
SMORE-DRL: Scalable Multi-Objective Robust and Efficient Deep Reinforcement Learning for Molecular Optimization

Aws Al-Jumaily

Department of Computer Science
Brock University
St. Catharines, ON, Canada
aa17qz@brocku.ca

Nicholas Aksamit

Department of Computer Science
Brock University
St. Catharines, ON, Canada
na16dg@brocku.ca

Yage Zhang

Department of Computer Science
Brock University
St. Catharines, ON, Canada
yz20cf@brocku.ca

Mohammad S. Ghaemi

Digital Technologies Research Centre
National Research Council Canada
Ottawa, ON, Canada
mohammadsajjad.ghaemi@nrc-cnrc.gc.ca

Jinqiang Hou

Department of Chemistry
Lakehead University
Thunder Bay, ON, Canada
jhou3@lakeheadu.ca

Hsu Kiang Ooi

Digital Technologies Research Centre
National Research Council Canada
Ottawa, ON, Canada
hsukiang.ooi@nrc-cnrc.gc.ca

Yifeng Li

Department of Computer Science
Brock University
St. Catharines, ON, Canada
yli2@brocku.ca

Abstract

The adoption of AI techniques within the domain of drug design provides an opportunity of systematic and efficient exploration of the vast chemical search space. In recent years, advancements in this domain have been driven by AI frameworks, including deep reinforcement learning (DRL). However, the scalability and performance of existing methodologies are constrained by prolonged training periods and inefficient sample data utilization. Furthermore, generalization capabilities of these models have not been fully investigated. To overcome these limitations, we take a multi-objective optimization perspective and introduce SMORE-DRL, a fragment and transformer-based multi-objective DRL architecture for the optimization of molecules across multiple pharmacological properties, including binding affinity to a cancer protein target. Our approach involves pretraining a transformer-encoder model on molecules encoded by a novel hybrid fragment-SMILES representation method. Fine-tuning is performed through a novel gradient-alignment-based DRL, where lead molecules are optimized by selecting and replacing their fragments with alternatives from a fragment dictionary, ultimately resulting in more desirable drug candidates. Our findings indicate that SMORE-DRL is superior to current DRL models for lead optimization in terms of quality, efficiency, scalability, and robust-

ness. Furthermore, SMORE-DRL demonstrates the capability of generalizing its optimization process to lead molecules that are not present during the pretraining or fine-tuning phases.

1 Introduction

Successfully developing a drug is a tremendously time-consuming, expensive and difficult endeavour. On average, it takes 10-15 years and costs \$1-2 billion USD to deliver a new drug to market (Sun et al., 2022). The objective of drug design is to identify molecules that exhibit multiple pharmacological properties characteristic of pharmaceutical-grade drugs, ensuring they are safe and efficacious. Drug design thus can be modelled as a multi-objective optimization (MOO) problem. One of the main challenges of drug design is effectively navigating the immense chemical search space, which is estimated to contain $10^{20} - 10^{200}$ possible drug-like molecules (Brown, 2015).

AI methods have demonstrated promise in addressing this challenge, primarily through molecular optimization (Bolcato, Heid, and Boström, 2022; Chen et al., 2021; Fu et al., 2022; Fu et al., 2021; Jin, Barzilay, and Jaakkola, 2018; Loeffler et al., 2024; Ståhl et al., 2019; Schneuing et al., 2022; Spiegel and Durrant, 2020; Zhou et al., 2019) and molecular generation (Ai et al., 2024; Bengio et al., 2021; De Cao and Kipf, 2018; Goel et al., 2021; Gottipati et al., 2021; Hoogeboom et al., 2022; Igashov et al., 2024; Liu et al., 2018; Pereira et al., 2021; Popova, Isayev, and Tropsha, 2018; Sattarov et al., 2019; Tang et al., 2023; Wang and Zhu, 2024; Yang et al., 2021a; Yang et al., 2021b; Zhu et al., 2024). Molecular optimization involves making minor modifications to a lead molecule to enhance its drug-like properties while maintaining structural similarity. Since molecules with similar structures are expected to exhibit comparable behaviors, this approach aims to prevent the generation of unrealistic or undesirable molecules. This differs from molecular generation, where the task is to generate novel and diverse compounds from scratch (Ståhl et al., 2019). These AI molecular optimization and generation methods span a wide range of frameworks, including those leveraging genetic algorithms (GAs) (Ahn et al., 2020; Fu et al., 2022; Spiegel and Durrant, 2020), autoencoders (Chen et al., 2021; Jin, Barzilay, and Jaakkola, 2018; Liu et al., 2018; Sattarov et al., 2019), generative adversarial networks (GANs) (De Cao and Kipf, 2018; Tang et al., 2023), flow-based models (Bengio et al., 2021; Zhu et al., 2024), transformer-based models (Ai et al., 2024; Liu et al., 2023; Yang et al., 2021a), equivariant diffusion-based models (Hoogeboom et al., 2022; Igashov et al., 2024; Schneuing et al., 2022), and deep reinforcement learning (DRL) models (Bolcato, Heid, and Boström, 2022; Goel et al., 2021; Gottipati et al., 2021; Loeffler et al., 2024; Pereira et al., 2021; Popova, Isayev, and Tropsha, 2018; Ståhl et al., 2019; Tang et al., 2023; Wang and Zhu, 2024; Yang et al., 2021b; Zhou et al., 2019). However, present methodologies are hindered by lengthy training requirements and sub-optimal use of training data, resulting in impaired scalability and performance.

In this work, we present SMORE-DRL (**Scalable Multi-Objective Robust and Efficient Deep Reinforcement Learning**), a gradient-alignment-based multi-objective DRL (MODRL) framework for molecular optimization. The key contributions of this research include: (1) molecular optimization is modelled naturally as a Pareto-based multi-objective reinforcement learning problem where the challenges of gradient dominance and conflict are addressed with gradient alignment inspired by a study from multi-task learning; (2) a novel molecular tokenization strategy is proposed to represent a molecule as a hybrid of fragments and SMILES, enabling efficient policy learning and effective representation of any new molecules; and (3) a synergistic integration of gradient alignment, hybrid fragment-SMILES representation, contrastive learning, and a transformer-encoder allows for scalability and generalization capability superior to existing MODRL methods. Moreover, SMORE-DRL demonstrated its ability to effectively scale and generalize its optimization process to new molecules after fine-tuning. This is a particularly notable aspect of our work, as existing DRL methods lack scalability and their generalization capacities are under-explored in current literature.

2 Related Work

The fundamental techniques of MODRL and aligned multi-task learning, closely related to this study, are reviewed below, followed by a brief overview of existing AI methods for molecular optimization. In addition, see Appendix A.1 for a review of transformer-encoder architectures and MLM.

2.1 Deep Reinforcement Learning

Reinforcement learning (RL) agents learn through a trial-and-error process guided by the Markov decision process (MDP). See Appendix A.2 for an introduction to basic RL concepts. When the RL task entails exploring a vast state or action space, as is often the case in drug design, learning an exact optimal policy or value function can become computationally intractable. Thus, DRL is used to approximate policies or value functions (Arulkumaran et al., 2017). The actor-critic framework approximates both and has been leveraged by various drug development frameworks (Al-Jumaily et al., 2023; Goel et al., 2021; Gottipati et al., 2021; Pereira et al., 2021; Popova, Isayev, and Tropsha, 2018; Ståhl et al., 2019; Tang et al., 2023; Wang and Zhu, 2024; Yang et al., 2021b). The actor model is responsible for learning a parameterized policy π_{θ_A} . This is guided by feedback known as temporal difference (TD) error from the critic model, which evaluates the actor’s actions based on the state. One approach to this is by learning the advantage function $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$, which measures the desirability of taking action a compared to alternative actions available from state s (Graesser and Keng, 2019).

2.2 Multi-Objective Deep Reinforcement Learning

MODRL is a domain within machine learning and also a family within MOO focused on simultaneously optimizing two or more objectives (Liu, Xu, and Hu, 2015). In an MODRL setting, the reward function is extended to a vector of size K , which represents K different objectives (Nguyen et al., 2020).

