

---

# Searching for High-Value Molecules Using Reinforcement Learning and Transformers

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Reinforcement learning (RL) over text representations can be effective for finding  
2 high-value policies that can search over graphs. However, RL requires careful  
3 structuring of the search space and algorithm design to be effective in this  
4 challenge. Through extensive experiments, we explore how different design  
5 choices for text grammar and algorithmic choices for training can affect an RL  
6 policy’s ability to generate molecules with desired properties. We arrive at a new  
7 RL-based molecular design algorithm (ChemRLformer) and perform a thorough  
8 analysis using 25 molecule design tasks, including computationally complex  
9 protein docking simulations. From this analysis, we discover unique insights  
10 in this problem space and show that ChemRLformer achieves state-of-the-art  
11 performance while being more straightforward than prior work by demystifying  
12 which design choices are actually helpful for text-based molecule design.

## 13 1 Introduction

14 Molecular discovery can have a significant impact on our society, however, the vast search space  
15 makes it challenging to find high-value molecules. The potential of reinforcement learning (RL)  
16 methods to discover new, high-value molecules has resulted in a series of research work performed  
17 by RL researchers focusing on learning policies as graph neural networks (GNNs) [You et al., 2018,  
18 Zhou et al., 2019, Jin et al., 2020, Fu et al., 2022, Yang et al., 2021, Bengio et al., 2021]. In this  
19 formulation, the RL policy is trained to add atoms and bonds to a molecular graph representation.  
20 In this formulation there is a one-to-one mapping between molecules and their graph representation,  
21 making it easier to construct state and action spaces with Markovian dynamics. However, the action  
22 space in the graph formulation is vast as it consists of the product of candidate attachment positions  
23 and candidate attachment sequences. Graph-based data structures (such as adjacency matrices,  
24 trees, etc.) are a powerful representation used to describe a number of design problems, including  
25 social networks [Tan et al., 2019], transportation networks [Wang and Tang, 2021], recommendation  
26 systems [Chen et al., 2021b], and combinatorial optimization problems [Khadka et al., 2020, Miret  
27 et al., 2022] have been popular in this design space. However, GNNs are often difficult to train  
28 [Chen et al., 2022] and cannot readily take advantage of large-scale text data sets that effectively  
29 describe molecular structures and properties.

30 In order to take advantage of the richness of text-based representations for molecules, one can  
31 formulate the molecular search problem as the construction of tokens in a sequence that become  
32 a molecular text. The molecular texts formulated by common text-based representations, such as  
33 SMILES [Weininger, 1988] and SELFIES [Krenn et al., 2020], can then be converted into molecular  
34 graphs with cheminformatics libraries [Landrum et al., 2013] using their respective encoding and  
35 decoding rules. However, the text-based representation can be more difficult to formulate as an  
36 MDP since there is not always an exact one-to-one mapping between texts and molecules. In

37 fact, the text-to-molecule conversion can be many-to-one, where the complexity of the dynamics  
 38 in the MDP given by many-to-one mappings is non-trivial. On the other hand, the action space in  
 39 molecular text design can be significantly reduced given the rules of text construction imposed by  
 40 a given representation. Moreover, formulating molecule discovery as sequence-generation has the  
 41 potential to capitalize on recent successes in natural language modeling [Brown et al., 2020a].

42 In this paper, we perform a detailed empirical study of molecular discovery using text-based RL  
 43 across more than 25 molecular properties relevant for drug-discovery, including docking simulation  
 44 for molecular ligands [García-Ortegón et al., 2022, Huang et al., 2022] and develop our own  
 45 algorithm (MoLRL) based on state-of-the-art literature as shown in Table 1. In our experiments, we  
 46 evaluate two molecular text representations (SMILES, SELFIES) and the use of three neural network  
 47 architectures (Multi-Layer Perceptron [Bengio et al., 2003], Recurrent Neural Network [Schmidt,  
 48 2019], Transformer [Vaswani et al., 2017]) pretrained on 5 datasets of varying quality and sizes.  
 49 We create ChemRLformer that achieves the highest performance across these tasks while being  
 50 much simpler than previous text-based RL algorithms [Blaschke et al., 2020a, Gao et al., 2022].  
 51 Via our detailed ablation study, we construct ChemRLformer and find that pretraining on *aligned*  
 52 datasets can significantly improve performance across all molecular design tasks, even exceeding  
 53 the performance of agents pretrained on 100 times larger datasets. We also show that targeted  
 54 algorithmic design, such as hill-climbing in the replay buffer and regularization, further increases  
 55 the performance of ChemRLformer. To the best of our knowledge, ChemRLformer is the largest  
 56 analysis of text-based RL methods for molecule discovery.

Table 1: Table showing conceptual comparisons of various text based molecular optimization methods. MoLRL combines the most successful elements of prior work.

Method	Text Representation	RL	Architecture	Pretraining	Algorithmic Components
Literature Methods					
SMILES-VAE [Gómez-Bombarelli et al., 2018]	SMILES	✗	VAE	✓	Maximum Likelihood
SMILES-LSTM [Brown et al., 2019]	SMILES	✗	LSTM	✓	Maximum Likelihood
BOSS [Moss et al., 2020]	SMILES	✗	VAE	✗	Bayesian Optimization
REINVENT [Blaschke et al., 2020a]	SMILES	✓	GRU	✓	Replay buffer, KL
REINVENT 2.0 [Blaschke et al., 2020b]	SMILES	✓	GRU	✓	HC-Replay buffer, Log p, KL
STONED [Nigam et al., 2021]	SELFIES	✗	FC	✗	Genetic algorithm
Pasithea [Shen et al., 2021]	SELFIES	✗	FC	✗	Deep dreaming
<b>ChemRLformer (Ours)</b>	SMILES, SELFIES	✓	Transformer, FC	✓	Replay buffer, KL

## 57 2 Related Work

58 **RL for Design and Discovery:** Many methods in diverse fields leverage RL to help augment a prior  
 59 design method to improve performance [Yu et al., 2018, Schaff et al., 2019]. Other methods have  
 60 explicitly included the design process in the RL loop by training design problems together [Chen  
 61 et al., 2021a, Ha, 2019, Luck et al., 2020, Kumar et al., 2022] with most prior work focusing on  
 62 robot and agent design, not molecular design. Our molecular design work creates an autoregressive  
 63 structure that grows the size of the state as the agent acts in the environment.

64 **Molecular Discovery Using Sequence-Based Methods:** Sequence-based methods treat molecular  
 65 design as a sequence of tokens that get concatenated in order. Generative models for sequence-based  
 66 methods span a diverse range, including variational autoencoders (VAEs) [Gómez-Bombarelli et al.,  
 67 2018, Alperstein et al., 2019], recurrent neural networks (RNNs) [Gupta et al., 2018, Bjerrum and  
 68 Threlfall, 2017, Grisoni et al., 2020, Flam-Shepherd et al., 2022] and transformer models [Wang  
 69 et al., 2019, Fabian et al., 2020, Edwards et al., 2022a, Zeng et al., 2022, Taylor et al., 2022]. The  
 70 general procedure for all the above methods is to perform self-supervised generative learning to  
 71 sample molecules similar to the original dataset. MoLRL can also make use of pretrained generative  
 72 models, which we then fine-tune using reinforcement learning to produce enhanced molecules.

73 **Molecular Discovery Using Search-Based Methods:** Although sequence-based molecule  
 74 generation methods often provide a more structured way of learning molecular distributions, search-  
 75 based methods generally have the advantage of being able to directly find molecules based on  
 76 a desired property. Although a wide range of graph-based RL methods [You et al., 2018, Zhou  
 77 et al., 2019, Jin et al., 2020, Fu et al., 2022, Yang et al., 2021, Bengio et al., 2021] for optimizing  
 78 molecules exist, graph-based state representations introduce significant complexity to the RL  
 79 problem formulation, both in the transition dynamics and action space. By contrast, text-based

80 methods are simpler and also relatively under-explored, motivating our focus on these methods in  
81 this work. Moreover, recent work [Cieplinski et al., 2021, Gao et al., 2022] has shown that an older  
82 text-based method REINVENT [Olivecrona et al., 2017] outperforms more complex graph-based  
83 RL methods. Some limited extensions to Olivecrona et al. [2017] have been explored, including  
84 experimenting with a newer molecular grammar designed for robust molecule generation [Gao et al.,  
85 2022]. However, there has been limited work proposing the use of language models and text-based  
86 RL for molecular discovery. Additionally, there have been limited efforts to incorporate recent  
87 advancements from the language modeling domain into these methods. For example, the a character-  
88 level LSTM network architecture used in Olivecrona et al. [2017], has not been revisited despite  
89 significant recent advances in sequence modeling [Vaswani et al., 2017, Brown et al., 2020b].

### 90 3 Background

91 The algorithms detailed in this paper are built on top of a foundation of reinforcement learning,  
92 text-based molecule representations, and language modeling.

93 **Reinforcement Learning:** Reinforcement learning can be used to learn policies for sequential  
94 decision-making problems. Policies are optimized based on an environment that is described as  
95 a Markov Decision Process (MDP). A discrete MDP is defined by the tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{T}, r, \gamma \rangle$  where  
96  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $T : \mathcal{S} \times \mathcal{A} \times \mathcal{S}' \rightarrow [0, 1]$  is the transition function,  
97  $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{R}$  is the reward function and  $\gamma$  is the discount rate.

98 For actions  $a_t \in \mathcal{A}$  and states  $s_t \in \mathcal{S}$ , the goal of reinforcement learning is to learn a policy  $\pi_\theta(a_t|s_t)$   
99 which maps states to actions, such that:

$$\pi_\theta(a_t|s_t) = \arg \max_{\theta} \mathbf{E}_{p(\tau|\theta)} \left[ \sum_{t=0}^T \gamma^t r(s_t, a_t) \right] \quad (1)$$

100 where  $p(\tau|\theta)$  is the distribution over trajectories induced by  $\pi_\theta$  and the transition function  $\mathcal{T}$ .

101 **Text representations for molecules:** Molecules are most naturally described using a graph  
102 structure of atoms and bonds. However, graph-based deep learning models can be difficult to train,  
103 especially at large scale [Dwivedi et al., 2022, Geisler et al., 2023]. Recent works have proposed a  
104 variety of text representations for molecules [Weininger, 1988, Krenn et al., 2020, Heller et al., 2013,  
105 Krenn et al., 2022, Cheng et al., 2023], each having their distinct advantages and shortcomings.  
106 In this study, we focus on the two most commonly used representations: SMILES [Weininger,  
107 1988] and SELFIES [Krenn et al., 2020]. Any text representation for molecules consists of a set  
108 of valid tokens, which may represent individual atoms or special characters that imply the presence  
109 of certain structures, as well as the encoding and decoding rules needed to convert between the text  
110 representation and the graph representation of a molecule. Valid texts under a grammar are those  
111 which respect both the vocabulary and the encoding/decoding rules for that grammar and, hence, can  
112 be converted into a graph representation of a molecule. SELFIES, which was developed in response  
113 to the tendency for SMILES-based deep learning models [Gómez-Bombarelli et al., 2018, Jin et al.,  
114 2018] to generate invalid molecular texts, has the useful property of providing a conversion for *any*  
115 text into a graph corresponding to a molecule, provided the tokens in the text respect the SELFIES  
116 vocabulary. For example, the text representation of Benzene in SMILES is C1=CC=CC=C1 while  
117 in SELFIES one possible representation is [C][=C][C][=C][C][=C][Ring1][=Branch1].

