathcal{L}_{\text{value}} \leftarrow \sum_i |y_i - Q_i|^2$ 
28:   $\mathcal{L}_{\text{auxil}} \leftarrow |P_i - R_{i+1}^{(1)}|^2$ 
29:   $\mathcal{L}_{\text{critic}} \leftarrow \mathcal{L}_{\text{value}} + \alpha \mathcal{L}_{\text{auxil}}$ 
30:   $\mathcal{L}_{\text{policy}} \leftarrow - \sum_i \text{CRITIC}(R_i^{(1)}, \text{ACTOR}(R_i^{(1)}))$ 
31:   $\mathcal{L}_{\text{auxil-actor}} \leftarrow - \sum_i (T_i^{(1)}, \log(f(R_i^{(1)})))$ 
32:   $\mathcal{L}_{\text{actor}} \leftarrow \mathcal{L}_{\text{policy}} + \beta \mathcal{L}_{\text{auxil-actor}}$ 
33:  min  $\mathcal{L}_{\text{actor}}, \mathcal{L}_{\text{critic}}$ 
34: procedure MAIN( $f, \pi, U, B, V$ )
35:   for episode = 1, M do
36:     sample  $R_0^{(1)}$ 
37:     for  $t = 0, N$  do
38:        $T_t, a_t, R_{mask_t} \leftarrow \text{Actor}(R_t^{(1)})$ 
39:        $R_{t+1}^{(1)}, r_t, done \leftarrow \text{ENV.STEP}(R_t^{(1)}, T_t, a_t)$ 
40:       store  $(R_t^{(1)}, T_t, a_t, R_{t+1}^{(1)}, r_t, R_{mask_t}, done)$  in
         buffer
41:        $minibatch \leftarrow \text{random sample from buffer}$ 
42:       BACKWARD}(minibatch)

```

---

### TAC-generic

333

334 There are a few challenges/differences to directly applying  
 335 the TAC formulation to environments that do not have/need  
 336 an inherent hierarchical structure in action space. In TAC-  
 337 generic, the MDP remains unchanged (i.e., the environment  
 338 still takes in the same state and action as inputs and returns

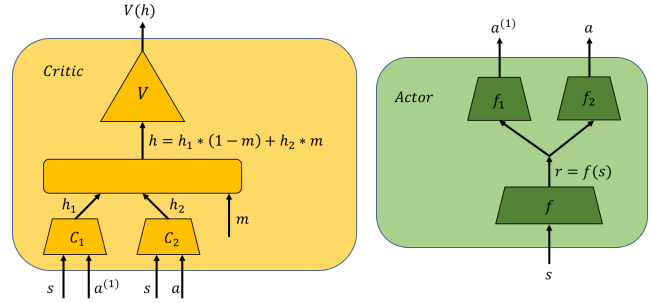


Figure 3: TAC-generic Architecture

the same reward and next state as the original MDP). Only  
 the actor and critic networks as well as the algorithm for  
 updating these networks are being changed. Let us introduce  
 a hypothetical action  $a^{(1)}$  (analogous to  $T$  in TAC-FS) of  
 dimensions the same as that of the actual action  $a$ . In the  
 actor network, unlike the drug discovery environment, we do  
 not need this other action type (hypothetical action) as inputs  
 to the higher levels of network. Thus, we can re-structure  
 the actor network as having some base layers followed by a  
 head-network  $f_1$ , and one main network  $f_2$  that outputs the  
 real action  $a$  spawning from the output of a base network  
 $f$  (which can be interpreted as a robust representation of  
 the state for the task under consideration). Thus, the actor  
 module consists of the base network  $f$  that takes in state  $s$   
 as input and computes some representation  $r = f(s)$ . The  
 first head-network  $f_1$  takes in the representation  $r$  as input  
 and computes the hypothetical action  $a^{(1)}$ . The other head  
 network  $f_2$  also takes in  $r$  as input and computes the action  
 $a = f_2(r)$ . Note that this hypothetical action  $a^{(1)}$  is not used  
 by the environment.

We construct the critic network in such a way that it takes  
 in state  $s$ , multiple action types ( $a^{(1)}$  and  $a$ ), and action mask  
 $M_a$  as inputs, and then computes the  $Q(s, a)$  value. Let us  
 consider a neural network  $C_1$  that takes in the current state  
 $s$  and the hypothetical action  $a^{(1)}$  as inputs and predicts the  
 hypothetical next state (technically, the output of  $C_1$  need not  
 be in the space of next states. It could be of any dimensions  
 and in any representation space):

$$h_1 = C_1(s, a^{(1)}).$$

One such ‘‘hypothetical next state prediction modules’’  
 (HyNeSP) can be assigned to process the action  $a$  as well:

$$h_2 = C_2(s, a).$$

The hypothetical next state  $h$  to be finally chosen from the  
 two hypothetical next states  $h_1$  and  $h_2$  is determined by the  
 action mask  $M_a$  which is a binary tensor. Similar to the drug  
 discovery setting, we introduce a new network  $V$  that takes  
 in a hypothetical next state  $h$  as input and predicts its value  
 function.