2.3 Aligned Multi-Task Learning

Two potential issues that arise when solving an MOO problem directly using gradient descent are dominating and conflicting gradients. A dominating objective gradient is characterized by the largest magnitude, which leads to a bias in the solution favouring the corresponding task (Senushkin et al., 2023). When two objective gradients are conflicting, an increase in the solution towards one objective decreases the solution for the conflicting objective. Conflicting gradients are characterized by having a negative cosine similarity (Yu et al., 2020). To address these challenges in the context of multi-task learning, Senushkin et al., 2023 propose aligned-multi-task learning (AMTL). Let $\mathcal{L}_k(\theta)$ represent the objective of task k , where there are $K > 1$ tasks that are associated with a set of model parameters θ . The training objective is to converge to a set of θ^* defined as follows:

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^m} \left\{ \mathcal{L}_0(\theta) \stackrel{\text{def}}{=} \sum_{k=1}^K \frac{1}{K} \mathcal{L}_k(\theta) \right\}. \quad (1)$$

To mitigate conflicting and dominating gradients, AMTL aligns the principal components of an initial linear system of gradients. This process can be interpreted as re-scaling the axes of a coordinate system that is determined by the principal components, such that the minimal singular value of the gradient matrices is identified and all other singular values are adjusted to match it, resulting in aligned gradients. Subsequently, these aligned gradients are combined into a common gradient (Senushkin et al., 2023). Inspired by AMTL which formulates a multi-task learning problem to an MOO problem, we integrate the same gradient alignment technique to solve MODRL for drug design in this study. A detailed description of AMTL-based MODRL can be found in Appendix A.3.

2.4 Existing AI Approaches for Molecular Optimization

AutoGrow4, proposed by Spiegel and Durrant (2020), uses a genetic algorithm (GA) for molecular optimization, where fitness is calculated as the docking score for a target protein. Given an initial set of lead molecules, in the form of SMILES strings, three operations are performed to evolve the molecules: elitism (selecting the fittest molecules to advance to the next generation), mutation (altering molecules through conducting chemical reactions), and crossover (joining two molecules into a new molecule). Fu et al. (2022) built upon AutoGrow4 with Reinforced Genetic Algorithm (RGA), which reinterprets the evolutionary process as a Markov Decision Process (MDP). Initially, $E(3)$ -equivariant neural networks (ENNs) are pre-trained to predict the binding affinity of 3D target-ligand complexes. For fine-tuning, the ENNs take the form of policy networks, which perform crossover and mutation operations on lead molecules across a set number of timesteps, with the reward being the binding affinity to the target protein.

The Junction Tree Variational Autoencoder (JT-VAE) framework, developed by Jin, Barzilay, and Jaakkola (2018), uses junction trees and graph representations to represent molecules. Molecules are encoded into a two-part latent representation, with a tree encoder capturing the junction tree structure and a graph encoder capturing the molecular graph. To decode the latent representation, a tree decoder reconstructs the junction tree, and then a graph decoder, conditioned on the decoded tree and graph latent representation, reconstructs the molecular graph. For the optimization task, a neural network is introduced to predict the property value of a molecule. This network is trained alongside JT-VAE to predict the property value for a given latent representation of a molecule. Starting from a lead molecule, a number of gradient ascent steps in the latent space are performed to improve the predicted score.

Introduced by Fu et al. (2021), MIMOSA utilizes Markov Chain Monte Carlo (MCMC) sampling to optimize molecules represented as graphs. Two GNNs are pre-trained for distinct tasks: predicting masked nodes within a molecular graph and determining molecular topology, specifically whether a node will expand. These GNNs guide the modification of lead molecules through substructure replacement, addition, or deletion. To sample molecules for the next generation, an unnormalized target distribution is constructed based on molecular property scores, which reflect the favourability of each molecule.

DiffSBDD, developed by Schneuing et al. (2022), is an $SE(3)$ -equivariant 3D-conditional diffusion model that generates molecules specifically conditioned on target protein pockets. Protein and ligand point clouds are represented as fully-connected graphs and processed using $SE(3)$ -equivariant graph neural networks (EGNNs). During training, varying levels of noise are introduced to the 3D structures of real ligands, and a neural network learns to reconstruct their original, noise-free features. For molecular optimization, lead molecules are partially noised and the trained model is tasked with denoising them, producing new drug candidates within the same area of the chemical space as the original leads.

Zhou et al., 2019 developed Molecule Deep Q-Networks (MolDQN), a multi-objective molecular optimization framework that implements double deep Q-learning (DDQN) and randomized value functions. For the optimization task, an episode starts with a seed lead molecule and in each timestep, MolDQN optimizes the molecule through one of the following actions: (1) atom addition, (2) bond addition, and (3) bond removal. A linear weighted sum method is used for MOO.

Deep Fragment-based Multi-Parameter Optimization (DeepFMPO), introduced by Ståhl et al., 2019, is an actor-critic multi-objective method for molecular optimization. In this work, a library of fragments is derived from a set of lead molecules and fragments are encoded using a balanced binary tree such that similar molecules have similar binary encoding. One modification step involves replacing a fragment in the lead molecule with a similar fragment from the fragment library. A constrained reward function is used, where a molecule is either assigned a constant positive reward for each objective achieved or a reward of zero. If all objectives are met, the reward is doubled. A dynamic reward mechanism is also implemented, where the model is penalized if it begins to under-perform compared to previous epochs.

Bolcato, Heid, and Boström, 2022 expand DeepFMPO to include 3D-shape and electrostatics in the similarity measurements. This extension was applied because a seemingly minor alteration to a SMILES string can significantly impact its 3D structure. Consequently, the revised representation of fragments is suggested to achieve a more precise similarity measure.

3 Methods

In this section, we discuss the data pre-processing, pretraining, and fine-tuning steps carried out in our SMORE-DRL framework.

3.1 Data Preparation: Fragments-SMILES Hybrid Tokenization Strategy

The dataset used for transformer-encoder pretraining is the MolGen task of the Therapeutics Data Commons (TDC) (Huang et al., 2021), a set of the ChEMBL, MOSES, and ZINC-250K datasets. We canonicalized all SMILES strings, and only kept strings with a maximum of 100 characters, resulting in a pretraining dataset of 4 million molecules. We then fragmented each molecule using the fragmentation method from HierVAE by Jin, Barzilay, and Jaakkola, 2020, which breaks single

bonds extending from ring atoms (Ståhl et al., 2019). This method is also used by DeepFMPO (Ståhl et al., 2019) and DeepFMPOv3D (Bolcato, Heid, and Boström, 2022).

While fragmentation is a good technique for reducing the chemical search space, it may result in a vast token dictionary size. Figure 5 in Appendix A.4 is a fragment frequency chart based on the 4-million molecule dataset, which shows that nearly 68,000 fragments were extracted. Most of these fragments are rarely encountered in the dataset, with 95% appearing fewer than 100 times. Moreover, building a token dictionary solely from the pretraining dataset will create obstacles during fine-tuning tasks. Given the sparse nature of fragment occurrences, it is likely that molecules used for fine-tuning will contain fragments not present in the dictionary, especially if the fine-tuning dataset differs from the one used for pretraining.

To reduce the dictionary size while still representing fragments that are absent or infrequently encountered, we propose a novel hybrid tokenization strategy that uses both fragments and SMILES. Following the construction of a fragment dictionary from the pretraining dataset, we append tokens for SMILES and exclude all fragments that appear less than twice. This reduces the dictionary size from 68,000 to approximately 41,130 tokens. As a result, if a molecule contains a fragment not found in the reduced dictionary, that fragment is represented atom by atom. See Appendix A.4 for a diagram of the hybrid tokenization strategy. As part of our ablation studies in Section 4.2, we demonstrate that further reducing the token dictionary by retaining only the most frequently encountered fragments (resulting in molecules being primarily represented by SMILES atoms) hinders training performance.

3.2 Pretraining

SMORE-DRL utilizes a transformer-encoder model inspired by architectural aspects of the Bidirectional Encoder Representations from Transformers (BERT) model introduced in MTL-BERT by Zhang et al., 2022, a multi-task learning model pretrained on SMILES strings and fine-tuned for downstream ADMET tasks. Rather than representing molecules by SMILES atom tokens as was done in MTL-BERT, we adopt our hybrid fragment token representation. SMORE-DRL employs a combination of two pretraining tasks: MLM and contrastive learning.

3.2.1 Masked Language Model (MLM)

Given a training batch of molecules, they are fragmented and encoded into their token representations. Unlike the static masking technique used for the MLM in the original BERT model, where sequences are masked once and reused throughout training, we employ the dynamic approach introduced by Liu et al., 2019 in Robustly Optimized BERT Pretraining Approach (RoBERTa). In each training batch, 20% of the tokens are randomly selected for masking, with 90% of those being masked. If a selected token is part of a sequence of atom tokens representing a fragment that does not exist in the token dictionary, 20% of that fragment’s atom token sequence is also masked. The masked molecule token sequence is then passed into the encoder model, which attempts to accurately reconstruct the original values of the masked tokens. See Appendix A.5 for a diagram of the MLM process.