118 **Language modeling:** Language modeling often relies on the self-supervised task of next-token  
119 prediction for model pretraining. The general framework for next-token prediction is to train a  
120 model to predict the next token in a sequence autoregressively, i.e. given the previous tokens in the  
121 sequence (left context). Many architectures to handle sequential data have been proposed: Recurrent  
122 Neural Networks (RNNs) [Hochreiter and Schmidhuber, 1997, Rumelhart and McClelland, 1987]  
123 are a class of models used in sequence modeling which use recursive connections in hidden layers  
124 to accumulate the left context for next-token prediction. Transformers are a more recent architecture  
125 that instead use a self-attention mechanism [Vaswani et al., 2017] to capture dependencies between  
126 all tokens in a sequence. For next-token prediction tasks, attention masking is used to enforce left  
127 context, meaning that representations for tokens later in the sequence are only allowed to attend to  
128 previous tokens in the sequence. In Section 4 we outline how we pretrain an autoregressive sequence  
129 model to predict sequences of known molecules.

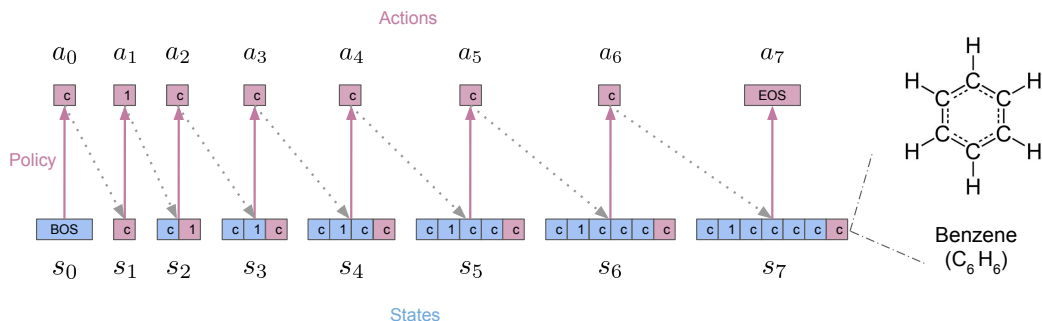


Figure 1: **Autoregressively generating a benzene molecule.**: An autoregressive model for sequence generation can be viewed as an RL policy where the actions  $a_t$  are the next tokens to append to the sequence and the state is the concatenation of all actions taken up to time  $t - 1$ . A special end-of-sequence token can terminate the episode early at time  $T$ . The text at time  $s_T$  is then converted into a molecule based on the text-representation grammar and then scored according to a scoring function that measures the alignment of the molecule with the desired properties informed by the application. Hydrogen atoms are added at the end to complete the structure.

#### 130 4 ChemRLformer Generating molecular strings via reinforcement learning

131 The molecular design space is complex but the benefit from finding improved options is great. In  
 132 this section, we describe ChemRLformer and how combining language models and tools from RL  
 133 produces a sota algorithm.

134 **MDP for molecule generation:** The vocabulary and grammar for text representations of  
 135 molecules can be interpreted as an MDP as described in Section 3, where the states  $s_t$  correspond to  
 136 a *variable length* text of accumulating tokens, and the actions  $a_t$  correspond to vocabulary defined  
 137 by the text-representation. The transition function is a deterministic function where the action  $a_t$   
 138 taken by the agent is appended to the end of the state  $s_t$  resulting in  $s_{t+1}$  using the dynamics  
 139  $s_{t+1} = [s_t, a_t] \leftarrow T(s_t, a_t)$ . However, the corresponding transition function induced in the graph  
 140 representation of molecules is more complex as it is determined by the encoding/decoding rules of  
 141 the chosen text representation. For example, in the SMILES grammar, a random concatenation of  
 142 tokens may not correspond to a valid molecule, while the SELFIES grammar is constructed such  
 143 that any ordering of its tokens is encoded as a valid molecule.

144 Finally, the reward function  $\mathcal{R}$  scores molecules according to their alignment with desired chemical  
 145 properties, which can involve complex material simulations. The underlying property computation  
 146 of the reward function further informs the dynamics of the MDP imposed by text representation.  
 147 For example, docking scores are used to estimate the binding affinity between ligands and protein  
 148 targets. We discuss reward functions for molecules in more detail in section Section 5.

149 **Pretraining policies for molecule discovery.** To advance effectively within this vast search space,  
 150 we make use of datasets containing a large number of drug-like molecules in text format [Irwin et al.,  
 151 2012, Sterling and Irwin, 2015b, Mendez et al., 2019]. This data is used to train an autoregressive  
 152 model to predict tokens that conform to the grammar for drug-like molecules, instead of the random  
 153 texts that are generated from a randomly initialized policy, thereby significantly simplifying the  
 154 exploration problem. In particular, we pretrain a network  $p_\phi$  on the self-supervised objective of  
 155 next-token prediction. Although large language models can be trained with other objectives, such as  
 156 corrupted text reconstruction [Edwards et al., 2022b], these models are not a good fit for our purposes  
 157 since they cannot generate diverse and valid molecules without access to carefully designed prompts.

$$\min_{\theta} \mathbb{E}_{A \sim D} \left[ \sum_{t=1}^H -\log p_{\theta}(a_t = A_t | A_{t-1}, \dots, A_0) \right]. \quad (2)$$

158 In practice a minibatch of sequences  $\{A^1, \dots, A^m\}$  are sampled from the prior dataset  $D$  to  
 159 evaluate the loss function in Equation 2, and the parameters are trained using gradient descent.

## 160 4.1 RL for molecule generation

161 To generate a molecule, a ChemRLformer policy  $\pi_\theta(a_t | s_t)$  is allowed to autoregressively sample  
162 tokens for a fixed number of timesteps  $H$ . The start state  $s_0$  is always a beginning-of-sequence token  
163 [BOS], and the agent can terminate early by taking the end-of-sequence action [EOS]. Figure 1,  
164 shows how an RL policy can construct a Benzene molecule. Since we are only interested in the  
165 properties of the final molecule, there are no intermediate rewards and the goal of the RL policy is to  
166 maximize the expected scalar reward corresponding to the final constructed molecule,  $r(s_T)$ . Thus,  
167 assuming a discount rate  $\gamma = 1$ , Equation 1 can be rewritten more simply as:

$$\max_{\theta} \mathbb{E}_{s_T \sim \pi_\theta} [r(s_T)] \quad (3)$$

where  $s_T = [\text{BOS}][a_0][a_1] \cdots [\text{EOS}]$ , is sampled autoregressively from the policy.

168 Our experiments use the policy gradient algorithm [Sutton et al., 1999b] to train the RL  
169 policy because it is known to achieve state-of-the-art performance amongst RL for molecular  
170 optimization [Olivecrona et al., 2017]. Deep RL policies are able to learn the non-linear global  
171 structures of molecular texts which, as we show in section 5, enables them to generalize to novel  
172 and diverse molecules. However, training RL policies from scratch is time-consuming and can  
173 make the exploration problem infeasibly difficult. Next, we explain how we adapt recent language  
174 modelling techniques to pretrain the RL policy.

175 **RL fine-tuning** The pretrained model can directly be used to sample novel drug-like molecules.  
176 These molecules, however, are not optimized for any particular property. Note that given our  
177 definition of the state  $s_t$  as the concatenated history of all previous actions, this pretrained network  
178 is exactly analogous to the policy network in Equation 1. Hence, by initializing  $\pi_\theta = p_\phi$ , and  
179  $\theta = \phi$ , we can fine-tune this pretrained network by optimizing Equation 1 via the policy gradient  
180 algorithm - REINFORCE [Sutton et al., 1999a]. We need only to define a reward function  $r(s_T)$   
181 which scores molecules according to their alignment with the desired properties. In the following  
182 experiments, we show that this fine-tuning is vital for ChemRLformer to sample better molecules.  
183 We also highlight the importance of pretraining and study how the size and quality of the prior data  
184 affect the downstream ability of RL to search for high-value molecules.

## 185 5 Experimental Results

186 Our proposed algorithm ChemRLformer uses the best combinations of choices resulting from  
187 assessing the performance across three dimensions: (1) what pretraining factors are important to  
188 improve RL for molecular discovery (Section 5.2), (2) how the use of recent text-based molecule  
189 grammars facilitates downstream RL exploration (Section 5.3); and, lastly, (3) which specific  
190 algorithmic changes are necessary to improve RL performance (Section 5.4).

### 191 5.1 Experimental Setup

192 **Tasks.** We evaluate ChemRLformer against five different **docking** targets [Alhossary et al., 2015]  
193 (fa7, parp1, 5ht1b, jak2, and braf) previously explored in the literature [Yang et al., 2021, Lee  
194 et al., 2023]. The docking scores used to estimate the binding affinity between ligands and  
195 protein targets are a complex function of the global molecular structure and have been proposed as  
196 chemically relevant benchmarks for molecule design algorithms [Cieplinski et al., 2021, Tripp et al.,  
197 2022]. In addition to the docking targets, we also evaluate on 22 pharmaceutically-relevant oracle  
198 functions [Huang et al., 2021, Gao et al., 2022, Brown et al., 2019] (**pytdc** tasks), which include  
199 tasks such as optimizing proxies of bioactivity, similarity to target molecules, and combinations of  
200 multiple physiochemical drug properties.

201 **Evaluation metrics.** We design our evaluation procedure with the final goal of identifying the  
202 best candidates to test in a wet lab. To discover such high-value candidate molecules, we use sota  
203 simulators that assign rewards to molecules by performing complex docking simulations [Alhossary  
204 et al., 2015] or using proxy models and chemical rules [Huang et al., 2021]. Previous works limit the  
205 number of molecules sampled during evaluation to around 3000 for **docking** tasks [García-Ortegón  
206 et al., 2022, Yang et al., 2021, Lee et al., 2023] and 10000 for **pytdc** tasks [Gao et al., 2022, Brown  
207 et al., 2019] due to the computational cost associated with these reward simulators. We allow up to

208 25000 unique oracle calls and up to 40000 total oracle calls (allowing repeats). We argue this better  
 209 reflects the lower cost and availability of computing resources relative to wet-lab resources. From  
 210 all the sampled molecules, the average score of the top- $k$  ( $k = 1, 10, 100$ ) molecules is used as a  
 211 performance metric. These top groups are an estimate of the algorithm’s ability to discover a group  
 212 of top-quality candidates that could be given to a wet lab for thorough testing. We report **pytdc**  
 213 scores on a normalized basis between zero and one by default. Next, we normalize all docking  
 214 scores by dividing them by -20 in our experiments. Additionally, we report *diversity*, defined as the  
 215 averaged internal distance of the top 100 molecules, and *redundancy*, defined as the total number of  
 216 oracle calls that an agent makes for an already evaluated molecule.