To summarize, the critic module takes in state  $s_t$ , hypothet-  
 ical action  $a^{(1)}$ , action  $a$  and action mask  $M_a$  and computes  
 the action value  $Q(s, a)$ . The  $C_i$  networks first predict the  
 hypothetical next state, These are then combined using the

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

action mask as follows:

$$h = h_1 \times M_a + h_2 \times (1 - M_a).$$

365 This hypothetical next state  $h$  is then passed through the  
366  $V$ -network to obtain the the corresponding value function  
367  $V(h)$ .

After sampling a random minibatch from the buffer, the actor and critic networks are updated as follows: The actions for the next time step  $a_{i+1}^{(1)}$  and  $a_{i+1}$  are computed by passing the next state  $s'$  through the target actor network:

$$a_{i+1}^{(1)}, a_{i+1} = \text{Actor-target}(s_{i+1}).$$

Unlike the drug discovery setting, the action mask  $M_a$  cannot be determined by the environment. Thus, we first set the action mask  $M_a = 0$  and compute the actual one-step TD targets and then compute the corresponding value loss:

$$y_i = r_i + \gamma \min_{i=1,2} \text{critic-target}(s_{i+1}, a_{i+1}^{(1)}, a_{i+1}, [0])$$

$$\mathcal{L}_{\text{value}} = \text{MSE}(y_i, \text{CRITIC}(s_i, a_i^{(1)}, a_i, [0])),$$

where MSE is the mean squared error. We then set the action mask  $M_a = 1$  and perform the same computations to get auxiliary loss  $\mathcal{L}_{\text{critic-auxil}}$ :

$$y_{i_{\text{auxil}}} = r_i + \gamma \min_{i=1,2} \text{critic-target}(s_{i+1}, a_{i+1}^{(1)}, a_{i+1}, [1])$$

$$\mathcal{L}_{\text{critic-auxil}} = \text{MSE}(y_{i_{\text{auxil}}}, \text{CRITIC}(s_i, a_i^{(1)}, a_i, [1])).$$

The overall critic loss is a weighted sum of value loss and critic’s auxiliary loss:

$$\mathcal{L}_{\text{critic}} = \mathcal{L}_{\text{value}} + \beta \times \mathcal{L}_{\text{critic-auxil}},$$

368 where  $\beta$  is a hyper parameter that could also be learned (for  
369 example, in a meta learning setup). The parameters of the  
370 critic networks ( $C_1$ ,  $C_2$ ,  $V$ ) are then updated by minimizing  
371 this critic loss  $\mathcal{L}_{\text{critic}}$ .

Similarly, the overall actor loss is computed as follows:

$$\mathcal{L}_{\text{policy}} = -\text{CRITIC}(s_i, \text{ACTOR}(s_i), [0])$$

$$\mathcal{L}_{\text{actor-auxil}} = -\text{CRITIC}(s_i, \text{ACTOR}(s_i), [1])$$

$$\mathcal{L}_{\text{actor}} = \mathcal{L}_{\text{policy}} + \beta \times \mathcal{L}_{\text{actor-auxil}}.$$

372 All the actor networks ( $f$ ,  $f_1$ ,  $f_2$ ) are updated by minimiz-  
373 ing the actor loss  $\mathcal{L}_{\text{actor}}$ . The complete algorithm is summa-  
374 rized in Algorithm 2. While the entire discussion so far has  
375 been only focused on having two action types, the algorithm  
376 can be used for any number of layers and is detailed in the  
377 Appendix.

## 378 Experiments

379 We first elaborate on the results in drug discovery setting and  
380 then go over results on several OpenAI Gym environments.