3.2.2 Contrastive Learning

We further refine the SMORE-DRL’s contextual understanding of molecules by allowing it to align its representations of similar molecules. This is particularly valuable during fine-tuning, where the model is tasked with optimizing lead molecules over multiple timesteps while ensuring that the optimized molecules in the earlier timesteps retain chemical similarity to the original lead molecule. To accomplish this, we introduce a straightforward contrastive learning technique that builds upon the MLM approach. Rather than directly masking tokens in the fragment sequence as previously described, a “separation” token is inserted at the end of the sequence, followed by an augmented version of that sequence. Augmentation involves randomly selecting a token to replace with a fragment token from the token dictionary. If the selected token belongs to a sequence of atom tokens representing a fragment that is not in the token dictionary, the entire SMILES atom sequence for that fragment is replaced with a randomly selected fragment token. The masking process consists of keeping the original molecule sequence fully visible to the model, while masking 25% of the augmented sequence using the same technique described earlier. See Appendix A.5 for a diagram of the contrastive learning process.

3.3 Fine-Tuning

For the fine-tuning phase, SMORE-DRL utilizes a novel multi-objective actor-critic framework with three pretrained encoder models: a masker, an actor, and a critic. The MDP cycle unfolds as follows: (1) starting from the initial state, represented by the token sequence of the lead molecule, the masker model modifies this sequence by masking certain tokens, resulting in a new state, (2) the masked sequence is processed by the actor model, which replaces the masked tokens with alternatives from the token dictionary, creating a further updated state; (3) this state is evaluated by the critic model, which assigns it a reward value; and (4) the updated state is passed to the reward system for the true reward value. This process repeats iteratively across all timesteps.

3.3.1 Agents

Masker Model: The masker model, denoted as π_{θ_M} , is responsible for selecting which tokens to mask from the lead molecule token sequence. At each timestep, it masks at least one token and up to 70% of the token sequence. The model is designed to prefer masking fragment tokens over SMILES atom tokens. The loss for the masker model is:

$$L(\theta_M) = \frac{1}{T} \sum_{t=0}^T \sum_{k=0}^K \left(-\hat{A}_k^{\pi_{\theta_A}}(s_t, a_t) \log \pi_{\theta_M}(a_t | s_t) \right). \quad (2)$$

Actor Model: After the lead molecule token sequence is masked by the masker model, it is passed to the actor model, denoted as π_{θ_A} . The actor model utilizes the same training head that was used during the pretraining phase. Hence, its task is to replace the masked tokens with tokens from the token dictionary. However, rather than focusing on recovering the original tokens, the actor’s task is to replace the masked tokens with new tokens so that the resulting sequence represents an optimized yet chemically similar version of the lead molecule. The actor model employs AMTL (Senushkin et al., 2023) to identify a common objective gradient, thereby avoiding conflicting and dominating gradients, which ensures that all molecular properties are optimized equally. This process involves obtaining a gradient matrix that collects all K objective gradients, represented as $\mathbf{G} = \{\mathbf{g}_1, \dots, \mathbf{g}_K\}$, where $\mathbf{g}_k = \nabla L_k(\theta_A)$ and

$$L_k(\theta_A) = \frac{1}{T} \sum_{t=0}^T \left(-\hat{A}_k^{\pi_{\theta_A}}(s_t, a_t) \log \pi_{\theta_A}(a_t | s_t) \right). \quad (3)$$

\mathbf{G} is then processed into the gradient matrix alignment algorithm to compute a common gradient, which is used to update the model. Our AMTL-based actor model optimization algorithm is given in Algorithm 1 of Appendix A.3. It is crucial for the masker and actor to work in tandem. If the actor performs well, this will be reflected in the the masker’s loss, as the masker utilizes the advantage function derived from the actor model’s policy. The fine-tuning process is illustrated in Figure 1.

Critic Model: The optimized token sequence is then fed into the critic model, V_{θ_C} , which generates a vector of size K , corresponding to K properties. The critic’s output reflects its assessment of the desirability of sequence as a potential drug candidate. The loss of the critic model is:

$$L(\theta_C) = \frac{1}{T} \sum_{t=0}^T \sum_{k=0}^K \left(r_{t,k} + \hat{V}_{\theta_C k}^{\pi_{\theta_A}}(s_{t+1}) - V_{\theta_C k}^{\pi_{\theta_A}}(s_t) \right)^2. \quad (4)$$

3.3.2 Reward System

To assess a molecule’s potential as a drug candidate, we use the following three properties: (1) logarithm of partition coefficient (ClogP), which impacts a drug’s administration, absorption, transport and excretion, (2) synthetic accessibility score (SAS), which measures the difficulty of synthesizing a molecule, and (3) binding affinity score (BAS) to LPA1, which quantifies the binding capability of a drug to a target protein (Brown, 2015; Ertl and Schuffenhauer, 2009; Li et al., 2019). However, the number and type of properties can be tailored to any specific optimization task. RDKit is used for ClogP and SAS calculations, and QuickVina2-GPU-2.1 (Tang et al., 2024) is used to calculate BAS. Lysophosphatidic acid receptor 1 (LPA1/LPAR1), a bioactive lipid mediator primarily derived from membrane phospholipids, is chosen as the target protein for BAS. LPA Receptors (LPARs) have been

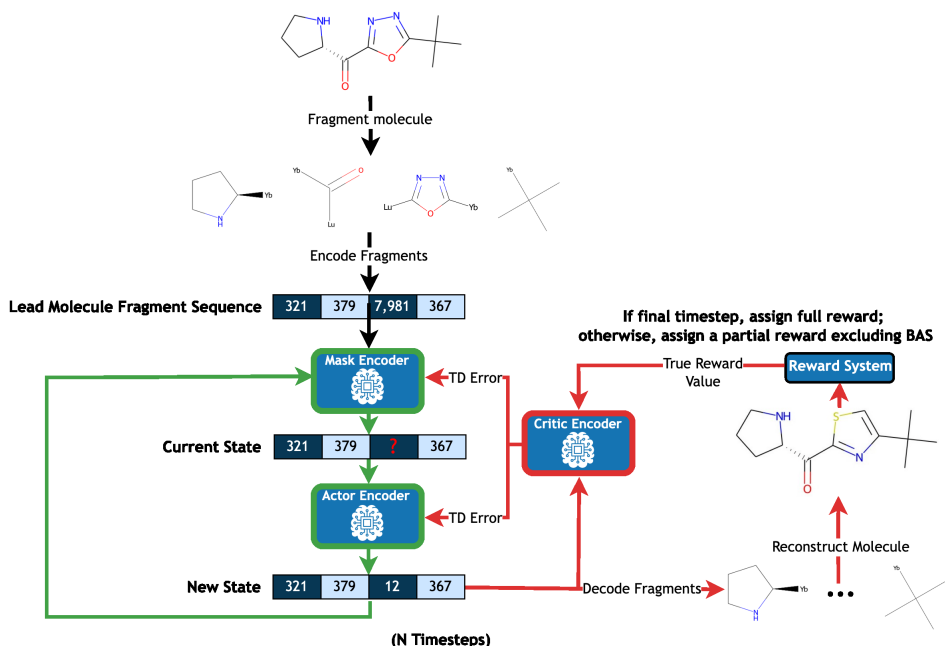


Figure 1: Our MODRL fine-tuning process for one molecule.

found to be over-expressed in multiple types of cancer, with LPA1 specifically expressed in ovarian cancer, breast cancer, liver cancer, gastric cancer, pancreatic cancer, lung cancer, glioblastoma and osteosarcoma. LPA1 promotes metastasis and tumor motility, making it a natural choice for targeting in efforts to inhibit cancer spread and cell movement (Lin, Lin, and Chen, 2021).