217 **Pretraining.** We study how the quality and size of prior data affect the downstream RL  
 218 performance of ChemRLformer by pretraining a GPT [Radford et al., 2018] style transformer model  
 219 on five datasets of varying sizes and quality and using the pretrained model as an initialization for the  
 220 RL agent’s policy network. See Table 2 for the name, size, and description of all datasets used in our  
 221 work. We also rank all datasets based on their quality on **docking** and **pytdc** tasks. We determine  
 222 the quality of a dataset by the performance of molecules sampled from the model pretrained on that  
 223 dataset. The quality of ChemRLformer’s pretrained model is evaluated using the *top-100* molecules  
 224 sampled by the pretrained model under the same evaluation setup in Appendix A.2. By default, these  
 225 open-sourced datasets contain a large number of drug-like molecules in SMILES format. For our  
 226 experiments, we also convert all datasets to the SELFIES format. Lastly, three different architectures  
 227 are compared: **fully-connected (FC)**, **recurrent (RNN)** - a GRU and **transformer** - GPT style  
 autoregressive model, and compare them on downstream RL tasks.

Table 2: **Description of molecular datasets used for pretraining:** Datasets are ranked according to procedure described in Section 5.1. Two datasets have the same rank if their average performance lies inside one standard error of the other. The datasets are drawn from a subset of the Zinc [Sterling and Irwin, 2015b, Irwin et al., 2022] and ChEMBL [Gaulton et al., 2012] databases.

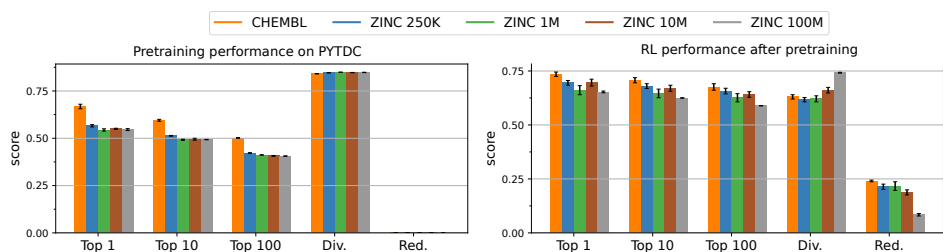
Dataset	Size	Docking Rank	Pytdc Rank	Description
CHEMBL	1.2 M	1	1	Manually curated database of bioactive molecules with drug-like properties [Gaulton et al., 2012].
ZINC 250K	250 K	2	2	ZINC database molecules curated for their pharmaceutical relevance and popularity [Gao et al., 2022].
ZINC 1M	1 M	3	3	Random molecules from $\approx 1.5$ billion
ZINC 10M	10 M	3	4	molecules from the ZINC database [Sterling and Irwin, 2015a].
ZINC 100M	100 M	3	4	ZINC 1M $\subset$ ZINC 10M $\subset$ ZINC100M.

228  
 229 All of our experiments on **pytdc** tasks are run across 5 random seeds. Since docking simulations are  
 230 expensive and time consuming, we run all **docking** experiments across 3 random seeds. Experiments  
 231 with different seeds use the same pretrained model which is only pretrained once for every dataset.  
 232 Additional details about the task rewards, evaluation metrics, and the pretraining datasets and models  
 233 are discussed in Appendix A.1, A.2, and A.3 respectively.

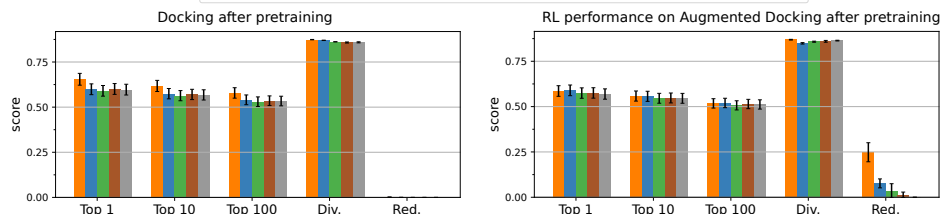
## 234 5.2 How does prior data affect the final performance of ChemRLformer?

235 In this section, we pretrain the REINFORCE policy on datasets of varying size and quality  
 236 from Table 2. Our datasets vary from small (250K) to very large (100M) sizes. Due to the  
 237 parallelizability of training on larger datasets, we use the transformer policy architecture for all  
 238 experiments in this section. In natural language processing (NLP), pretraining transformer models  
 239 on large and diverse unlabelled datasets have been found to perform well on downstream tasks  
 240 using few-shot labeled data [Brown et al., 2020b]. Yet, our results in Figure 2b indicate that the  
 241 quality of the prior dataset matters more than its size. Figure 2, shows that the distribution of the  
 242 *ChEMBL* dataset is more aligned with both the **pytdc** and the **docking** tasks. As a result, the RL  
 243 agent pretrained on the *ChEMBL* dataset outperforms all other agents, including the ones trained on  
 244 100 times more data.

245 Results may seem surprising from an NLP perspective, but they make sense when viewed from an  
 246 RL perspective. Pretraining using next token prediction Equation (2), is analogous to behavior  
 247 cloning in this context, where the performance depends largely on the quality of the offline  
 248 dataset [Ross et al., 2011, Ho and Ermon, 2016]. These results suggest that ChemRLformer might



(a) Performance on SMILES-based molecular design with perturbing (left) and with pretraining and RL (right).



(b) performance on SMILES-based molecular docking with perturbing (left) and with pretraining and RL (right). Section 5.2 describes augmented docking setting with additional experiments shown in Appendix A.3.

Figure 2: On the left pretrained performance on SMILES-based ChemRLformer. Higher-quality datasets, such as ChemBL lead to higher-performance for both **pytdc** and augmented docking. On the right is the performance after RL training. RL has a substantial benefit for **pytdc** tasks, while for docking tasks an *augmented docking* score is used to avoid reward hacking, see Figure 4 for details.

249 benefit from better pretraining objectives, that go beyond simple imitation learning, when trained on  
 250 large and diverse offline datasets [Kumar et al., 2023, Farebrother et al., 2023].

### 251 5.3 Text representations and architectures for ChemRLformer

252 Starting with a REINFORCE agent, we isolate the effect of various text representations for  
 253 molecules and policy network architectures on performance. All experiments in this section use  
 254 **ZINC-250k** dataset for pretraining. Similar results obtained for other datasets are shared in the  
 255 following sections. Whenever we show normalized results across different experiments, we add the  
 256 individual plots in Appendix C.1.

257 **Text representations.** In Figure 3 we compare ChemRLformer agents using different  
 258 architectures and tasks across environments that base their dynamics on SELFIES and SMILES.  
 259 The results show normalized scores across all architectures. Consistent with prior work [Gao et al.,  
 260 2022] we find that SMILES-based policies generally outperforms SELFIES-based policies. On all  
 261 **pytdc** tasks and architectures, ChemRLformer agents based on SMILES consistently achieve better  
 262 rewards when compared to SELFIES-based agents across all reward metrics. Although more subtle,  
 263 we observe a similar theme in the **docking** tasks where SMILES achieves higher rewards than  
 264 SELFIES on all top-K metrics. Another consistent theme in the results is that even though the  
 265 diversity of top-100 molecules obtained by SELFIES is higher, the redundancy of SELFIES agents  
 266 is higher as well. This means that SELFIES-based ChemRLformer agents explore a much smaller  
 267 region of the molecular space. These results suggest that the rules which allow SELFIES strings to  
 268 always be converted into a valid molecule can actually be detrimental to the agent’s exploration and  
 269 search for high-value molecules, more details in Appendix C.1.

270 **Architectures.** The results in Figure 4 show that the **transformer** and **RNN** have similar  
 271 performance on all tasks. On the **pytdc** tasks, **FC** achieves worse performance than other  
 272 architectures specially made to handle strings, as expected. However, on **docking** tasks, **FC**  
 273 obtains unusually high rewards. We find that this method performs a type of *reward function*  
 274 *hacking* [Amodei et al., 2016, Skalse et al., 2022, Everitt, 2019] by exploiting a corner case of  
 275 the docking-based reward function which provides high rewards for long strings of Carbon and  
 276 Nitrogen atoms together. To evade the reward hacking of docking scores, we constructed an  
 277 *augmented docking* score function with commonly used oracles (QED and SA scores) based on

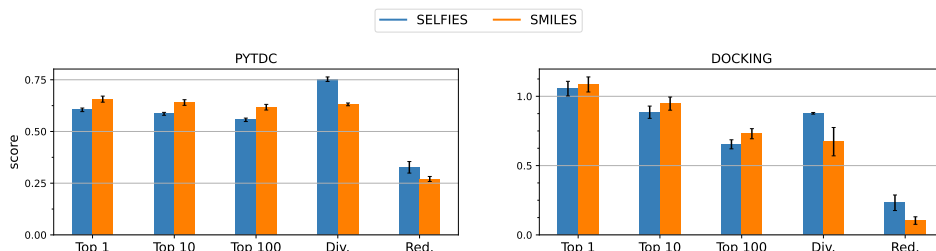


Figure 3: **Comparison between SELFIES and SMILES:** The SELFIES representation makes it relatively difficult for ChemRLformer agents to explore effectively leading to generally lower performance on pytdc and docking while scoring higher on diversity. Scores are reported for the transformer model and are averaged across all reward functions.

278 previous work [Lee et al., 2023] (See Appendix A.3 for more details). This finding shows that the  
 279 REINFORCE agent can search the space well and, in this case, can be used to expose issues with  
 280 the current design of reward functions.

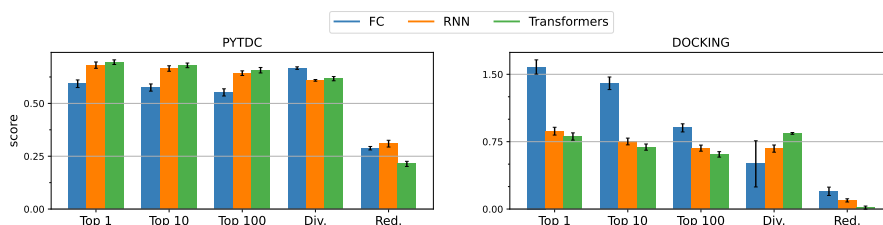


Figure 4: **Comparison of different policy architectures:** No single architecture clearly outperforms for molecular ChemRLformer. Although FC does better on the docking tasks, our analysis shows that it learns to exploit the docking function as opposed to designing high-value molecules. More details about ways to tackle this issue are given in Appendix C.3. Additional experiments for comparing transformers and RNNs are shown in Appendix C.5. These experiments use the smiles text representation.

## 281 5.4 Revisiting RL Algorithm Design Choices for ChemRLformer.

282 Previous experiments identify key design and algorithmic choices for efficient text-based RL. But  
 283 recent text-based RL algorithms [Olivecrona et al., 2017, Bjerrum et al., 2023, Thomas et al.,  
 284 2022] also employ many additional algorithmic components like replay buffers, hill climbing, KL  
 285 regularisation towards the pretrained policy, and likelihood penalties. We perform an ablation  
 286 study across these components to understand which ones are beneficial for the performance of  
 287 ChemRLformer. All experiments in this section are performed on **pytdc** tasks, using an RNN  
 288 architecture as it is most commonly used in text-based RL. See Appendix C.4 for experiments with  
 289 other architectures and reward functions.