---

### Algorithm 2 TAC-generic

---

```

1: procedure ACTOR( $s$ )
2:    $r \leftarrow \text{base}(s)$ 
3:    $a^{(1)} \leftarrow f_1(r)$ 
4:    $a \leftarrow f_2(r)$ 
5:   return  $a^{(1)}, a$ 
6: procedure CRITIC( $s, a^{(1)}, a, M_a$ )
7:    $h^{(1)} \leftarrow C_1(a^{(1)}, s)$ 
8:    $h^{(2)} \leftarrow C_2(a, s)$ 
9:    $h = h^{(1)} \times M_a + h^{(2)} \times (1 - M_a)$ 
10:  return  $V(h')$ 
11: procedure ENV.STEP( $s_t, a_t$ )
12:  return  $s_{t+1}, r_t, \text{done}$ 
13: procedure BACKWARD(buffer minibatch)
14:   $a_{i+1}^{(1)}, a_{i+1} \leftarrow \text{Actor-target}(s_{i+1})$ 
15:   $y_i \leftarrow r_i + \gamma \min_{i=1,2} \text{critic-target}(s_{i+1}, a_{i+1}^{(1)}, a_{i+1}, [0])$ 
16:   $\mathcal{L}_{\text{value}} \leftarrow \text{MSE}(y_i, \text{CRITIC}(s_i, a_i^{(1)}, a_i, [0]))$ 
17:   $y_{i_{\text{auxil}}} \leftarrow r_i + \gamma \min_{i=1,2} \text{critic-target}(s_{i+1}, a_{i+1}^{(1)}, a_{i+1}, [1])$ 
18:   $\mathcal{L}_{\text{critic-auxil}} \leftarrow \text{MSE}(y_{i_{\text{auxil}}}, \text{CRITIC}(s_i, a_i^{(1)}, a_i, [1]))$ 
19:   $\mathcal{L}_{\text{critic}} = \mathcal{L}_{\text{value}} + \beta \times \mathcal{L}_{\text{critic-auxil}}$ 
20:   $\mathcal{L}_{\text{policy}} \leftarrow -\text{CRITIC}(s_i, \text{ACTOR}(s_i), [0])$ 
21:   $\mathcal{L}_{\text{actor-auxil}} \leftarrow -\text{CRITIC}(s_i, \text{ACTOR}(s_i), [1])$ 
22:   $\mathcal{L}_{\text{actor}} \leftarrow \mathcal{L}_{\text{policy}} + \beta \times \mathcal{L}_{\text{actor-auxil}}$ 
23:  min  $\mathcal{L}_{\text{actor}}, \mathcal{L}_{\text{critic}}$ 
24: procedure MAIN( $f, f_1, f_2, C_1, c_2, V$ )
25:  for episode = 1, M do
26:    for  $t = 0, N$  do
27:       $a_t^{(1)}, a_t \leftarrow \text{ACTOR}(s_t)$ 
28:       $s_{t+1}, r_t, \text{done} \leftarrow \text{ENV.STEP}(s_t, a_t)$ 
29:      store ( $s_t, a_t^{(1)}, a_t, s_{t+1}, r_t, \text{done}$ ) in buffer
30:       $\text{minibatch} \leftarrow$  random sample from buffer
31:      BACKWARD( $\text{minibatch}$ )

```

---

### De Novo Drug Design

382 We compare the performance of the proposed TAC-FS algo-  
383 rithm with five of the existing approaches: GCPN - graph  
384 convolutional policy network (You et al. 2018a), JT-VAE  
385 - junction tree variational auto encoder (Jin, Barzilay, and  
386 Jaakkola 2018), MSO - molecular swarm optimization (Win-  
387 ter et al. 2019b), PGFS, RS - random search baseline also in-  
388 troduced by (Gottipati et al. 2020b), and the Enamine dataset  
389 containing all the initial reactants used in this study. It is  
390 important to note that only RS, PGFS and TAC-FS incorpo-  
391 rates the reaction based synthesis in its molecule generation  
392 framework, thus ensuring that all the generated molecules are  
393 theoretically synthesizable based on the synthesis path that  
394 is simultaneously generated by these methods. While QED  
395 (measures drug-likeness) and clogP (measures lipophilicity)  
396 are the standard metrics that are used to measure the perfor-  
397 mance of any molecule generation algorithm, Gottipati et al.  
398 (2020b) introduced three new performance metrics (based  
399 on HIV targets) that attempt to mimic a real world use case  
400 of drug discovery. We thus compare the performance on all

381

382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400

401 these five reward metrics.

402 To ensure a fair comparison, TAC-FS is trained for only  
 403 400,000 time steps, same as that of PGFS and the results  
 404 included in the Tables 1 and 2 were obtained by fixing  $\alpha = 0$   
 405 and  $\beta = 1$  for all the TAC-FS experiments. We noted much  
 406 better performance of TAC-FS when the parameters  $\alpha$  and  
 407  $\beta$  are tuned and are reported only in the Appendix. Finally,  
 408 while in PGFS, the authors performed hyper-parameter tuning  
 409 over the state feature representations, action feature repre-  
 410 sentations, policy noise and noise clip, we tuned only for pol-  
 411 icy noise and noise clip parameters. In Table 1, we computed  
 412 the highest reward achieved over the entire training process  
 413 of 400,000 time steps. We can observe that TAC-FS improved  
 414 over the existing state-of-the-art on all the tasks. In Table 2  
 415 we computed the mean and standard deviation of the top-100  
 416 rewards achieved over the entire training process. The val-  
 417 ues reported in Table 1 and Table 2 for ENAMINEBB, RS,  
 418 GCPN(You et al. 2018a), JT-VAE(Jin, Barzilay, and Jaakkola  
 419 2018), MSO (Winter et al. 2019a), PGFS (Gottipati et al.  