To convert a property value to a reward, we treat all properties to be minimized and normalize property values. The reward for molecule m , where property p is ClogP or SAS, is defined as:

$$r_{p,m} = \frac{p_{thresh} - p_m}{p_{thresh} - p_{true_min}}, \quad (5)$$

where p_{thresh} is the target maximum parameter that is set for p , p_m is the p score for molecule m , and p_{true_min} is the true minimum of p . The p_{thresh} parameter controls the difficulty of the optimization task for a given property, with a lower threshold demanding molecules of higher quality. When p is BAS, the absolute mean BAS score of the lead molecule set ($|p_{lead_mean}|$) is used, as no true minimum exists for BAS:

$$r_{p,m} = \frac{p_{thresh} - p_m}{|p_{lead_mean}|}. \quad (6)$$

A widely known challenge of utilizing molecular docking is that it requires significant time (Thafar et al., 2022). To overcome this, for all intermediate timesteps, we assign a partial reward to the molecule that only includes ClogP and SAS, and in the final timestep, BAS is calculated and a full reward comprised of ClogP, SAS and BAS is given.

The following three criteria are also examined for each epoch: (1) validity: the ratio of chemically valid optimized molecules, checked using RDKit, (2) novelty: the ratio of optimized molecules that are different from the lead molecules from which they were derived, and (3) uniqueness: the ratio of unique molecules in the optimized molecules (Mukaidaisi et al., 2022). In each timestep, if a molecule is valid, unique, and novel, it is assigned its reward, else it is assigned a reward of -1 for each objective, for a total reward of -3, if it is the final timestep. If all properties are achieved by a molecule, it is provided with extra reinforcement by doubling the final reward.

Thus, the output of the reward function is $\mathbf{r}^K = [r^1, r^2, \dots, r^K]$, where $K = 3$ for the final timestep and $K = 2$ for all intermediate timesteps. The values of the reward system are listed in Table 1 of the next section.

4 Experiments & Results

4.1 Pretraining

The encoder model was pretrained for six consecutive epochs on a combination of MLM and contrastive learning tasks, where the same training head was used throughout learning. While the main experiment employs the 2-frequency token dictionary (41,130 tokens), we also investigate the effect that different dictionaries have on pretraining and fine-tuning. The pretraining results using the 2-frequency, 100-frequency (3,460 tokens), and 1000-frequency (790 tokens) token dictionaries can be found in Appendix A.6.

4.2 Fine-Tuning

In this section, we demonstrate SMORE-DRL’s molecular optimization performance against three other DRL methods, as well as its scalability and generalization abilities. For the masker, actor and critic models, encoder weights are not frozen. Additionally, the actor model uses the same head from the pretraining phase. We show that these configurations achieve optimal results in our experiments.

4.2.1 Performance Comparison of SMORE-DRL against other DRL Methods for Molecular Optimization

We compare SMORE-DRL’s optimization performance with three other MODRL optimization frameworks: (1) DeepFMPOv3D, (2) DeepFMPO and (3) MolDQN. A primary goal of this paper is to present the scalability of SMORE-DRL. However, it is not feasible to conduct large-scale optimization using thousands of lead molecules to compare with the other models, as these benchmarks lack the efficiency for scalability. As described in their papers, DeepFMPOv3D, DeepFMPO and MolDQN optimized a set of 138, 387 and 800 lead molecules, respectively. To facilitate comparison with these methods, a small-scale dataset was utilized. Scalability and generalizability of SMORE-DRL are demonstrated in the following sections. Challenging property values were selected for the optimization task, as noted in Table 1. The results of this comparative study represent the mean scores of three separate runs for all models.

Table 1: Targeted Molecular Properties and Their Maximal Thresholds and True Minimum/Lead Mean Score

Property	Target Value	True Min/Lead Mean
ClogP	<3	-3 (True Min)
SAS	<2.5 (2.75 for Testing)	1 (True Min)
BAS (LPA1)	<-6	-5.27 (Lead Mean)

All models were trained on a subset of 1,000 lead molecules from the DrugBank database (Wishart et al., 2018) that do not satisfy all properties. SMORE-DRL trained for 70 epochs, during which all lead molecules were optimized over 4 timesteps per epoch. The molecules optimized in the final timestep of the last epoch are used for our comparisons. DeepFMPOv3D and DeepFMPO trained for 1,000 epochs, optimizing random batches of 512 unique molecules per epoch. DeepFMPOv3D performed optimization over 4 timesteps and DeepFMPO used 8 timesteps. In the final epoch, all lead molecules were optimized, and results from this epoch are used for our comparisons. MolDQN was trained for 6,000 epochs. For the first 5,000, a random molecule from the dataset is selected and optimized for 20 timesteps. The final 1,000 epochs focus on optimizing each lead molecule, with the resulting optimized molecules utilized for comparison. To calculate property rewards, all models use the normalization method described in Section 3.3.2. In DeepFMPOv3D and DeepFMPO, a single cumulative reward for all objectives is assigned in the final timestep. For MolDQN, a partial reward excluding BAS is assigned at each intermediate timestep, while a full reward is assigned in the final timestep.

Figure 2 compares each model’s learning progression while Table 2 displays the target property percentages achieved by the lead molecules and the optimized outputs of the various models. MolDQN is excluded from Figure 2 as it optimizes one lead molecule per epoch. Figure 9 of Appendix A.7 depicts the property-wise distributions of the final epoch. The optimization capabilities of SMORE-DRL clearly surpass those of the other models. While all other models struggled heavily with the

optimization task, SMORE-DRL maintained stability throughout training. Further, it successfully optimized 23.54% of molecules to meet all properties, while maintaining comparable computation time. The next best model, DeepFMPO, managed only 0.92%. MolDQN had the worst overall performance, failing to produce a single molecule that achieved all properties.

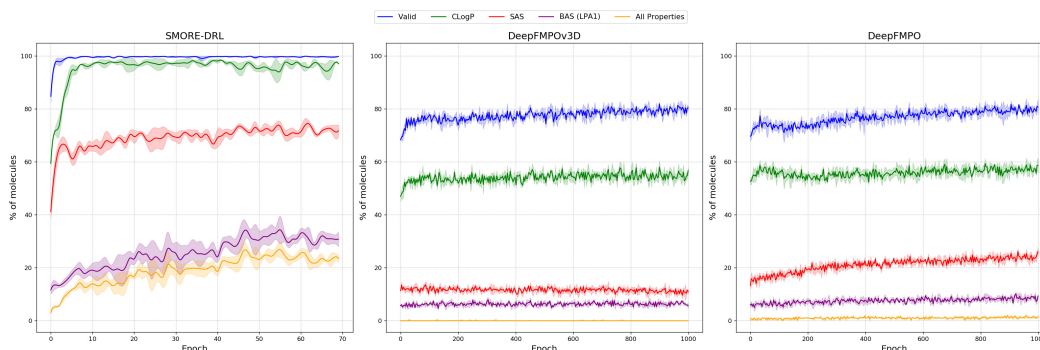


Figure 2: Percentages of valid molecules and those achieving target properties through training for SMORE-DRL, DeepFMPOv3D, and DeepFMPO.

Table 2: Percentage of molecules that satisfy each property from the 1,000 lead molecules and the molecules optimized in the last training epoch by SMORE-DRL, DeepFMPOv3D, DeepFMPO, and MolDQN.

Property	Lead Molecules	SMORE-DRL	DeepFMPOv3D	DeepFMPO	MolDQN
Compute Time	-	~6 hrs	~4.5 hrs	~6 hrs	~5.5 hrs
Validity	-	99.80% (± 0.00)	80.00% (± 2.16)	77.33% (± 3.09)	100% (± 0.00)
Novelty	-	98.52% (± 0.05)	80.01% (± 1.85)	82.99% (± 0.27)	100% (± 0.00)
Uniqueness	-	93.63% (± 1.47)	79.68% (± 2.09)	82.71% (± 0.30)	100% (± 0.00)
ClogP	73.94%	97.25% (± 0.26)	54.90% (± 0.83)	61.78% (± 2.44)	87.6% (± 0.94)
SAS	38.93%	71.83% (± 2.29)	10.96% (± 0.60)	22.86% (± 3.82)	0.70% (± 0.29)
BAS (LPA1)	8.65%	30.85% (± 2.77)	6.60% (± 0.17)	7.86% (± 1.43)	1.73% (± 0.66)
All Properties	0%	23.54% (± 1.56)	0.03% (± 0.05)	0.92% (± 0.42)	0% (± 0.00)

4.2.2 Scalability of SMORE-DRL

To examine scalability, SMORE-DRL was trained for 70 epochs using 10,000 molecules-5,000 from the DrugBank database (Wishart et al., 2018) and 5,000 from the Collection of Open Natural Products (COCONUT) database (Sorokina et al., 2021). COCONUT molecules were included to attempt to test the model’s robustness, as they are different from those typically encountered during pretraining. All lead molecules were optimized over 4 timesteps per epoch, with those from the final timestep of the last epoch used for comparisons. While the baseline model was run five times, ablation studies were also conducted. These included: (1) without the use of AMTL, (2) with a weight emphasis on the BAS reward (0.25 for ClogP, 0.25 for SAS, and 0.5 for BAS), (3) with the freezing of encoder weights for all agents, (4) with the use of pretraining on the 100-frequency dictionary, and (5) with the use of pretraining on the 1000-frequency dictionary. Figure 3 demonstrates that freezing encoder weights and pretraining on the 100-frequency and 1000-frequency dictionaries significantly impair the model’s learning progress. As such, these experiments were limited to a single run of 25 epochs and excluded from further analysis. To evaluate the impact of omitting AMTL and placing greater emphasis on the BAS reward, these model variations were run three times for 70 epochs, while the baseline model ran five times. Presented results are based on run averages.