290 **Replay buffers and hill climbing.** In  
 291 off-policy deep RL, a replay buffer is  
 292 generally used to store and reuse previous  
 293 trajectories for training. Although text-based  
 294 RL algorithms are trained on-policy, prior  
 295 work has proposed using a replay buffer to  
 296 improve performance [Mnih et al., 2013].  
 297 Standard replay buffers throw away the oldest  
 298 trajectories as newer ones arrive. But many  
 299 text-based RL algorithms propose to use hill-  
 300 climb replay buffers, that randomly sample a batch of molecules from the highest scoring molecules  
 301 seen so far and add them to the current mini batch. In Figure 5, we see that using the hill-climb  
 302 buffer results in a significant performance boost for ChemRLformer, whereas using a standard  
 303 buffer does not contribute much. Notably, the use of a hill-climb replay buffer reduces diversity

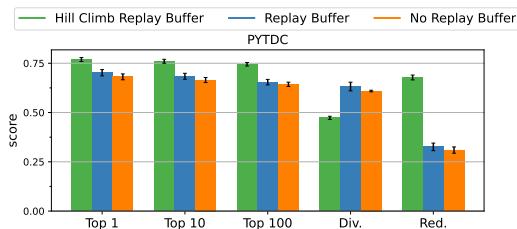


Figure 5: Hill climbing buffer lead to 13% improvement in Top-100 rewards.



304 and increases redundancy quite substantially. The following two experiments involve combining  
305 regularisation terms with the RL objective in Equation 3. The coefficients on these extra terms can  
306 largely affect the final performance. To make a fair comparison, we perform hyper-parameter tuning  
307 over six different values for every new regularisation term with details provided in Appendix B.

### 308 **Should the policy to stay close to the pretrained model?**

309 Pretrained models carry information on how to build valid drug-like molecules. To ensure that ChemRLformer  
310 agents do not stray far away from the space of valid drug-like molecules during exploration,  
311 Olivecrona et al. [2017], Gao et al. [2022] constrain the KL divergence between the policy  
312 and the pretrained model by adding a KL penalty to the policy gradient loss function in  
313 Equation (3). Prior works show that adding this penalty helps the agent achieve better sample  
314 efficiency [Gao et al., 2022]. Yet, our results in Figure 6 suggest that, when you increase the  
315 number of oracle calls in simulation, adding this penalty does not yield any additional benefit while  
316 substantially increasing the GPU memory requirement, especially when using larger models. Since  
317 invalid molecules correspond to zero rewards, the ChemRLformer agent is able to learn to avoid  
318 invalid structures on its own merit.  
319  
320  
321  
322  
323

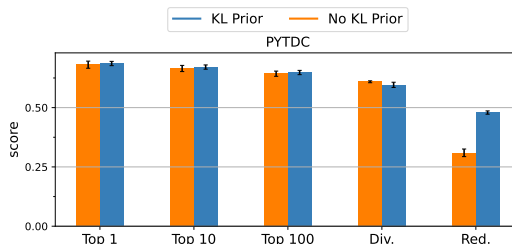


Figure 6: Our experiments show little difference in performance for multiple KL regularization terms.

### 324 **Regularizing the policy’s likelihood for exploration.**

325 RL agents classically face an exploration-exploitation dilemma, which can lead to agents getting stuck in  
326 sub-optimal local maxima when not well balanced. ChemRLformer agents are not immune to this dilemma.  
327 Upon encountering good, but sub-optimal molecules, an agent may adjust its policy to increase the likelihood  
328 of sampling these sub-optimal molecules and, without sufficient exploration, fail to discover higher-value  
329 regions of policy space. This can be particularly detrimental during the initial learning stages.  
330  
331  
332  
333  
334

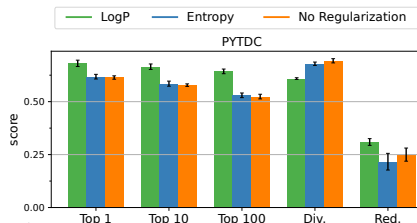


Figure 7: Different likelihood penalization for exploration. Log P regularization is a better choice for efficient exploration for ChemRLformer.

335 To combat this issue, entropy regularisation, which adds a  $\log \pi(s)$  term to the RL loss, has been  
336 proposed [Haarnoja et al., 2018]. This encourages the RL policy to explore states with lower  
337 likelihood values. Similarly [Olivecrona et al., 2017] adds a Log p regularizer, which penalizes  
338 higher likelihood values by adding a  $-1/\log \pi(s)$  term to the RL loss. In Figure 7, our results  
339 show that although an entropy regularizer leads to lesser redundancy, the Log p regularizer boosts  
340 performance significantly by exploring more efficiently. The Log p regularizer only penalizes the  
341 agent for being extremely certain (likelihood  $\xrightarrow{\text{tends to}} 1$ ) about its actions, and is mostly agnostic for  
342 lower likelihood values. This penalty is a much better choice for ChemRLformer as it only activates  
343 when stuck in a local optimum of molecular space.

## 344 **6 Conclusion and Future Work**

345 We present ChemRLformer that resulted from our empirical study of multiple algorithmic  
346 components of text-based molecular design. For future practitioners, our method suggests the  
347 following philosophy: (1) Using SMILES is a better choice than SELFIES. (2) When collecting  
348 data for pretraining, the quality of molecules matter much more than the number of molecules.  
349 (3) Both transformer and RNN architectures achieve similar performance across all tasks using  
350 current datasets. (4) Incorporating components such as a hill-climb buffer and Log P regularization  
351 yields substantial performance improvements. Conversely, introducing KL regularization or opting  
352 for more intricate actor-critic algorithms may result in diminished performance, at the cost of more  
353 hyperparameters and memory resources. While our analysis addresses many questions, it also shows  
354 that RL agents were able to *hack* the reward functions suggesting that there is space to improve on  
355 the metrics used for molecule quality.

## 356 References

- 357 Amr Alhossary, Stephanus Daniel Handoko, Yuguang Mu, and Chee-Keong Kwoh. Fast, accurate,  
358 and reliable molecular docking with QuickVina 2. *Bioinformatics*, 31(13):2214–2216, 02 2015.  
359 ISSN 1367-4803. doi: 10.1093/bioinformatics/btv082. URL [https://doi.org/10.1093/  
360 bioinformatics/btv082](https://doi.org/10.1093/bioinformatics/btv082).
- 361 Zaccary Alperstein, Artem Cherkasov, and Jason Tyler Rolfe. All smiles variational autoencoder.  
362 *arXiv preprint arXiv:1905.13343*, 2019.
- 363 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané.  
364 Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- 365 Emmanuel Bengio, Moksh Jain, Maksym Korablyov, Doina Precup, and Yoshua Bengio. Flow  
366 network based generative models for non-iterative diverse candidate generation, 2021. URL  
367 <https://arxiv.org/abs/2106.04399>.
- 368 Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Janvin. A neural probabilistic  
369 language model. *J. Mach. Learn. Res.*, 3(null):1137–1155, mar 2003. ISSN 1532-4435.
- 370 Mostapha Benhenda. Chemgan challenge for drug discovery: can ai reproduce natural chemical  
371 diversity?, 2017.
- 372 Esben Jannik Bjerrum and Richard Threlfall. Molecular generation with recurrent neural networks  
373 (rnns). *arXiv preprint arXiv:1705.04612*, 2017.
- 374 Esben Jannik Bjerrum, Christian Margreitter, Thomas Blaschke, and Raquel López-Ríos de Castro.  
375 Faster and more diverse de novo molecular optimization with double-loop reinforcement learning  
376 using augmented smiles, 2023.
- 377 Thomas Blaschke, Josep Arús-Pous, Hongming Chen, Christian Margreitter, Christian Tyrchan, Ola  
378 Engkvist, Kostas Papadopoulos, and Atanas Patronov. Reinvent 2.0: an ai tool for de novo drug  
379 design. *Journal of chemical information and modeling*, 60(12):5918–5922, 2020a.
- 380 Thomas Blaschke, Josep Arús-Pous, Hongming Chen, Christian Margreitter, Christian Tyrchan, Ola  
381 Engkvist, Kostas Papadopoulos, and Atanas Patronov. Reinvent 2.0: An ai tool for de novo drug  
382 design. *Journal of Chemical Information and Modeling*, 60, 10 2020b. doi: 10.1021/acs.jcim.  
383 0c00915.
- 384 Nathan Brown, Marco Fiscato, Marwin H.S. Segler, and Alain C. Vaucher. GuacaMol:  
385 Benchmarking models for de novo molecular design. *Journal of Chemical Information and  
386 Modeling*, 59(3):1096–1108, mar 2019. doi: 10.1021/acs.jcim.8b00839. URL [https://doi.  
387 org/10.1021%2Facs.jcim.8b00839](https://doi.org/10.1021%2Facs.jcim.8b00839).
- 388 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
389 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel  
390 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel  
391 Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz  
392 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec  
393 Radford, Ilya Sutskever, and Dario Amodei. Language Models are Few-Shot Learners.  
394 In H Larochelle, M Ranzato, R Hadsell, M F Balcan, and H Lin, editors, *Advances in  
395 Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates,  
396 Inc., 2020a. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/file/  
397 1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf).
- 398 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,  
399 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel  
400 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.  
401 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin,  
402 Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford,  
403 Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020b.

- 404 Tianjian Chen, Zhanpeng He, and Matei Ciocarlie. Hardware as policy: Mechanical and  
405 computational Co-Optimization using deep reinforcement learning. In Jens Kober, Fabio Ramos,  
406 and Claire Tomlin, editors, *Proceedings of the 2020 Conference on Robot Learning*, volume 155  
407 of *Proceedings of Machine Learning Research*, pages 1158–1173. PMLR, 2021a.
- 408 Tianlong Chen, Kaixiong Zhou, Keyu Duan, Wenqing Zheng, Peihao Wang, Xia Hu, and Zhangyang  
409 Wang. Bag of tricks for training deeper graph neural networks: A comprehensive benchmark  
410 study. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- 411 Xiaocong Chen, Lina Yao, Julian McAuley, Guanglin Zhou, and Xianzhi Wang. A survey of deep  
412 reinforcement learning in recommender systems: A systematic review and future directions. *arXiv*  
413 *preprint arXiv:2109.03540*, 2021b.
- 414 Austin H Cheng, Andy Cai, Santiago Miret, Gustavo Malkomes, Mariano Phielipp, and Alán  
415 Aspuru-Guzik. Group selfies: a robust fragment-based molecular string representation. *Digital*  
416 *Discovery*, 2023.
- 417 Tobiasz Cieplinski, Tomasz Danel, Sabina Podlowska, and Stanislaw Jastrzebski. We should at least  
418 be able to design molecules that dock well, 2021.
- 419 Vijay Prakash Dwivedi, Chaitanya K. Joshi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio, and  
420 Xavier Bresson. Benchmarking graph neural networks, 2022.
- 421 Carl Edwards, Tuan Lai, Kevin Ros, Garrett Honke, Kyunghyun Cho, and Heng Ji. Translation  
422 between molecules and natural language. In *Proceedings of the 2022 Conference on Empirical*  
423 *Methods in Natural Language Processing*, pages 375–413, Abu Dhabi, United Arab Emirates,  
424 December 2022a. Association for Computational Linguistics. URL <https://aclanthology.org/2022.emnlp-main.26>.
- 425
- 426 Carl Edwards, Tuan Lai, Kevin Ros, Garrett Honke, Kyunghyun Cho, and Heng Ji. Translation  
427 between molecules and natural language, 2022b.
- 428 Tom Everitt. *Towards safe artificial general intelligence*. PhD thesis, The Australian National  
429 University (Australia), 2019.
- 430 Benedek Fabian, Thomas Edlich, H el ena Gaspar, Marwin Segler, Joshua Meyers, Marco Fiscato,  
431 and Mohamed Ahmed. Molecular representation learning with language models and domain-  
432 relevant auxiliary tasks. *arXiv preprint arXiv:2011.13230*, 2020.
- 433 Jesse Farebrother, Joshua Greaves, Rishabh Agarwal, Charline Le Lan, Ross Goroshin,  
434 Pablo Samuel Castro, and Marc G. Bellemare. Proto-value networks: Scaling representation  
435 learning with auxiliary tasks, 2023.
- 436 Daniel Flam-Shepherd, Kevin Zhu, and Al an Aspuru-Guzik. Language models can learn complex  
437 molecular distributions. *Nature Communications*, 13(1):3293, 2022.
- 438 Tianfan Fu, Cao Xiao, Lucas M. Glass, and Jimeng Sun. Moler: Incorporate molecule-level reward  
439 to enhance deep generative model for molecule optimization. *IEEE Transactions on Knowledge*  
440 *and Data Engineering*, 34(11):5459–5471, 2022. doi: 10.1109/TKDE.2021.3052150.
- 441 Wenhao Gao, Tianfan Fu, Jimeng Sun, and Connor W. Coley. Sample efficiency matters: A  
442 benchmark for practical molecular optimization, 2022.
- 443 Miguel Garc a-Orteg on, Gregor NC Simm, Austin J Tripp, Jos e Miguel Hern andez-Lobato, Andreas  
444 Bender, and Sergio Bacallado. Dockstring: easy molecular docking yields better benchmarks for  
445 ligand design. *Journal of chemical information and modeling*, 62(15):3486–3502, 2022.
- 446 Anna Gaulton, Louisa J. Bellis, A. Patricia Bento, Jon Chambers, Mark Davies, Anne Hersey,  
447 Yvonne Light, Shaun McGlinchey, David Michalovich, Bissan Al-Lazikani, and John P.  
448 Overington. ChEMBL: A large-scale bioactivity database for drug discovery. *Nucleic Acids*  
449 *Research*, 40(D1):D1100–D1107, January 2012. ISSN 0305-1048. doi: 10.1093/nar/gkr777.
- 450 Simon Geisler, Tobias Schmidt, Hakan  irin, Daniel Z ugner, Aleksandar Bojchevski, and Stephan  
451 G unnemann. Robustness of graph neural networks at scale, 2023.