 420 2020b) are taken from Gottipati et al. (2020b). While for  
 421 HIV-RT, both TAC-FS and PGFS achieved similar results,  
 422 we see that for HIV-INT and HIV-CCR5, TAC-FS performed  
 423 significantly better than existing state-of-the-art. To further  
 424 prove the effectiveness of the proposed approach, we provide  
 425 the statistics of molecules obtained under the applicability  
 426 domain of the models for HIV-targets, inference reward re-  
 427 sults, plots of actor loss, critic loss, auxiliary losses, rewards,  
 428 inference rewards in the supplementary material.

Table 1: Comparison of maximum reward achieved over the entire course of training.

Method	QED	clogP	RT	INT	CCR5
ENAMINEBB	<b>0.948</b>	5.51	7.49	6.71	8.63
RS	<b>0.948</b>	8.86	7.65	7.25	8.86
GCPN	<b>0.948</b>	7.98	7.45	6.45	8.62
JT-VAE	0.925	5.30	7.58	7.25	8.23
MSO	<b>0.948</b>	26.10	7.76	7.28	8.77
PGFS	<b>0.948</b>	27.22	7.89	7.55	9.05
<b>TAC-FS</b>	<b>0.948</b>	<b>28.97</b>	<b>7.92</b>	<b>7.75</b>	<b>9.17</b>

Table 2: Comparison of mean (and standard deviation) of top-100 rewards achieved over the entire course of training

Method	RT	INT	CCR5
ENAMINEBB	6.87 $\pm$ 0.11	6.32 $\pm$ 0.12	7.10 $\pm$ 0.27
RS	7.39 $\pm$ 0.10	6.87 $\pm$ 0.13	8.65 $\pm$ 0.06
GCPN	7.07 $\pm$ 0.10	6.18 $\pm$ 0.09	7.99 $\pm$ 0.12
JT-VAE	7.20 $\pm$ 0.12	6.75 $\pm$ 0.14	7.60 $\pm$ 0.16
MSO	7.46 $\pm$ 0.12	6.85 $\pm$ 0.10	8.23 $\pm$ 0.24
PGFS	<b>7.81</b> $\pm$ 0.03	7.16 $\pm$ 0.09	8.96 $\pm$ 0.04
<b>TAC-FS</b>	7.79 $\pm$ 0.018	<b>7.54</b> $\pm$ 0.022	<b>9.07</b> $\pm$ 0.05

## 429 TAC-generic

430 We tested on the following environments: BipedalWalker-  
 431 v3, LunarLanderContinuous-v2, MountainCarContinuous-v0,  
 432 Pendulum-v0, Hopper (PyBullet), Ant (PyBullet), Reacher

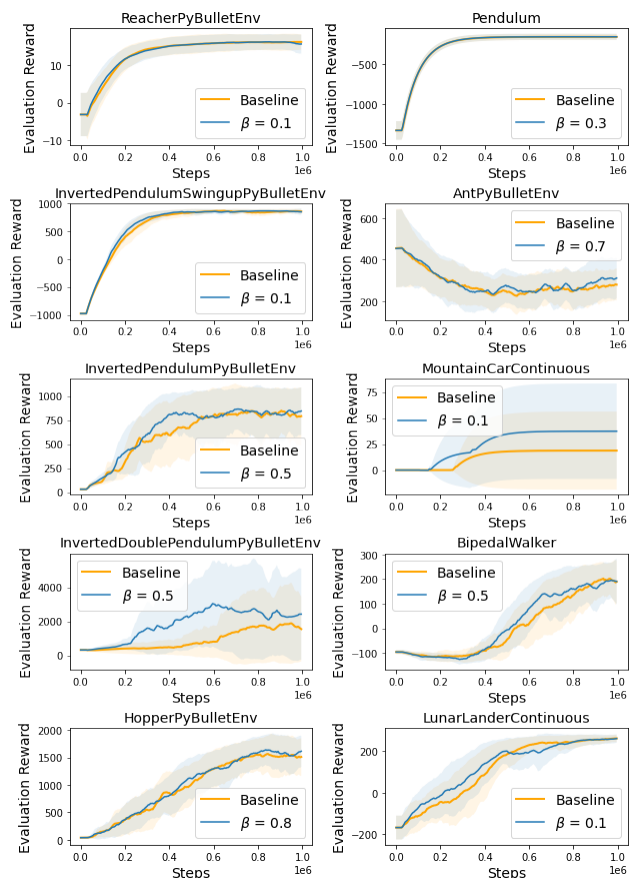


Figure 4: Performance comparison of generic TAC vs. TD3 averaged over 5 random seeds on 10 gym environments. We can observe that generic TAC performed better than, or same as TD3 on all the environments.