As displayed in Figure 3, incorporating AMTL improves training stability and enhances the overall quality of optimized molecules. Findings in Table 3 support this by demonstrating that omitting AMTL significantly impairs most properties, namely uniqueness.

The baseline (SMORE-DRL) model can effectively scale to optimize thousands of lead molecules in a timely manner, even if the molecules are structurally distinct from those used during pretraining. This demonstrates its scalability, efficiency and robustness. Figure 10 (Appendix A.7) exhibits the

property-wise distributions, while examples of lead molecules optimized by SMORE-DRL from these experiments are presented in Appendix A.8.

Interestingly, emphasizing the BAS reward did not necessarily produce molecules that were more optimized for BAS compared to the baseline version of SMORE-DRL. A possible explanation for this is that doing so may constrict the model’s exploration of the search space, leading it to focus primarily on BAS. This narrow focus may result in the model converging to a local minimum, hindering its ability to discover more optimal solutions in other areas of the search space. A more effective approach would be to implement a dynamic weighting system, initially assigning equal weights to encourage exploration. Over time, these weights could be adjusted to prioritize specific properties.

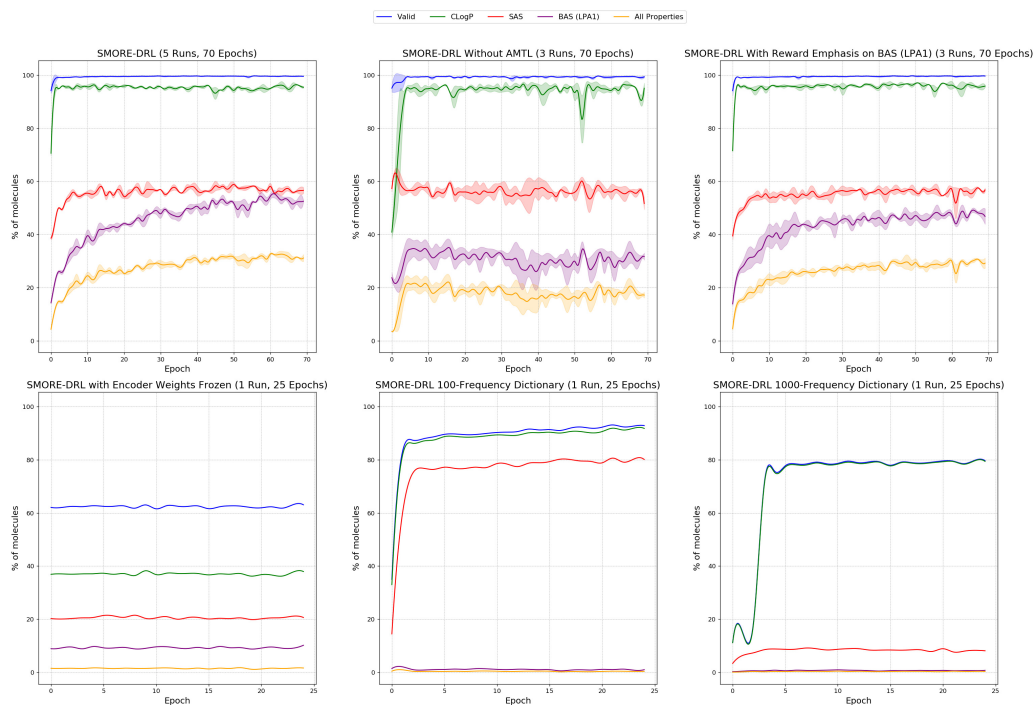


Figure 3: Percentages of valid molecules and those achieving target properties through training for different versions of SMORE-DRL.

Table 3: Percentage of molecules that satisfy each property from the 10,000 lead molecules and the molecules optimized in the last training epoch by SMORE-DRL, SMORE-DRL without AMTL, and SMORE-DRL with a reward emphasis on BAS.

Property	Lead Molecules	SMORE-DRL	No AMTL	BAS Reward Focus
Compute Time	-	~60 hrs	~60 hrs	~60 hrs
Validity	-	99.54% (± 0.08)	99.33% (± 0.34)	99.64% (± 0.10)
Novelty	-	99.98% (± 0.01)	99.97% (± 0.02)	99.99% (± 0.00)
Uniqueness	-	89.31% (± 1.59)	67.30% (± 3.88)	90.49% (± 1.91)
ClogP	62.06%	94.57% (± 1.35)	95.04% (± 1.90)	95.83% (± 0.50)
SAS	25.64%	57.08% (± 0.92)	51.67% (± 2.17)	56.74% (± 0.70)
BAS (LPA1)	9.79%	53.87% (± 2.05)	31.71% (± 1.30)	46.79% (± 2.90)
All Properties	0%	32.22% (± 1.15)	17.20% (± 1.03)	29.20% (± 1.92)

4.2.3 Generalization Performance of SMORE-DRL

While many MODRL drug design frameworks focus on optimization tasks, their ability to generalize and optimize molecules that they have not encountered before remains unexplored. The weights of the baseline SMORE-DRL model from the scalability experiments were frozen, with their optimization

process tested on 40,000 molecules from the COCONUT dataset that differ from those used in the scalability experiments. The following are the average results of the five SMORE-DRL model runs from the fine-tuning phase. To encourage similarity to lead molecules, optimization was restricted to two timesteps. Additionally, the SAS target maximum parameter was increased from 2.5 to 2.75.

SMORE-DRL took 1.25 hours to optimize a test set of 40,000 lead molecules, none of which originally achieved all target properties. 19% of the resulting molecules met all target properties, and all properties were significantly improved (see Table 4). Property-wise distributions are seen in Figure 11 of Appendix A.7, and examples of lead molecule optimized are presented in Appendix A.9.

Table 4: Generalization results – percentage of molecules that satisfy each property from the 40,000 test lead molecule set and the molecules optimized by SMORE-DRL over two timesteps.

Property	Lead Molecules	SMORE-DRL
Avg Compute Time	-	~1.25 hrs
Validity	-	99.38% (± 0.14)
Novelty	-	99.76% (± 0.10)
Uniqueness	-	89.85% (± 0.94)
ClogP	51.13%	85.01% (± 2.60)
SAS	38.91%	51.54% (± 1.47)
BAS (LPA1)	16.98%	38.19% (± 1.69)
All Properties	0%	18.50% (± 0.51)

4.3 Discussion

In this paper, we introduce SMORE-DRL, a novel transformer-based MODRL model for molecular optimization. Three sets of experiments were conducted to evaluate the model’s performance: (1) a comparative study against DeepFMPO, DeepFMPOv3D, and MolDQN, three MODRL molecular optimization models, tasked with optimizing 1,000 lead molecules, (2) a scalability study, where the model was tasked with optimizing 10,000 lead molecules, and (3) a generalization study to assess how well the model, after training in the scalability study, can optimize 40,000 lead molecules in a test scenario.

SMORE-DRL demonstrated outstanding performance in all experiments. In the comparative study, it significantly outperformed all other models. In the scalability study, SMORE-DRL performed efficiently, optimizing a set of lead molecules that did not achieve all properties such that one third of produced molecules satisfied all properties. Additionally, SMORE-DRL’s robustness allowed it to successfully generalize its optimization approach to unseen molecules. With just two modification steps, it improved the lead molecules from 0% to 19% achieving all target properties. The inclusion of AMTL has proven to be a vital component of SMORE-DRL, enhancing training stability and improving the overall performance.