- 452 Rafael Gómez-Bombarelli, Jennifer N. Wei, David Duvenaud, José Miguel Hernández-Lobato,  
453 Benjamín Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D. Hirzel,  
454 Ryan P. Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven  
455 continuous representation of molecules. *ACS Central Science*, 4(2):268–276, jan 2018. doi:  
456 10.1021/acscentsci.7b00572. URL <https://doi.org/10.1021/acscentsci.7b00572>.
- 457 Rafael Gómez-Bombarelli, Jennifer N Wei, David Duvenaud, José Miguel Hernández-Lobato,  
458 Benjamín Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D Hirzel,  
459 Ryan P Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven  
460 continuous representation of molecules. *ACS central science*, 4(2):268–276, 2018.
- 461 Francesca Grisoni, Michael Moret, Robin Lingwood, and Gisbert Schneider. Bidirectional molecule  
462 generation with recurrent neural networks. *Journal of chemical information and modeling*, 60(3):  
463 1175–1183, 2020.
- 464 Anvita Gupta, Alex T Müller, Berend JH Huisman, Jens A Fuchs, Petra Schneider, and Gisbert  
465 Schneider. Generative recurrent networks for de novo drug design. *Molecular informatics*, 37  
466 (1-2):1700111, 2018.
- 467 David Ha. Reinforcement learning for improving agent design. *Artif. Life*, 25(4):352–365,  
468 November 2019.
- 469 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy  
470 maximum entropy deep reinforcement learning with a stochastic actor, 2018.
- 471 Stephen Heller, Alan McNaught, Stephen Stein, Dmitrii Tchekhovskoi, and Igor Pletnev. Inchi-the  
472 worldwide chemical structure identifier standard. *Journal of cheminformatics*, 5(1):1–9, 2013.
- 473 Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning, 2016.
- 474 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):  
475 1735–1780, nov 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.8.1735. URL <https://doi.org/10.1162/neco.1997.9.8.1735>.
- 477 Kexin Huang, Tianfan Fu, Wenhao Gao, Yue Zhao, Yusuf Roohani, Jure Leskovec, Connor W.  
478 Coley, Cao Xiao, Jimeng Sun, and Marinka Zitnik. Therapeutics data commons: Machine  
479 learning datasets and tasks for drug discovery and development, 2021.
- 480 Kexin Huang, Tianfan Fu, Wenhao Gao, Yue Zhao, Yusuf Roohani, Jure Leskovec, Connor W Coley,  
481 Cao Xiao, Jimeng Sun, and Marinka Zitnik. Artificial intelligence foundation for therapeutic  
482 science. *Nature Chemical Biology*, 18(10):1033–1036, 2022.
- 483 John J Irwin, Teague Sterling, Michael M Mysinger, Erin S Bolstad, and Ryan G Coleman. Zinc:  
484 a free tool to discover chemistry for biology. *Journal of chemical information and modeling*, 52  
485 (7):1757–1768, 2012.
- 486 Ross Irwin, Spyridon Dimitriadis, Jiazhen He, and Esben Jannik Bjerrum. Chemformer: a pre-  
487 trained transformer for computational chemistry. *Machine Learning: Science and Technology*, 3  
488 (1):015022, 2022.
- 489 Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Junction tree variational autoencoder for  
490 molecular graph generation, 2018. URL <https://arxiv.org/abs/1802.04364>.
- 491 Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Multi-objective molecule generation using  
492 interpretable substructures, 2020.
- 493 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,  
494 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language  
495 models, 2020.
- 496 Shauharda Khadka, Estelle Aflalo, Mattias Marder, Avrech Ben-David, Santiago Miret, Shie  
497 Mannor, Tamir Hazan, Hanlin Tang, and Somdeb Majumdar. Optimizing memory placement  
498 using evolutionary graph reinforcement learning. *arXiv preprint arXiv:2007.07298*, 2020.

499 Mario Krenn, Florian Häse, AkshatKumar Nigam, Pascal Friederich, and Alan Aspuru-Guzik.  
500 Self-referencing embedded strings (SELFIES): A 100% robust molecular string representation.  
501 *Machine Learning: Science and Technology*, 1(4):045024, oct 2020. doi: 10.1088/2632-2153/  
502 aba947. URL <https://doi.org/10.1088/2632-2153/2Faba947>.

503 Mario Krenn, Qianxiang Ai, Senja Barthel, Nessa Carson, Angelo Frei, Nathan C Frey, Pascal  
504 Friederich, Théophile Gaudin, Alberto Alexander Gayle, Kevin Maik Jablonka, et al. Selfies and  
505 the future of molecular string representations. *Patterns*, 3(10):100588, 2022.

506 Aviral Kumar, Amir Yazdanbakhsh, Milad Hashemi, Kevin Swersky, and Sergey Levine. Data-  
507 driven offline optimization for architecting hardware accelerators, 2022.

508 Aviral Kumar, Rishabh Agarwal, Xinyang Geng, George Tucker, and Sergey Levine. Offline q-  
509 learning on diverse multi-task data both scales and generalizes, 2023.

510 Greg Landrum et al. Rdkit: A software suite for cheminformatics, computational chemistry, and  
511 predictive modeling. *Greg Landrum*, 8, 2013.

512 Seul Lee, Jaehyeong Jo, and Sung Ju Hwang. Exploring chemical space with score-based out-of-  
513 distribution generation, 2023.

514 Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts, 2017.

515 Kevin Sebastian Luck, Heni Ben Amor, and Roberto Calandra. Data-efficient Co-Adaptation of  
516 morphology and behaviour with deep reinforcement learning. In Leslie Pack Kaelbling, Danica  
517 Kracic, and Komei Sugiura, editors, *Proceedings of the Conference on Robot Learning*, volume  
518 100 of *Proceedings of Machine Learning Research*, pages 854–869. PMLR, 2020.

519 David Mendez, Anna Gaulton, A Patrícia Bento, Jon Chambers, Marleen De Veij, Eloy Félix,  
520 María Paula Magariños, Juan F Mosquera, Prudence Mutowo, Michał Nowotka, et al. ChEMBL:  
521 towards direct deposition of bioassay data. *Nucleic acids research*, 47(D1):D930–D940, 2019.

522 Santiago Miret, Vui Seng Chua, Mattias Marder, Mariano Phiellip, Nilesh Jain, and  
523 Somdeb Majumdar. Neuroevolution-enhanced multi-objective optimization for mixed-precision  
524 quantization. In *Proceedings of the Genetic and Evolutionary Computation Conference*, pages  
525 1057–1065, 2022.

526 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan  
527 Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning, 2013.

528 Henry B. Moss, Daniel Beck, Javier Gonzalez, David S. Leslie, and Paul Rayson. Boss: Bayesian  
529 optimization over string spaces, 2020.

530 AkshatKumar Nigam, Robert Pollice, and Alan Aspuru-Guzik. Janus: Parallel tempered genetic  
531 algorithm guided by deep neural networks for inverse molecular design, 2021.

532 Marcus Olivecrona, Thomas Blaschke, Ola Engkvist, and Hongming Chen. Molecular de novo  
533 design through deep reinforcement learning, 2017.

534 Emilio Parisotto, H. Francis Song, Jack W. Rae, Razvan Pascanu, Caglar Gulcehre, Siddhant M.  
535 Jayakumar, Max Jaderberg, Raphael Lopez Kaufman, Aidan Clark, Seb Noury, Matthew M.  
536 Botvinick, Nicolas Heess, and Raia Hadsell. Stabilizing transformers for reinforcement learning,  
537 2019.

538 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor  
539 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward  
540 Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner,  
541 Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance  
542 deep learning library, 2019.

543 Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language  
544 understanding by generative pre-training. 2018.

545 Stephane Ross, Geoffrey J. Gordon, and J. Andrew Bagnell. A reduction of imitation learning and  
546 structured prediction to no-regret online learning, 2011.