(PyBullet), Inverted Pendulum Swingup, Inverted Pendulum  
 and Inverted Double Pendulum. For each environment, we  
 ran the experiment for 1 million time steps and 5 random  
 seeds. Evaluation reward is computed after every 5000 steps  
 by taking the average reward achieved over 10 episodes,  
 initialized with a fixed random seed that is not used during  
 training. For a fair comparison, we did not perform any hyper-  
 parameter tuning other than for the  $\beta$  parameter. The baseline  
 comparison we used is the TD3 algorithm with exactly same  
 network architecture as that of TAC, and exactly same param-  
 eters (discount factor, policy noise, noise clip etc.). Thus, the  
 only difference is the  $\beta$  parameter. For TD3, it is fixed at zero,  
 where as for TAC, it is tuned to compute an appropriate value  
 for the task under consideration. As we can observe from  
 Figure 4, TAC performed same or better than the baseline on  
 all the 10 environments we tested.

## 449 Conclusion and Future work

450 In this work, we introduced a novel RL algorithm that solved  
 451 the severe limitation of the existing reaction based molecule  
 452 generation framework and demonstrated better results than  
 453 existing state-of-the-art on all the five tasks. We further ex-

454 amined the algorithm on MDPs that do not have any inherent  
455 hierarchical structure and demonstrated comparable or better  
456 performance over TD3 on a wide range of Gym environments.  
457 While this work was developed using the concepts of TD3  
458 algorithm, there is a potential for improvement when some  
459 of the latest advances from deep RL literature are incorpo-  
460 rated such as (Haarnoja et al. 2018; Kuznetsov et al. 2020;  
461 Wang et al. 2019; Schaul et al. 2015; Fedus et al. 2020; Sinha  
462 et al. 2020; Fujimoto, Meger, and Precup 2020). Another in-  
463 teresting avenue is to experiment on procedurally generated  
464 environments (Chevalier-Boisvert, Willems, and Pal 2018).

## 465 References

466 Ahn, S.; Kim, J.; Lee, H.; and Shin, J. 2020. Guiding Deep  
467 Molecular Optimization with Genetic Exploration.

468 Anthony, T.; Tian, Z.; and Barber, D. 2017. Thinking Fast  
469 and Slow with Deep Learning and Tree Search. *CoRR*  
470 abs/1705.08439. URL <http://arxiv.org/abs/1705.08439>.

471 Bacon, P.-L.; Harb, J.; and Precup, D. 2016. The Option-  
472 Critic Architecture.

473 Barreto, A.; Dabney, W.; Munos, R.; Hunt, J. J.; Schaul,  
474 T.; van Hasselt, H. P.; and Silver, D. 2017. Successor Fea-  
475 tures for Transfer in Reinforcement Learning. In Guyon, I.;  
476 Luxburg, U. V.; Bengio, S.; Wallach, H.; Fergus, R.; Vish-  
477 wanathan, S.; and Garnett, R., eds., *Advances in Neural In-*  
478 *formation Processing Systems 30*, 4055–4065. Curran Asso-  
479 ciates, Inc. URL [http://papers.nips.cc/paper/6994-successor-](http://papers.nips.cc/paper/6994-successor-features-for-transfer-in-reinforcement-learning.pdf)  
480 [features-for-transfer-in-reinforcement-learning.pdf](http://papers.nips.cc/paper/6994-successor-features-for-transfer-in-reinforcement-learning.pdf).

481 Borsa, D.; Barreto, A.; Quan, J.; Mankowitz, D. J.; Munos,  
482 R.; van Hasselt, H.; Silver, D.; and Schaul, T. 2018. Universal  
483 Successor Features Approximators. *CoRR* abs/1812.07626.  
484 URL <http://arxiv.org/abs/1812.07626>.

485 Bradshaw, J.; Paige, B.; Kusner, M. J.; Segler, M. H. S.;  
486 and Hernández-Lobato, J. M. 2019. A Model to Search for  
487 Synthesizable Molecules. *CoRR* abs/1906.05221.

488 Brown, N.; Fiscato, M.; Segler, M. H.; and Vaucher, A. C.  
489 2019. Guacamol: benchmarking models for de novo molec-  
490 ular design. *Journal of chemical information and modeling*  
491 59(3): 1096–1108.

492 Brown, N.; McKay, B.; Gilardoni, F.; and Gasteiger, J. 2004.  
493 A graph-based genetic algorithm and its application to the  
494 multiobjective evolution of median molecules. *Journal of*  
495 *chemical information and computer sciences* 44(3): 1079–  
496 1087.

497 Chappot, D. S.; Parisotto, E.; and Salakhutdinov, R. 2018.  
498 Active Neural Localization. In *International Conference*  
499 *on Learning Representations*. URL [https://openreview.net/](https://openreview.net/forum?id=ry6-G_66b)  
500 [forum?id=ry6-G\\_66b](https://openreview.net/forum?id=ry6-G_66b).

501 Chevalier-Boisvert, M.; Willems, L.; and Pal, S. 2018. Min-  
502 imalistic Gridworld Environment for OpenAI Gym. [https://](https://github.com/maximecb/gym-minigrid)  
503 [github.com/maximecb/gym-minigrid](https://github.com/maximecb/gym-minigrid).