As discussed, the primary objective of molecular optimization is developing a novel molecule similar to a lead molecule, aiming to have both molecules exhibit comparable qualities. As such, the progression of SMORE-DRL’s optimized molecules were analyzed by comparing their similarity to lead molecules and their corresponding rewards across all timesteps. Figure 4 depicts the average similarity and reward for each of the four optimization timesteps performed on 1,000 molecules during the scalability study. To measure similarity, we utilize the method described in DeepFMPO (Stahl et al., 2019), which employs a combination of maximum common substructure Tanimoto similarity and Levenshtein distance. A similarity score greater than or equal to 0.7 indicates high similarity, while a score between 0.5 and 0.7 is considered medium similarity (Loeffler et al., 2024). While SMORE-DRL does not achieve high similarity, it still presents strong results. As seen in Figure 4, there is an inverse correlation between average similarity and average sum of rewards across all objectives, where as similarity decreases, reward increases. This represents a trade-off: restricting the optimization process to minimal modifications of a lead molecule may result in high similarity, but will likely restrict exploration and hinder the development of superior candidates. Nonetheless, the next iteration of SMORE-DRL should balance exploration with maintaining similarity to lead molecules, aiming to generate high-quality compounds without sacrificing similarity. One possible approach involves incorporation of a dynamic similarity component into the reward function, allowing

for exploration in the initial training epochs while penalizing molecules with low similarity to lead molecules in the later epochs.

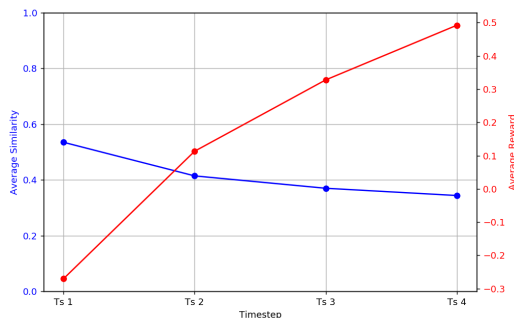


Figure 4: An analysis of: (1) the average similarity to lead molecules and (2) the average sum of rewards across all properties over four optimization timesteps for 1,000 molecules optimized by SMORE-DRL.

5 Conclusion

In this work, we present SMORE-DRL, a scalable gradient-alignment-based MODRL framework for molecular optimization. A novel hybrid fragment-SMILES representation to depict molecules enables SMORE-DRL to select and replace fragments in the lead molecules with alternatives from the fragment dictionary, resulting in improved drug candidates. This is achieved by using three agents: a masker, actor and critic, all pretrained on MLM and contrastive learning tasks. SMORE-DRL excelled as a lead molecule optimizer, significantly outperforming other MODRL models while demonstrating scalability. Furthermore, when evaluated on new molecules post fine-tuning, SMORE-DRL effectively generalized its optimization process. The next development of SMORE-DRL will include additional measures to encourage the model to produce molecules that are as effective as those in the current version, but with greater similarity to lead compounds. The implementation of SMORE-DRL is available at <https://anonymous.4open.science/r/SMORE-DRL-F38B>.

Acknowledgments

This work is supported in part by funds from (1) the AI for Design Challenge Program, National Research Council Canada (AI4D-108 to YL), (2) the Discovery Grant Program, Natural Sciences and Engineering Research Council of Canada (RGPIN 2021-03879 to YL), (3) Canada Research Chair Program (to YL), (4) Canada Foundation for Innovation (to YL), (5) Ontario Research Fund – Small Infrastructure Fund (to YL), and (6) the Vector Scholarship for AI from the Vector Institute (to AAJ)

References

- Ahn, Sungsoo et al. (2020). “Guiding deep molecular optimization with genetic exploration”. In: *Advances in neural information processing systems* 33, pp. 12008–12021.
- Ai, Chengwei et al. (2024). “MTMol-GPT: De novo multi-target molecular generation with transformer-based generative adversarial imitation learning”. In: *PLOS Computational Biology* 20.6, e1012229.
- Arulkumaran, Kai et al. (Aug. 2017). “A Brief Survey of Deep Reinforcement Learning”. In: *IEEE Signal Processing Magazine* 34. DOI: 10.1109/MSP.2017.2743240.
- Bengio, Emmanuel et al. (2021). “Flow network based generative models for non-iterative diverse candidate generation”. In: *Advances in Neural Information Processing Systems* 34, pp. 27381–27394.
- Bolcato, Giovanni, Esther Heid, and Jonas Boström (2022). “On the Value of Using 3D Shape and Electrostatic Similarities in Deep Generative Methods”. In: *Journal of Chemical Information and Modeling* 62.6, pp. 1388–1398.
- Brown, Nathan (2015). *In silico Medicinal Chemistry: Computational Methods to Support Drug Design*. Royal Society of Chemistry.
- Chen, Ziqi et al. (2021). “A deep generative model for molecule optimization via one fragment modification”. In: *Nature machine intelligence* 3.12, pp. 1040–1049.
- De Cao, Nicola and Thomas Kipf (2018). “MolGAN: An implicit generative model for small molecular graphs”. In: *ICML 2018 workshop on Theoretical Foundations and Applications of Deep Generative Models*.
- Devlin, Jacob et al. (June 2019). “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Ed. by Jill Burstein, Christy Doran, and Thamar Solorio. Minneapolis, Minnesota: Association for Computational Linguistics, pp. 4171–4186.
- Ertl, Peter and Ansgar Schuffenhauer (2009). “Estimation of synthetic accessibility score of drug-like molecules based on molecular complexity and fragment contributions”. In: *Journal of Cheminformatics* 1.1, pp. 1–11.
- Fu, Tianfan et al. (2021). “Mimosa: Multi-constraint molecule sampling for molecule optimization”. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 35. 1, pp. 125–133.
- Fu, Tianfan et al. (2022). “Reinforced genetic algorithm for structure-based drug design”. In: *Advances in Neural Information Processing Systems* 35, pp. 12325–12338.
- Goel, Manan et al. (2021). “MoleGuLAR: Molecule generation using reinforcement learning with alternating rewards”. In: *Journal of Chemical Information and Modeling* 61.12, pp. 5815–5826.
- Gottipati, Sai Krishna et al. (2021). “Towered actor critic for handling multiple action types in reinforcement learning for drug discovery”. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 35, pp. 142–150.
- Graesser, Laura and Wah Loon Keng (2019). *Foundations of Deep Reinforcement Learning: Theory and Practice in Python*. Addison-Wesley Professional.
- Hoogeboom, Emiel et al. (2022). “Equivariant diffusion for molecule generation in 3d”. In: *International conference on machine learning*. PMLR, pp. 8867–8887.
- Huang, Kexin et al. (2021). “Therapeutics Data Commons: Machine Learning Datasets and Tasks for Drug Discovery and Development”. In: *NeurIPS Datasets and Benchmarks*.
- Igashov, Iliia et al. (2024). “Equivariant 3D-conditional diffusion model for molecular linker design”. In: *Nature Machine Intelligence*, pp. 1–11.
- Jin, Wengong, Regina Barzilay, and Tommi Jaakkola (2018). “Junction tree variational autoencoder for molecular graph generation”. In: *International conference on machine learning*. PMLR, pp. 2323–2332.
- (2020). “Hierarchical generation of molecular graphs using structural motifs”. In: *International Conference on Machine Learning*, pp. 4839–4848.
- Al-Jumaily, Aws et al. (2023). “Examining multi-objective deep reinforcement learning frameworks for molecular design”. In: *Biosystems* 232, p. 104989.
- Li, Yanjun et al. (2019). “DeepAtom: A framework for protein-ligand binding affinity prediction”. In: *2019 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*. IEEE, pp. 303–310.
- Lin, Yu-Hsuan, Yueh-Chien Lin, and Chien-Chin Chen (2021). “Lysophosphatidic acid receptor antagonists and cancer: the current trends, clinical implications, and trials”. In: *Cells* 10.7, p. 1629.