- 547 David E. Rumelhart and James L. McClelland. *Learning Internal Representations by Error*  
548 *Propagation*, pages 318–362. MIT Press, 1987.
- 549 Charles Schaff, David Yunis, Ayan Chakrabarti, and Matthew R Walter. Jointly learning to construct  
550 and control agents using deep reinforcement learning. In *2019 International Conference on*  
551 *Robotics and Automation (ICRA)*, pages 9798–9805. [ieeexplore.ieee.org](https://ieeexplore.ieee.org), May 2019.
- 552 Robin M Schmidt. Recurrent neural networks (rnns): A gentle introduction and overview. *arXiv*  
553 *preprint arXiv:1912.05911*, 2019.
- 554 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
555 optimization algorithms, 2017.
- 556 Cynthia Shen, Mario Krenn, Sagi Eppel, and Alán Aspuru-Guzik. Deep molecular dreaming:  
557 inverse machine learning for de-novo molecular design and interpretability with surjective  
558 representations. *Machine Learning: Science and Technology*, 2(3):03LT02, jul 2021. doi:  
559 10.1088/2632-2153/ac09d6. URL <https://doi.org/10.1088/2632-2153/ac09d6>.
- 560 Joar Skalse, Nikolaus HR Howe, Dmitrii Krashennnikov, and David Krueger. Defining and  
561 characterizing reward hacking. *arXiv preprint arXiv:2209.13085*, 2022.
- 562 Teague Sterling and John Irwin. Zinc 15 - ligand discovery for everyone. *Journal of chemical*  
563 *information and modeling*, 55, 10 2015a. doi: 10.1021/acs.jcim.5b00559.
- 564 Teague Sterling and John J Irwin. Zinc 15–ligand discovery for everyone. *Journal of chemical*  
565 *information and modeling*, 55(11):2324–2337, 2015b.
- 566 Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient  
567 methods for reinforcement learning with function approximation. In S. Solla, T. Leen, and  
568 K. Müller, editors, *Advances in Neural Information Processing Systems*, volume 12. MIT  
569 Press, 1999a. URL [https://proceedings.neurips.cc/paper\\_files/paper/1999/file/](https://proceedings.neurips.cc/paper_files/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf)  
570 [464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf).
- 571 Richard S. Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods  
572 for reinforcement learning with function approximation. In *Proceedings of the 12th International*  
573 *Conference on Neural Information Processing Systems*, NIPS’99, page 1057–1063, Cambridge,  
574 MA, USA, 1999b. MIT Press.
- 575 Qiaoyu Tan, Ninghao Liu, and Xia Hu. Deep representation learning for social network analysis.  
576 *Frontiers in big Data*, 2:2, 2019.
- 577 Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia,  
578 Andrew Poulton, Viktor Kerkez, and Robert Stojnic. Galactica: A large language model for  
579 science. *arXiv preprint arXiv:2211.09085*, 2022.
- 580 Morgan Thomas, Noel O’Boyle, Andreas Bender, and Chris Graaf. Augmented hill-climb increases  
581 reinforcement learning efficiency for language-based de novo molecule generation. *Journal of*  
582 *Cheminformatics*, 14, 10 2022. doi: 10.1186/s13321-022-00646-z.
- 583 Austin Tripp, Wenlin Chen, and José Miguel Hernández-Lobato. An evaluation framework for the  
584 objective functions of de novo drug design benchmarks. In *ICLR2022 Machine Learning for Drug*  
585 *Discovery*, 2022. URL <https://openreview.net/forum?id=W1tcNQNG1S>.
- 586 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
587 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*  
588 *processing systems*, 30, 2017.
- 589 Qi Wang and Chunlei Tang. Deep reinforcement learning for transportation network combinatorial  
590 optimization: A survey. *Knowledge-Based Systems*, 233:107526, 2021.
- 591 Sheng Wang, Yuzhi Guo, Yuhong Wang, Hongmao Sun, and Junzhou Huang. Smiles-bert: large  
592 scale unsupervised pre-training for molecular property prediction. In *Proceedings of the 10th*  
593 *ACM international conference on bioinformatics, computational biology and health informatics*,  
594 pages 429–436, 2019.

- 595 David Weininger. Smiles, a chemical language and information system. 1. introduction to  
596 methodology and encoding rules. *J. Chem. Inf. Comput. Sci.*, 28(1):31–36, feb 1988. ISSN  
597 0095-2338. doi: 10.1021/ci00057a005. URL <https://doi.org/10.1021/ci00057a005>.
- 598 Soojung Yang, Doyeong Hwang, Seul Lee, Seongok Ryu, and Sung Ju Hwang. Hit and lead  
599 discovery with explorative rl and fragment-based molecule generation, 2021.
- 600 Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Mastering visual continuous  
601 control: Improved data-augmented reinforcement learning, 2021.
- 602 Jiaxuan You, Bowen Liu, Rex Ying, Vijay Pande, and Jure Leskovec. Graph convolutional policy  
603 network for goal-directed molecular graph generation, 2018. URL [https://arxiv.org/abs/  
604 1806.02473](https://arxiv.org/abs/1806.02473).
- 605 Wenhao Yu, C. Karen Liu, and Greg Turk. Policy transfer with strategy optimization, 2018.
- 606 Zhenni Zeng, Yuan Yao, Zhiyuan Liu, and Maosong Sun. A deep-learning system bridging molecule  
607 structure and biomedical text with comprehension comparable to human professionals. *Nature  
608 communications*, 13(862), 2022.
- 609 Zhenpeng Zhou, Steven Kearnes, Li Li, Richard N. Zare, and Patrick Riley. Optimization of  
610 molecules via deep reinforcement learning. *Scientific Reports*, 9(1), jul 2019. doi: 10.1038/  
611 s41598-019-47148-x. URL <https://doi.org/10.1038/s41598-019-47148-x>.

612 **Outline of Appendices.** In Appendix A we provide details about the experimental setup.  
613 In Appendix B we describe our hyperparameter tuning strategy. In Appendix C we include  
614 additional results from our experiments.

## 615 A Experimental setup

616 In this section, we provide additional details about the tasks, evaluation metrics and pretraining  
617 models and data used in our work.

### 618 A.1 Tasks

619 **Pydtc tasks.** These tasks are a set of 21 pharmaceutically-relevant oracle functions, which have  
620 been commonly used in prior work [Brown et al., 2019, Gao et al., 2022, Huang et al., 2021] for  
621 evaluating performance across molecular discovery algorithms:

- 622 • QED: A quantitative estimate of drug-likeness calculated using a set of rules.
- 623 • DRD2, GSK3 $\beta$ , and JNK3: Classical machine learning models (SVMs and random forests)  
624 that provide an estimate of properties like target affinity or susceptibility towards a disorder.
- 625 • Celecoxib, Troglitazone, and Thiothixene rediscovery: An estimate of smiles text  
626 similarity, based on tanimoto metric, towards a target molecule.
- 627 • Albuterol and Mestranol similarity: Generate molecules similar to a target molecule.
- 628 • Isomers\_c7h8n2o2 and isomers\_c9h10n2o2pf2cl: Generate molecules corresponding to a  
629 target molecular formula.
- 630 • Median1 and Median2: Generate molecules that are maximally similar to several target  
631 molecules.
- 632 • Osimertinib\_mpo, fexofenadine\_mpo, ranolazine\_mpo, perindopril\_mpo, amlodipine\_mpo,  
633 sitagliptin\_mpo, zaleplon\_mpo: Generate molecules that maximize multiple properties of a  
634 targeted drug.
- 635 • valsartan\_smarts: Generate molecules that contain a certain SMARTS pattern and certain  
636 physicochemical properties.

637 Most of these tasks are from the GuacaMol benchmark [Brown et al., 2019]. All oracles are  
638 calculated using the Python API provided by Therapeutics Data Commons [Huang et al., 2021]  
639 and more details for these tasks can be found on their website.

640 **Docking tasks.** We used QuickVina 2 [Alhossary et al., 2015] for calculating docking scores using  
641 the same default configuration parameters as prior works [Yang et al., 2021, Lee et al., 2023]. For  
642 example, we used exhaustiveness = 1, and modes = 10. We choose 5 different protein targets  
643 to calculate docking scores: fa7 (FA7), parp1 (PARP-1), 5ht1b (5-HT1B), jak2 (JAK-2), and braf  
644 (BRAF). These targets were chosen by [Yang et al., 2021, Lee et al., 2023] because the docking  
645 simulators for these targets work fairly well when compared to the ground truth. In our experiments  
646 in Section 5.3, we found that text-based RL algorithms were easily able to produce chemically  
647 trivial molecules that have very high docking scores. To understand the complexity of computing  
648 docking scores, we report the time taken to dock 1000 molecules in parallel using 12 CPUs table 3.  
649 We also provide the time taken to run the RL algorithm on these 1000 molecules after their docking  
650 scores are available.

Table 3: **Time complexity of docking score evaluation** : More than half the running time is spent evaluating the docking scores.

Number of molecules	Docking time	RL update time
1000	130 seconds	74 seconds



651 **Augmented docking tasks.** In our results for the standard **docking** tasks (Figure 3 and Figure 4),  
652 we found that using the simulated docking scores as rewards did not lead to chemically relevant  
653 molecules. Text-based RL algorithms were able to exploit their state and action spaces to design  
654 chemically trivial molecules that have very high docking scores. To tackle this issue of undesirable  
655 reward hacking, we tried a reward function based on prior works [García-Ortegón et al., 2022, Lee  
656 et al., 2023] that combine objectives for drug-like, and synthesizable molecules with docking scores.  
657 We call tasks corresponding to this new reward function as **augmented docking** tasks. Concretely,  
658 we chose the same reward function from [Lee et al., 2023]

$$r(s) = -DS(s)/20 \times QED(s) \times (10 - SA(s))/9, \quad (4)$$

659 Where DS is the docking score, QED and SA are quantitative estimates of drug likeness and  
660 synthesizability respectively.

## 661 A.2 Evaluation metrics

662 Most of the metrics we use are described in detail in Section 5.1. Here, we provide additional  
663 details about the diversity metric. We calculate the diversity of the top 100 molecules sampled by  
664 the algorithm, where higher diversity is considered better given that it increases the chances for  
665 success in further wet lab experiments. In our experiments, we use the diversity evaluator from  
666 TDC [Huang et al., 2021], which defines the diversity of a set of molecules as the average pairwise  
667 Tanimoto similarity between Morgan fingerprints of the molecules. See Section 2 of [Benhenda,  
668 2017] for exact details of how Tanimoto similarity is calculated.

## 669 A.3 Pretraining

670 In this section, we provide more details about the pretraining datasets and models used in our  
671 experiments.

672 **Pretraining datasets.** The ZINC 250k dataset contains approximately 250k molecules from the  
673 ZINC database [Irwin et al., 2012], chosen for their pharmaceutical relevance, moderate size, and  
674 popularity [Gao et al., 2022]. The ChEMBL dataset [Mendez et al., 2019] consists of approximately  
675 2M manually curated drug-like molecules. The other 3 datasets consist of randomly selected  
676 subsets of the ZINC-15 dataset [Sterling and Irwin, 2015b] that obey some chemically imposed  
677 mild constraints [Irwin et al., 2022]. We test three subsets of different sizes: (1) ZINC 1M (2) ZINC  
678 10M, and (3) ZINC 100M, to test the impact of scaling the size of pre-training data. These datasets  
679 and data-subsets, including their vocabularies, will be shared in an easily accessible format upon  
680 acceptance.

681 Removing outliers and unusual non drug-like compounds helps to keep the vocabulary small and  
682 improves the quality of the generative model [Blaschke et al., 2020a]. To achieve this, we filter all  
683 datasets by removing molecules which contain 1) less than 10 or more than 50 heavy atoms and  
684 2) molecules *other than* Carbon, Nitrogen, Oxygen, Fluorine, Silicon, Chlorine and Bromine. We  
685 also canonicalize and sanitize all molecules using RDKit [Landrum et al., 2013]. For experiments  
686 that apply SELFIES, we convert all datasets to SELFIES using the Python API provided by [Krenn  
687 et al., 2020] (Version: 2.1.1).

688 Apart from the experiments shown in the main paper, Appendix C.2 contains additional experiments  
689 comparing text-based RL agents across different pretraining datasets.

690 **Pretraining models.** In Table 4 we provide details about the pretraining modes which we use in  
691 our experiments. Upon acceptance, we will open-source our code and release the pretrained weights  
692 to support reproducible research.

693 We select network sizes that have been commonly used in RL [Blaschke et al., 2020a, Yarats et al.,  
694 2021]. Although conducting a study of scaling the model size [Kaplan et al., 2020] is out of the  
695 scope of our work, we believe that it is a promising direction for future.