504 Dayan, P. 1993. Improving Generalization for Temporal  
505 Difference Learning: The Successor Representation. *Neural*  
506 *Computation* 5(4): 613–624.

Fedus, W.; Ramachandran, P.; Agarwal, R.; Bengio, Y.;  
507 Larochelle, H.; Rowland, M.; and Dabney, W. 2020. Re-  
508 visiting Fundamentals of Experience Replay. 509

Fujimoto, S.; Meger, D.; and Precup, D. 2020. An Equiva-  
510 lence between Loss Functions and Non-Uniform Sampling  
511 in Experience Replay. 512

Fujimoto, S.; van Hoof, H.; and Meger, D. 2018. Address-  
513 ing Function Approximation Error in Actor-Critic Methods.  
514 *CoRR* abs/1802.09477. 515

Gao, W.; and Coley, C. W. 2020. The Synthesizability  
516 of Molecules Proposed by Generative Models. *Journal of*  
517 *Chemical Information and Modeling* ISSN 1549-960X. doi:  
518 10.1021/acs.jcim.0c00174. URL <http://dx.doi.org/10.1021/>  
519 [acs.jcim.0c00174](http://dx.doi.org/10.1021/acs.jcim.0c00174). 520

Gómez-Bombarelli, R.; Wei, J. N.; Duvenaud, D.; Hernández-  
521 Lobato, J. M.; Sánchez-Lengeling, B.; Sheberla, D.; Aguilera-  
522 Iparraquiere, J.; Hirzel, T. D.; Adams, R. P.; and Aspuru-  
523 Guzik, A. 2018. Automatic chemical design using a data-  
524 driven continuous representation of molecules. *ACS central*  
525 *science* 4(2): 268–276. 526

Gottipati, S. K.; Pathak, Y.; Nuttall, R.; Sahir; Chunduru, R.;  
527 Touati, A.; Subramanian, S. G.; Taylor, M. E.; and Chandar,  
528 S. 2020a. Maximum Reward Formulation In Reinforcement  
529 Learning . 530

Gottipati, S. K.; Sattarov, B.; Niu, S.; Pathak, Y.; Wei, H.;  
531 Liu, S.; Thomas, K. M. J.; Blackburn, S.; Coley, C. W.; Tang,  
532 J.; Chandar, S.; and Bengio, Y. 2020b. Learning To Nav-  
533 igate The Synthetically Accessible Chemical Space Using  
534 Reinforcement Learning. 535

Griffiths, R.-R.; and Hernández-Lobato, J. M. 2020. Con-  
536 strained Bayesian optimization for automatic chemical de-  
537 sign using variational autoencoders. *Chem. Sci.* 11: 577–586.  
538 doi:10.1039/C9SC04026A. URL <http://dx.doi.org/10.1039/>  
539 [C9SC04026A](http://dx.doi.org/10.1039/C9SC04026A). 540

Haarnoja, T.; Zhou, A.; Abbeel, P.; and Levine, S. 2018.  
541 Soft Actor-Critic: Off-Policy Maximum Entropy Deep Re-  
542 inforcement Learning with a Stochastic Actor. *CoRR*  
543 abs/1801.01290. 544

Hansen, S.; Dabney, W.; Barreto, A.; de Wiele, T. V.; Warde-  
545 Farley, D.; and Mnih, V. 2019. Fast Task Inference with Vari-  
546 ational Intrinsic Successor Features. *CoRR* abs/1906.05030.  
547 URL <http://arxiv.org/abs/1906.05030>. 548

Jensen, J. H. 2019. A graph-based genetic algorithm and  
549 generative model/Monte Carlo tree search for the exploration  
550 of chemical space. *Chemical science* 10(12): 3567–3572. 551

Jin, W.; Barzilay, R.; and Jaakkola, T. 2018. Junction  
552 Tree Variational Autoencoder for Molecular Graph Gener-  
553 ation. In Dy, J.; and Krause, A., eds., *Proceedings of the*  
554 *35th International Conference on Machine Learning*, vol-  
555 ume 80 of *Proceedings of Machine Learning Research*, 2323–  
556 2332. Stockholm, Sweden: PMLR. URL  
557 <http://proceedings.mlr.press/v80/jin18a.html>. 558

Kahn, G.; Abbeel, P.; and Levine, S. 2020. Badgr: An au-  
559 tonomous self-supervised learning-based navigation system.  