- Liu, Bo et al. (2021). “Conflict-averse gradient descent for multi-task learning”. In: *Advances in Neural Information Processing Systems* 34, pp. 18878–18890.
- Liu, Chunming, Xin Xu, and Dewen Hu (2015). “Multiobjective reinforcement learning: A comprehensive overview”. In: *IEEE Transactions on Systems, Man, and Cybernetics: Systems* 45.3, pp. 385–398.
- Liu, Qi et al. (2018). “Constrained graph variational autoencoders for molecule design”. In: *Advances in neural information processing systems* 31.
- Liu, Xuhan et al. (2023). “DrugEx v3: scaffold-constrained drug design with graph transformer-based reinforcement learning”. In: *Journal of Cheminformatics* 15.1, p. 24.
- Liu, Yinhan et al. (2019). *RoBERTa: A Robustly Optimized BERT Pretraining Approach*. arXiv: 1907.11692 [cs.CL].
- Loeffler, Hannes H et al. (2024). “Reinvent 4: Modern AI-driven generative molecule design”. In: *Journal of Cheminformatics* 16.1, p. 20.
- Mukaidaisi, Muhetaer et al. (2022). “Multi-objective drug design based on graph-fragment molecular representation and deep evolutionary learning”. In: *Frontiers in Pharmacology* 13, p. 920747.
- Nguyen, Thanh Thi et al. (2020). “A multi-objective deep reinforcement learning framework”. In: *Engineering Applications of Artificial Intelligence* 96, p. 103915.
- Pereira, Tiago et al. (2021). “Diversity oriented Deep Reinforcement Learning for targeted molecule generation”. In: *Journal of Cheminformatics* 13.1, pp. 1–17.
- Popova, Mariya, Olexandr Isayev, and Alexander Tropsha (2018). “Deep reinforcement learning for de novo drug design”. In: *Science Advances* 4.7, eaap7885.
- Sattarov, Boris et al. (2019). “De novo molecular design by combining deep autoencoder recurrent neural networks with generative topographic mapping”. In: *Journal of chemical information and modeling* 59.3, pp. 1182–1196.
- Schneuing, Arne et al. (2022). “Structure-based drug design with equivariant diffusion models”. In: *arXiv preprint arXiv:2210.13695*.
- Senushkin, Dmitry et al. (2023). “Independent component alignment for multi-task learning”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 20083–20093.
- Shreyashree, S. et al. (2022). “A Literature Review on Bidirectional Encoder Representations from Transformers”. In: *Inventive Computation and Information Technologies*. Ed. by S. Smys, Valentina Emilia Balas, and Ram Palanisamy. Singapore: Springer Nature Singapore, pp. 305–320.
- Sorokina, Maria et al. (2021). “COCONUT online: Collection of open natural products database”. In: *Journal of Cheminformatics* 13.1, p. 2.
- Spiegel, Jacob O and Jacob D Durrant (2020). “AutoGrow4: an open-source genetic algorithm for de novo drug design and lead optimization”. In: *Journal of cheminformatics* 12, pp. 1–16.
- Ståhl, Niclas et al. (2019). “Deep Reinforcement Learning for Multiparameter Optimization in de novo Drug Design”. In: *Journal of Chemical Information and Modeling* 59.7, pp. 3166–3176.
- Sun, Duxin et al. (2022). “Why 90% of clinical drug development fails and how to improve it?” In: *Acta Pharmaceutica Sinica B* 12.7, pp. 3049–3062.
- Sutton, Richard S and Andrew G Barto (2018). *Reinforcement Learning: An Introduction*. MIT press.
- Tang, Huidong et al. (2023). “EarlGAN: An enhanced actor-critic reinforcement learning agent-driven GAN for de novo drug design”. In: *Pattern Recognition Letters* 175, pp. 45–51.
- Tang, Shidi et al. (2024). “Vina-GPU 2.1: Towards Further Optimizing Docking Speed and Precision of AutoDock Vina and Its Derivatives”. In: *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, pp. 1–13.
- Thafar, Maha et al. (Mar. 2022). “Affinity2Vec: Drug-target binding affinity prediction through representation learning, graph mining, and machine learning”. In: *Scientific Reports* 12.
- Vaswani, Ashish et al. (2017). “Attention is all you need”. In: *Neural Information Processing Systems*.
- Wang, Jing and Fei Zhu (2024). “ExSelfRL: An exploration-inspired self-supervised reinforcement learning approach to molecular generation”. In: *Expert Systems with Applications* 260, p. 125410. ISSN: 0957-4174.
- Wettig, Alexander et al. (May 2023). “Should You Mask 15% in Masked Language Modeling?” In: *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*. Ed. by Andreas Vlachos and Isabelle Augenstein. Dubrovnik, Croatia: Association for Computational Linguistics, pp. 2985–3000.
- Wishart, David S. et al. (2018). “DrugBank 5.0: A major update to the DrugBank database for 2018”. In: *Nucleic Acids Research* 46.Database, pp. D1074–D1082.

- Yang, Lijuan et al. (2021a). “Transformer-based generative model accelerating the development of novel BRAF inhibitors”. In: *ACS omega* 6.49, pp. 33864–33873.
- Yang, Soojung et al. (2021b). “Hit and lead discovery with explorative RL and fragment-based molecule generation”. In: *Advances in Neural Information Processing Systems* 34.
- Yu, Tianhe et al. (2020). “Gradient surgery for multi-task learning”. In: *Advances in Neural Information Processing Systems* 33, pp. 5824–5836.
- Zhang, Xiao-Chen et al. (2022). “Pushing the Boundaries of Molecular Property Prediction for Drug Discovery with Multitask Learning BERT Enhanced by SMILES Enumeration”. In: *Research 2022*, p. 0004.
- Zhou, Zhenpeng et al. (2019). “Optimization of molecules via deep reinforcement learning”. In: *Scientific Reports* 9.1, pp. 1–10.
- Zhu, Yiheng et al. (2024). “Sample-efficient multi-objective molecular optimization with gflownets”. In: *Advances in Neural Information Processing Systems* 36.

A Appendix

A.1 Transformer-Encoder

Transformer-encoder models are made for pretraining on unlabeled data in a bidirectional fashion (Vaswani et al., 2017; Devlin et al., 2019; Shreyashree et al., 2022). To extract features, an embedding layer transforms the input fragment tokens $x = (x_1, x_2, \dots, x_n)$ into learnable embedding vectors $w = (w_1, w_2, \dots, w_n)$, with the addition of a sinusoidal positional encoding vectors to reflect sequential location information. This is done using an embedding dictionary $\mathbf{D} \in \mathbb{R}^{V \times F}$, where $w_i \in \mathbb{R}^F$, V is the vocabulary size, and F is the embedding vector size. As an input feature matrix $\mathbf{Y} \in \mathbb{R}^{N \times F}$ is passed through the multi-head self-attention layer, it is linearly transformed into the following $h = 1, 2, \dots, H$ matrices: (1) the query matrix $\mathbf{Q}_h = \mathbf{Y}\mathbf{W}_h^Q$, (2) the key matrix $\mathbf{K}_h = \mathbf{Y}\mathbf{W}_h^K$, and (3) the value matrix $\mathbf{V}_h = \mathbf{Y}\mathbf{W}_h^V$, where \mathbf{W}_h^Q , \mathbf{W}_h^K , and \mathbf{W}_h^V are model weight matrices. The scaled dot-product attention is then computed for each linear projection, producing the output for a single attention head: $\mathbf{O}_h = \text{softmax}\left(\frac{\mathbf{Q}_h\mathbf{K}_h^T}{\sqrt{d_k}}\right)\mathbf{V}_h$, where $\sqrt{d_k}$ is a scaling factor. To get the final attention output, all attention heads $\mathbf{O}_1, \mathbf{O}_2, \dots, \mathbf{O}_H$ are then concatenated and fed into a linear layer. Finally, during pretraining, the attention output is processed by a feed-forward network, referred to as the ‘‘pretraining head.’’ This head is typically replaced with task-specific head during the fine-tuning stage (Zhang et al., 2022).