696 Since the fully connected model can only take fixed length inputs, we always input a molecular text  
697 padded to a certain maximum length (we used length 100 in our experiments). This padding is done

Table 4: Description of model architectures used for pretraining

Model	Number of Parameters	Description
FC	$1.07 \times 10^7$	FC is a fully connected neural network with 3 hidden layers of 1024 size each.
RNN	$4.17 \times 10^6$	RNN is a recurrent network which consists of 3 GRU layers of hidden sizes 512 each.
TRANSFORMER	$4.78 \times 10^6$	GPT [Brown et al., 2020b] style transformer with 6 layers, 16 heads and 256 embedding dimensions.

698 using a special token [PAD] to convey that corresponding tokens should not be considered while  
 699 deciding the value of the text.

700 **Pretraining experimental details.** We pretrain FC, RNN and transformer architectures on the  
 701 ZINC 250K dataset and pretrain a transformer on all other datasets. All models are pretrained using  
 702 the PyTorch [Paszke et al., 2019] framework. All models used an initial learning rate of  $1e-3$ , with  
 703 a cosine learning rate schedule [Loshchilov and Hutter, 2017]. FC and RNNs used a batch size of  
 704 128 and were trained for 10 epochs. All transformers were trained for 5 epochs, with the largest  
 705 batch size that we could fit in the memory of a single NVIDIA RTX A6000 GPU, for example, a  
 706 batch size of 2048 for pretraining the transformer on ZINC 100M dataset. We made sure that all  
 707 models were trained until convergence. On the ZINC 250K SMILES dataset, the FC, the RNN and  
 708 the transformer model achieved a validation loss of 29.417, 22.507, and 22.923 respectively.

#### 709 A.4 RL finetuning

710 The pretrained model is further trained using the policy gradient algorithm, REINFORCE [Sutton  
 711 et al., 1999b]. Given the reward function  $r(s_H)$  corresponding to the text  $s_T$ , this algorithm  
 712 optimizes the loss function

$$\min_{\theta} J(\theta) = - \left[ \sum_{t=1}^H \log p_{\theta}(a_t = A_t \mid A_{t-1}, \dots, A_0) r(s_H = [A_0, \dots, A_H]) \right], \quad (5)$$

713 where  $A_t$  is the token sampled by the agent at time-step  $t$ .

## 714 B Hyperparameter tuning

715 We conduct a common hyperparameter tuning strategy for all experiments. Specifically, we conduct  
 716 hyperparameter tuning for

- 717 • Learning rate for different architectures Figure 4 and text grammars Figure 3.
- 718 • Coefficients for different likelihood regularizations Figure 7.
- 719 • Coefficients for KL regularization loss term Figure 6.

720 We select three tasks from the **pytdc** tasks, i.e., troglitazone\_rediscovery, sitagliptin\_mpo, and  
 721 median2 for hyperparameter tuning. For each hyperparameter, we select a set of 5 evenly spaced  
 722 realistic values and run 5 random seeds of RL experiments per hyperparameter value. We select the  
 723 hyperparameter value that achieves the best average score of the top-100 molecules as the final value  
 724 for running all the experiments. We report the hyperparameters used for the policy gradient training  
 725 in table 5.

## 726 C Results

### 727 C.1 Text representations and architectures for RL

728 Here, we present additional results from subsection 5.3. Figure 8 shows that SMILES are a better  
 729 molecular grammar when compared to SELFIES across all architectures, for the text based RL

Table 5: Hyperparameters

Name	Value
Maximum number of unique molecules	25000
Learning rate	$5.00 \times 10^{-4}$ RNN and FC $1.00 \times 10^{-4}$ Transformer
Batch size	64
Log p coefficient	5
KL coefficient	$1.00 \times 10^{-3}$

730 algorithms that we consider. Figure 9 compares various architectures, while keeping the molecular  
 731 grammar fixed to SELFIES. The results in Figure 9 reflect our findings in Figure 4 that no single  
 732 architecture clearly outperforms for molecular text-based RL. It also shows the reward hacking  
 733 behavior of the **docking** tasks by the FC based RL agent.

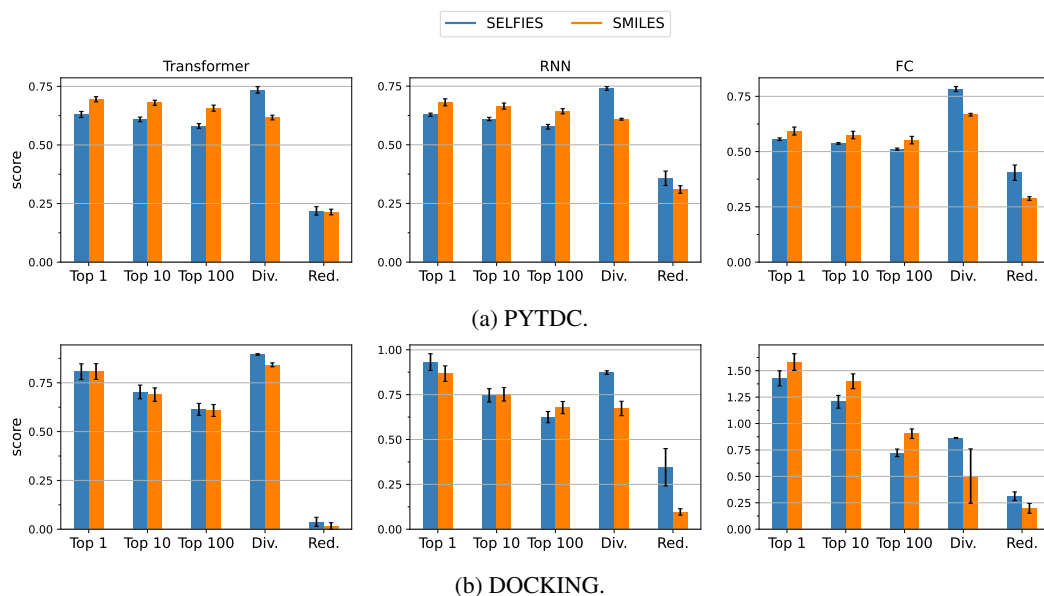


Figure 8: Comparison between SELFIES and SMILES across different architectures. These figures are the individual plots corresponding to the normalised plot show in Figure 3.

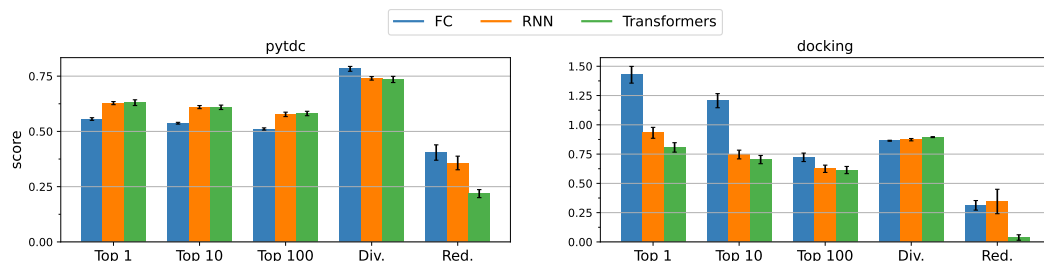


Figure 9: **Comparison of different policy architectures (SELFIES):** No single architecture clearly outperforms for molecular text-based RL. Although FC does better on the docking tasks, our analysis shows that it learns to exploit the docking function as opposed to designing high-value molecules.

734 The reason for lower value molecules for SELFIES environments can be explained by the SELFIES  
 735 grammar that induces a flat optimization landscape. Many SELFIES strings can correspond to the

736 same molecule, and in fact, once an *invalid* action is taken, any subsequent sequence of tokens will  
 737 be ignored on the resulting molecule. This makes exploration of new molecules difficult [Krenn  
 738 et al., 2020, Gao et al., 2022]. On the other hand, the benefit of SELFIES over SMILES in  
 739 eliminating invalid molecule generation is mitigated by our pretraining process, which initializes  
 740 SMILES-based policies with a strong bias toward generating valid molecules. Overall, we find  
 741 that SMILES-based policies, when combined with pretraining, are more effective at exploring and  
 742 finding high-value molecules.

## 743 C.2 Pretraining for RL

744 Figure 2 (right) shows the top docking scores obtained by RL agents pre-trained on different datasets  
 745 when trained with on the **augmented docking** tasks. In Figure 10, we show the actual augmented  
 746 rewards obtained by the RL agent. These results suggest that the augmented docking score is a  
 747 complex reward function as the RL agent is achieved minimal improvement over the prior agent. To  
 748 verify this hypothesis, we increased the molecule budget of the RL agent by 10 times. We indeed  
 749 see that RL agents corresponding to all prior-datasets exhibit considerable improvement. Text-based  
 750 RL algorithms learn to search more efficiently when provided with more compute.

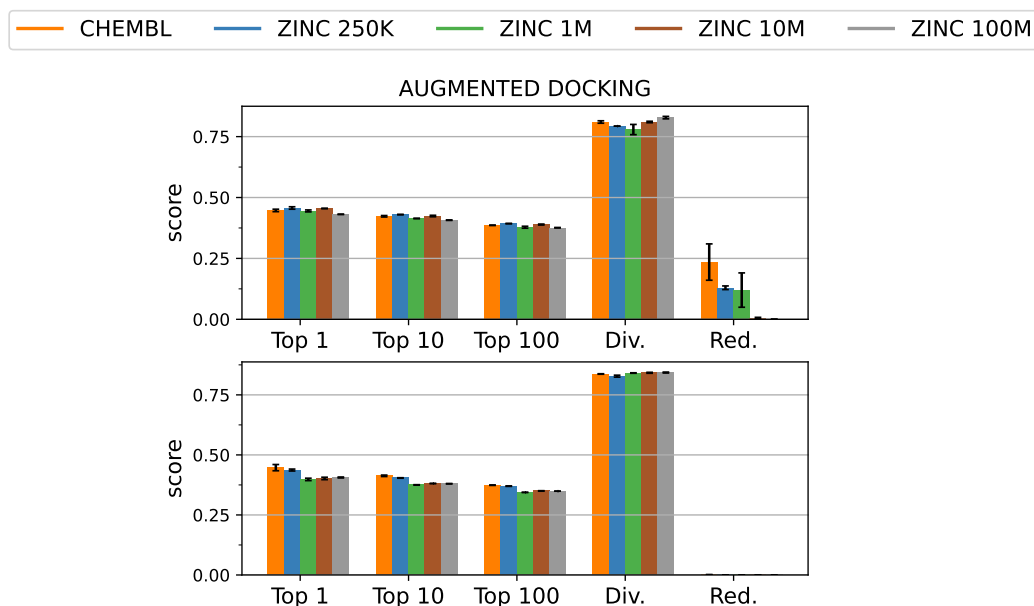


Figure 10: This figure shows the augmented rewards obtained by the RL agents (Top) and data quality (Bottom) of different datasets. See subsection A.1 for how the augmented reward is calculated.