560 *arXiv preprint arXiv:2002.05700* . 561



- 562 Korovina, K.; Xu, S.; Kandasamy, K.; Neiswanger, W.; Poczos, B.; Schneider, J.; and Xing, E. P. 2019. ChemBO: Bayesian Optimization of Small Organic Molecules with Synthesizable Recommendations. *arXiv:1908.01425 [physics, stat]* URL <http://arxiv.org/abs/1908.01425>. ArXiv: 1908.01425. 615
- 563 616
- 564 617
- 565 618
- 566 619
- 567 620
- 568 Kumar, P.; and Precup, D. 2017. Multi-Timescale, Gradient Descent, Temporal Difference Learning with Linear Options. *ArXiv abs/1703.06471*. 621
- 569 622
- 570 623
- 571 Kuznetsov, A.; Shvechikov, P.; Grishin, A.; and Vetrov, D. 2020. Controlling Overestimation Bias with Truncated Mixture of Continuous Distributional Quantile Critics. 624
- 572 625
- 573 626
- 574 Lakshminarayanan, A. S.; Krishnamurthy, R.; Kumar, P.; and Ravindran, B. 2016. Option Discovery in Hierarchical Reinforcement Learning using Spatio-Temporal Clustering. *arXiv: Learning*. 627
- 575 628
- 576 629
- 577 630
- 578 Levy, A.; Konidaris, G.; Platt, R.; and Saenko, K. 2017. Learning Multi-Level Hierarchies with Hindsight. 631
- 579 632
- 580 Lillicrap, T. P.; Hunt, J. J.; Pritzel, A.; Heess, N. M. O.; Erez, T.; Tassa, Y.; Silver, D.; and Wierstra, D. 2015. Continuous control with deep reinforcement learning. *CoRR abs/1509.02971*. 633
- 581 634
- 582 635
- 583 636
- 584 Machado, M. C.; Bellemare, M. G.; and Bowling, M. 2018. Count-Based Exploration with the Successor Representation. *CoRR abs/1807.11622*. URL <http://arxiv.org/abs/1807.11622>. 637
- 585 638
- 586 639
- 587 Masson, W.; and Konidaris, G. D. 2015. Reinforcement Learning with Parameterized Actions. *CoRR abs/1509.01644*. URL <http://arxiv.org/abs/1509.01644>. 640
- 588 641
- 589 642
- 590 Olivecrona, M.; Blaschke, T.; Engkvist, O.; and Chen, H. 2017. Molecular de-novo design through deep reinforcement learning. *Journal of cheminformatics* 9(1): 48. 643
- 591 644
- 592 645
- 593 OpenAI; Akkaya, I.; Andrychowicz, M.; Chociej, M.; Litwin, M.; McGrew, B.; Petron, A.; Paino, A.; Plappert, M.; Powell, G.; Ribas, R.; Schneider, J.; Tezak, N.; Tworek, J.; Welinder, P.; Weng, L.; Yuan, Q.; Zaremba, W.; and Zhang, L. 2019. Solving Rubik's Cube with a Robot Hand. 646
- 594 647
- 595 648
- 596 649
- 597 650
- 598 Popova, M.; Isayev, O.; and Tropsha, A. 2018. Deep reinforcement learning for de novo drug design. *Science advances* 4(7): eaap7885. 651
- 599 652
- 600 653
- 601 Röder, F.; Epe, M.; Nguyen, P. D. H.; and Wermter, S. 2020. Curious Hierarchical Actor-Critic Reinforcement Learning. 654
- 602 655
- 603 656
- 604 657
- 605 Segler, M. H.; Kogej, T.; Tyrchan, C.; and Waller, M. P. 2017. Generating focused molecule libraries for drug discovery with recurrent neural networks. *ACS central science* 4(1): 120–131. 658
- 606 659
- 607 660
- 608 661
- 609 Silver, D.; Hubert, T.; Schrittwieser, J.; Antonoglou, I.; Lai, M.; Guez, A.; Lanctot, M.; Sifre, L.; Kumaran, D.; Graepel, T.; Lillicrap, T. P.; Simonyan, K.; and Hassabis, D. 2017. Mastering Chess and Shogi by Self-Play with a General Reinforcement Learning Algorithm. *CoRR abs/1712.01815*. URL <http://arxiv.org/abs/1712.01815>. 662
- 610 663
- 611 664
- 612 665
- 613 666
- 614 667
- Silver, D.; Lever, G.; Heess, N.; Degris, T.; Wierstra, D.; and Riedmiller, M. 2014. Deterministic Policy Gradient Algorithms. In *Proceedings of the 31st International Conference on International Conference on Machine Learning - Volume 32, ICML'14*, I–387–I–395. JMLR.org. 615
- 616 617
- 618 619
- 620 621
- 622 623
- 624 625
- 626 627
- 628 629
- 629 630
- 630 631
- 631 632
- 632 633
- 633 634
- 634 635
- 635 636
- 636 637