## 751 C.3 Reward Hacking

752 Figure 4 and Figure 9 show that text-based RL agents that are trained using fully connected neural  
 753 networks are able to obtain unusually high rewards. This is probably because it is easier for  
 754 FC agents to find actions that exploit the local structure of the reward function as RNNs and  
 755 Transformers are inductively biased to find global solutions. This highlights an undesirable type  
 756 of reward function hacking by the FC agent which provides high rewards for molecules with long  
 757 strings of Carbon and Nitrogen atoms together. Similar to prior work [Lee et al., 2023], we augment  
 758 the docking scores with objectives for drug-like and synthesizable molecules. See Appendix A for  
 759 details of this task and Figure 2 and Figure 10 for results corresponding to this task. Our initial  
 760 results on this task (Figure 2 and Figure 10) suggested that the augmented reward function was  
 761 more aligned towards chemically relevant molecules. We also noticed that the RL agents were not  
 762 able to improve a lot over the prior baselines for this task. To verify whether the low performance  
 763 of RL agents was because of less training data or the augmented reward function was indeed a more  
 764 realistic and robust reward function, we repeated the experiments in Figure 10 with a ten times

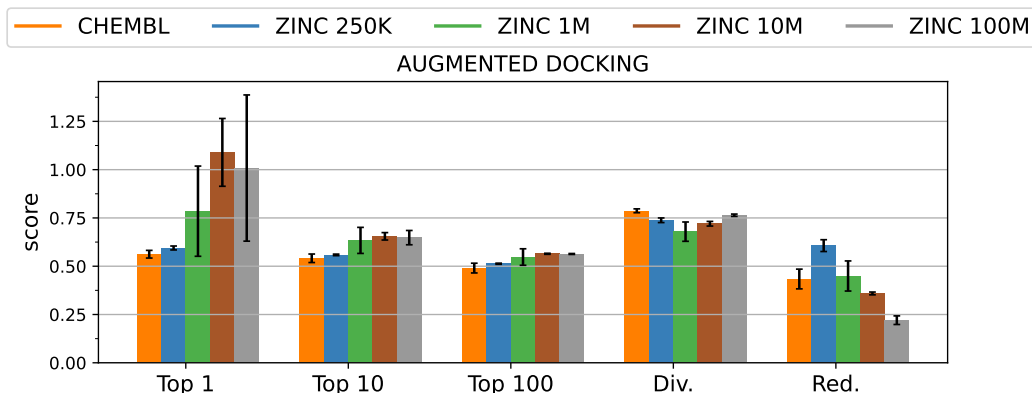


Figure 11: This figure shows the augmented rewards obtained by the RL agents trained for 10 times more molecules.

765 higher training budget. Given more data, all RL agents showed considerable improvements over  
 766 the priors. This experiment also revealed that the agents pre-trained on ZINC 1M, ZINC 10M and  
 767 ZINC 100M, were able to exploit the reward function to generate unrealistic yet highly rewarding  
 768 molecules. These molecules have unusually low docking scores (less than -20). Our results highlight  
 769 the need for an aligned and a more robust reward function to generate molecules for docking protein  
 770 targets.

#### 771 C.4 Additional results for the importance of algorithmic choices for text-based RL.

772 Section 5.4 compares various algorithmic components like replay buffers, hill climbing, KL  
 773 regularisation towards the pretrained policy, and likelihood penalties and show results for **PYDC**  
 774 tasks. In this section, we repeat all the experiments from Section 5.4 on **augmented docking** tasks as  
 775 well and reach the same conclusions. In Figure 14 we see that using the hill-climb buffer results in a  
 776 significant performance boost, whereas using a standard buffer does not contribute much. Figure 15  
 777 shows that Log P regularization is a better choice for efficient exploration when it comes to text-  
 778 based RL algorithms. In Figure 16 show that penalising the policy to move away from the pretrained  
 779 policy does not improve performance.

#### 780 C.5 Instability of transformers for online RL.

781 Many works [Parisotto et al., 2019] have pointed out the instability of training transformers using  
 782 online reinforcement learning. To understand this in the context of text based RL, we compare a  
 783 transformer and an RNN based agent on the augmented docking task. To probe whether pronounced  
 784 effects of this instability are seen, we train both agents for 10 times more molecules (250K  
 785 molecules). In figure Figure 12, we see that both agents perform comparably across all docking  
 786 targets.

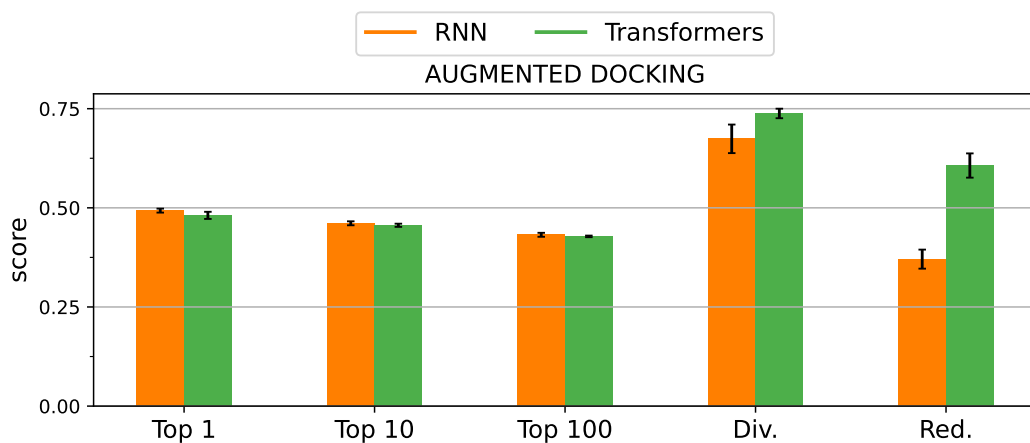


Figure 12: A transformer and an RNN based RL agent trained for 10 times more molecules on augmented docking scores.

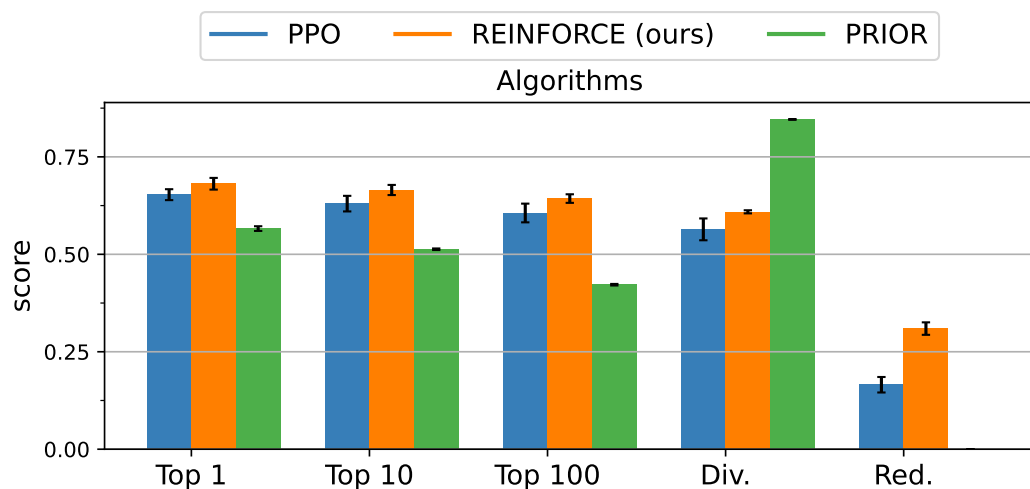


Figure 13: **Comparison with another RL algorithm PPO.** In control tasks complex algorithms like PPO [Schulman et al., 2017] are known to outperform the vanilla policy gradient algorithm. But on the molecular optimization tasks of PyTDC, our results indicate that vanilla policy gradient algorithms are more stable than actor critic algorithms like PPO and achieve higher performance. This resonates with the findings of previous work in molecular optimization [Cieplinski et al., 2021, Gao et al., 2022].

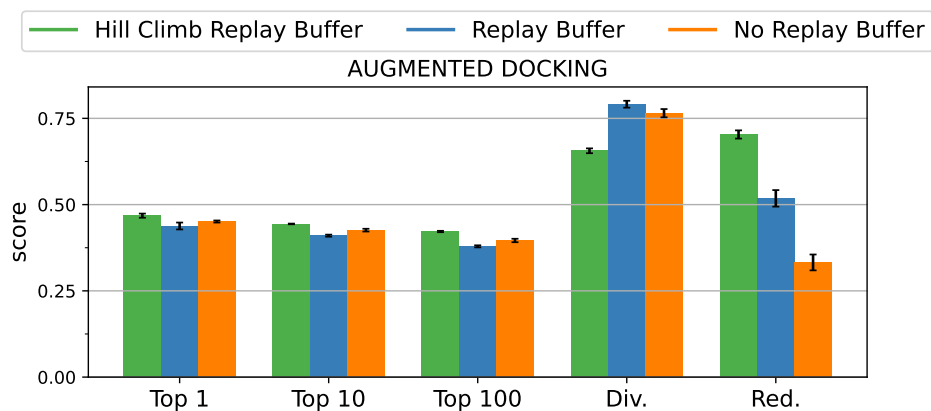


Figure 14: **Do replay buffers help?**

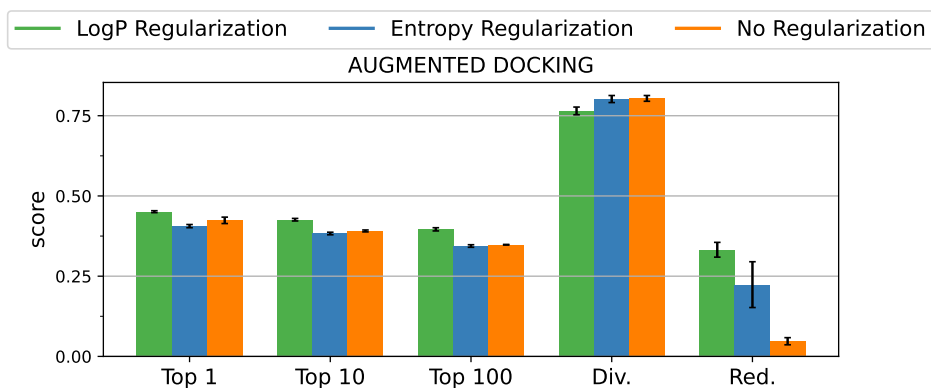


Figure 15: **Comparison of different likelihood penalization for efficient exploration**

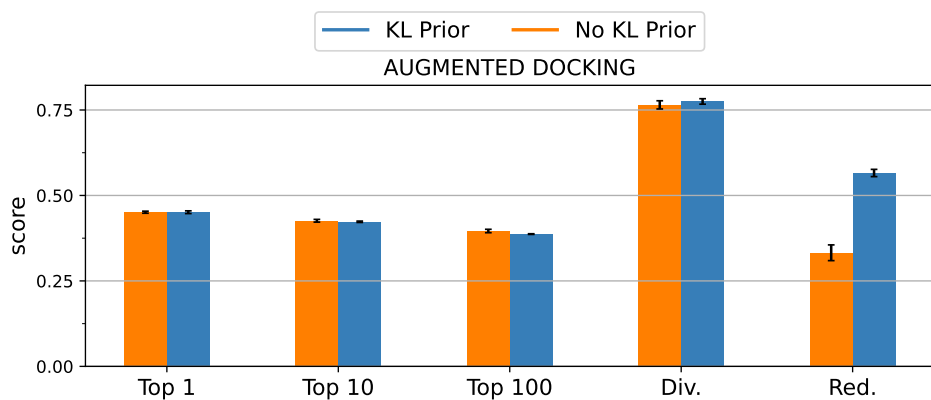


Figure 16: **Is KL regularisation with a prior necessary?**