- 637 638
- 638 639
- 639 640
- 640 641
- 641 642
- 642 643
- 643 644
- 644 645
- 645 646
- 646 647
- 647 648
- 648 649
- 649 650
- 650 651
- 651 652
- 652 653
- 653 654
- 654 655
- 655 656
- 656 657
- 657 658
- 658 659
- 659 660
- 660 661
- 661 662
- 662 663
- 663 664
- 664 665
- 665 666
- 666 667
- 667 668
- 668 669
- 669 670
- 670 671
- 671 672
- 672 673
- 673 674
- 674 675
- 675 676
- 676 677
- 677 678
- 678 679
- 679 680
- 680 681
- 681 682
- 682 683
- 683 684
- 684 685
- 685 686
- 686 687
- 687 688
- 688 689
- 689 690
- 690 691
- 691 692
- 692 693
- 693 694
- 694 695
- 695 696
- 696 697
- 697 698
- 698 699
- 699 700
- 700 701
- 701 702
- 702 703
- 703 704
- 704 705
- 705 706
- 706 707
- 707 708
- 708 709
- 709 710
- 710 711
- 711 712
- 712 713
- 713 714
- 714 715
- 715 716
- 716 717
- 717 718
- 718 719
- 719 720
- 720 721
- 721 722
- 722 723
- 723 724
- 724 725
- 725 726
- 726 727
- 727 728
- 728 729
- 729 730
- 730 731
- 731 732
- 732 733
- 733 734
- 734 735
- 735 736
- 736 737
- 737 738
- 738 739
- 739 740
- 740 741
- 741 742
- 742 743
- 743 744
- 744 745
- 745 746
- 746 747
- 747 748
- 748 749
- 749 750
- 750 751
- 751 752
- 752 753
- 753 754
- 754 755
- 755 756
- 756 757
- 757 758
- 758 759
- 759 760
- 760 761
- 761 762
- 762 763
- 763 764
- 764 765
- 765 766
- 766 767
- 767 768
- 768 769
- 769 770
- 770 771
- 771 772
- 772 773
- 773 774
- 774 775
- 775 776
- 776 777
- 777 778
- 778 779
- 779 780
- 780 781
- 781 782
- 782 783
- 783 784
- 784 785
- 785 786
- 786 787
- 787 788
- 788 789
- 789 790
- 790 791
- 791 792
- 792 793
- 793 794
- 794 795
- 795 796
- 796 797
- 797 798
- 798 799
- 799 800
- 800 801
- 801 802
- 802 803
- 803 804
- 804 805
- 805 806
- 806 807
- 807 808
- 808 809
- 809 810
- 810 811
- 811 812
- 812 813
- 813 814
- 814 815
- 815 816
- 816 817
- 817 818
- 818 819
- 819 820
- 820 821
- 821 822
- 822 823
- 823 824
- 824 825
- 825 826
- 826 827
- 827 828
- 828 829
- 829 830
- 830 831
- 831 832
- 832 833
- 833 834
- 834 835
- 835 836
- 836 837
- 837 838
- 838 839
- 839 840
- 840 841
- 841 842
- 842 843
- 843 844
- 844 845
- 845 846
- 846 847
- 847 848
- 848 849
- 849 850
- 850 851
- 851 852
- 852 853
- 853 854
- 854 855
- 855 856
- 856 857
- 857 858
- 858 859
- 859 860
- 860 861
- 861 862
- 862 863
- 863 864
- 864 865
- 865 866
- 866 867
- 867 868
- 868 869
- 869 870
- 870 871
- 871 872
- 872 873
- 873 874
- 874 875
- 875 876
- 876 877
- 877 878
- 878 879
- 879 880
- 880 881
- 881 882
- 882 883
- 883 884
- 884 885
- 885 886
- 886 887
- 887 888
- 888 889
- 889 890
- 890 891
- 891 892
- 892 893
- 893 894
- 894 895
- 895 896
- 896 897
- 897 898
- 898 899
- 899 900
- 900 901
- 901 902
- 902 903
- 903 904
- 904 905
- 905 906
- 906 907
- 907 908
- 908 909
- 909 910
- 910 911
- 911 912
- 912 913
- 913 914
- 914 915
- 915 916
- 916 917
- 917 918
- 918 919
- 919 920
- 920 921
- 921 922
- 922 923
- 923 924
- 924 925
- 925 926
- 926 927
- 927 928
- 928 929
- 929 930
- 930 931
- 931 932
- 932 933
- 933 934
- 934 935
- 935 936
- 936 937
- 937 938
- 938 939
- 939 940
- 940 941
- 941 942
- 942 943
- 943 944
- 944 945
- 945 946
- 946 947
- 947 948
- 948 949
- 949 950
- 950 951
- 951 952
- 952 953
- 953 954
- 954 955
- 955 956
- 956 957
- 957 958
- 958 959
- 959 960
- 960 961
- 961 962
- 962 963
- 963 964
- 964 965
- 965 966
- 966 967
- 967 968
- 968 969
- 969 970
- 970 971
- 971 972
- 972 973
- 973 974
- 974 975
- 975 976
- 976 977
- 977 978
- 978 979
- 979 980
- 980 981
- 981 982
- 982 983
- 983 984
- 984 985
- 985 986
- 986 987
- 987 988
- 988 989
- 989 990
- 990 991
- 991 992
- 992 993
- 993 994
- 994 995
- 995 996
- 996 997
- 997 998
- 998 999
- 999 1000