The Masked Language Model (MLM) task, a denoising-based auto-encoding technique, is often used to pretrain encoder models (Devlin et al., 2019). The goal is to reconstruct a noisy token sequence, where some tokens are masked, back to its original form. The model achieves this by using the surrounding visible tokens to build context for predicting the masked tokens (Zhang et al., 2022). More formally, given an input token sequence x , a noisy version \hat{x} is generated by masking a percentage m of its tokens (Wettig et al., 2023). The model’s task is to predict on the masked token set \mathcal{M} of \hat{x} to recover x :

$$L(\mathcal{C}) = \mathbb{E}_{x \in \mathcal{C}} \mathbb{E}_{\substack{\mathcal{M} \subset x \\ |\mathcal{M}|=m|x}} \left[\sum_{x_i \in \mathcal{M}} \log p(x_i | \hat{x}) \right]. \tag{7}$$

A.2 Reinforcement Learning

The MDP is defined by the tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$, where \mathcal{S} and \mathcal{A} represent the state and action spaces, \mathcal{P} is the state transition probability distribution $\mathcal{P}(s_{t+1}|s_t, a_t)$, \mathcal{R} is the reward distribution $\mathcal{R}(r_t|s_t, a_t)$, and γ is the discount factor used to control the trade-off between immediate rewards and future rewards, where t is the current timestep, and r_t is a scalar reward function for t (Al-Jumaily et al., 2023; Graesser and Keng, 2019). The goal of an RL agent is to learn a policy distribution $\pi(a_t|s_t)$ that maximizes long-term cumulative rewards through exploration of the environment over multiple timesteps. This is accomplished by the agent starting at state s_t , selecting action a_t , receiving reward r_t , and transitioning to the new state s_{t+1} (Sutton and Barto, 2018). To assess the value of states and actions with respect to expected long-term returns, two functions are formulated: $V^\pi(s)$, which measures the desirability of s : $V^\pi(s) = \mathbb{E}_{s_0=s, \tau \sim \pi} \left[\sum_{t=0}^T \gamma^t r_t \right]$, and $Q^\pi(s, a)$, which measure the desirability of taking action a given state s : $Q^\pi(s, a) = \mathbb{E}_{s_0=s, a_0=a, \tau \sim \pi} \left[\sum_{t=0}^T \gamma^t r_t \right]$ (Graesser and Keng, 2019).

When the RL task entails exploring a vast state or action space, as is often the case in drug design, learning an exact optimal policy or value function can become computationally intractable. Thus, DRL is used to approximate policies or value functions (Arulkumaran et al., 2017). The actor-critic framework approximates both and has been leveraged by various drug development frameworks (Al-Jumaily et al., 2023; Goel et al., 2021; Gottipati et al., 2021; Pereira et al., 2021; Popova, Isayev, and Tropsha, 2018; Ståhl et al., 2019; Tang et al., 2023; Wang and Zhu, 2024; Yang et al., 2021b). The actor model is responsible for learning a parameterized policy π_{θ_A} , guided by feedback, known as temporal difference (TD) error from the critic model, which evaluates the actor’s actions based on the state. One approach to this evaluation is by learning the advantage function $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$, which measures the desirability of taking action a compared to alternative actions available from state s (Graesser and Keng, 2019). However, having the critic

model learn both $Q^\pi(s, a)$ and $V^\pi(s)$ is computationally expensive. Therefore, in practice, the critic model only learns $V^\pi(s)$ and combines it with reward information from the trajectory to estimate the advantage function:

$$\begin{aligned} A^\pi(s_t, a_t) &= Q^\pi(s_t, a_t) - V^\pi(s_t) \\ &\approx r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^n r_{t+n} + \gamma^{n+1} \hat{V}^\pi(s_{t+n+1}) - \hat{V}^\pi(s_t). \end{aligned} \quad (8)$$

Thus, the value function is parameterized as $V_{\theta_C}^\pi(s)$ and is updated using loss function:

$$L_{\text{val}}(\theta_C) = \frac{1}{T} \sum_{t=0}^T \left(r_t + \hat{V}_{\theta_C}^\pi(s_{t+1}) - V_{\theta_C}^\pi(s_t) \right)^2, \quad (9)$$

while the loss function for the actor is given by:

$$L_{\text{pol}}(\theta_A) = \frac{1}{T} \sum_{t=0}^T \left(-\hat{A}^\pi(s_t, a_t) \log \pi_{\theta_A}(a_t | s_t) \right). \quad (10)$$

A.3 AMTL-based MODRL Algorithm

For the MODRL training, we aim to use the gradient modulation method AMTL (Senushkin et al., 2023) for policy learning. AMTL specifically addresses the multi-task optimization challenges, i.e., gradient dominance and gradient conflicts, by aligning principal components of a gradient matrix. The existence of conflicting or dominating gradients disrupts the stability of the training process and leads to a deterioration in overall performance.

It is acknowledged that the gradient dominance can be measured with a gradient magnitude similarity (Yu et al., 2020), and a cosine distance between vectors can measure the gradient conflicts (Liu et al., 2021). However, the two metrics cannot offer a comprehensive assessment if taken in isolation. One of the key components of AMTL is the proposal of the condition number, a stability criterion that can indicate the presence of both challenges. The value of the condition number is the ratio of the maximum and minimum singular values of the corresponding matrix. Minimizing the condition number of the linear system of gradients, a linear combination of gradients for all objectives, mitigates dominance and conflicts within this system. If we apply singular value decomposition (SVD), we can have

$$\mathbf{G} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T, \quad (11)$$

where $\mathbf{\Sigma} = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_K)$ and the eigen-values are arranged in decreasing order. One can easily obtain that

$$\mathbf{G}^T \mathbf{G} = \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \mathbf{V}\mathbf{\Sigma}\mathbf{\Sigma}\mathbf{V}^T = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T, \quad (12)$$

where $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_K)$ and we know that $\sigma_k = \sqrt{\lambda_k}$. Thus, the singular values in the SVD of \mathbf{G} correspond to the squared roots of the eigen-values from the eigen-decomposition of the Gram matrix $\mathbf{G}^T \mathbf{G}$. According to AMTL, a gradient matrix with a minimal condition number (i.e., the singular values are equal to the last positive singular value) can be decomposed as:

$$\hat{\mathbf{G}} = \mathbf{U}\hat{\mathbf{\Sigma}}\mathbf{V}^T = \mathbf{U}\sigma\mathbf{I}\mathbf{V}^T = \sigma\mathbf{U}\mathbf{V}^T = \sigma\mathbf{G}\mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{V}^T, \quad (13)$$

where $\sigma = \sqrt{\lambda_K}$ and $\mathbf{U} = \mathbf{G}\mathbf{V}\mathbf{\Sigma}^{-1}$ because of Equation (11), and $\hat{\mathbf{G}}$ is the aligned gradient matrix. A linear combination of the aligned objective-specific gradient vectors using the objective importance would be $\hat{\mathbf{G}}\boldsymbol{\omega} = \sum_{k=1}^K \omega_k \hat{\mathbf{g}}_k$. The gist of AMTL is to align the gradient matrix by conducting an SVD to the original gradient matrix and rescaling the singular values to match the smallest singular value. The pseudocode for the MODRL fine-tuning algorithm proposed in this work to align the language model is given in Algorithm 1.

Algorithm 1: Multi-Objective Deep Reinforcement Learning (MODRL) Pseudocode

Require: π_0 : original policy; K : number of objectives; ω : task importance (all objectives are deemed equal importance in this work); η : learning rate;

```
1 Let  $\pi_\phi = \pi_0$ ;  
2 foreach epoch do  
3   foreach minibatch do  
4     foreach  $k = 1, 2, \dots, K$  do  
5       Compute loss  $\mathcal{L}_k(\phi)$ ;  
6       Compute gradient  $\mathbf{g}_k = \nabla_\phi \mathcal{L}_k(\phi)$ ;  
7     end  
8     Get the gradient matrix  $\mathbf{G} = \{\mathbf{g}_1, \dots, \mathbf{g}_K\}$ ; // playing objective-specific  
       gradient vectors as columns in  $\mathbf{G}$   
9     Compute task space Gram matrix  $\mathbf{M} \leftarrow \mathbf{G}^T \mathbf{G}$ ;  
10    Get eigen-values and eigen-vectors  $(\boldsymbol{\lambda}, \mathbf{V}) \leftarrow \text{eigen}(\mathbf{M})$ ; // eigen-decomposition  
       such that  $\mathbf{M} = \mathbf{V} \boldsymbol{\Lambda} \mathbf{V}^T$  where  $\boldsymbol{\Lambda} = \text{diag}(\boldsymbol{\lambda})$   
11     $\boldsymbol{\Sigma}^{-1} \leftarrow \text{diag}\left(\sqrt{\frac{1}{\lambda_1}}, \dots, \sqrt{\frac{1}{\lambda_K}}\right)$ ;  
12    Balance transformation  $\mathbf{B} \leftarrow \sqrt{\lambda_\eta} \mathbf{V} \boldsymbol{\Sigma}^{-1} \mathbf{V}^T$ ;  
13    Get new aligned gradient matrix  $\hat{\mathbf{G}} = \mathbf{G} \mathbf{B}$ ; Updated gradient  $\nabla \phi = \hat{\mathbf{G}} \omega$ ;  
14    Update policy parameter  $\phi = \phi - \eta \nabla \phi$ ;  
15  end  
16 end  
17 Return policy  $\pi_\phi$ ;
```

A.4 Fragments-SMILES Hybrid Tokenization Strategy Figures

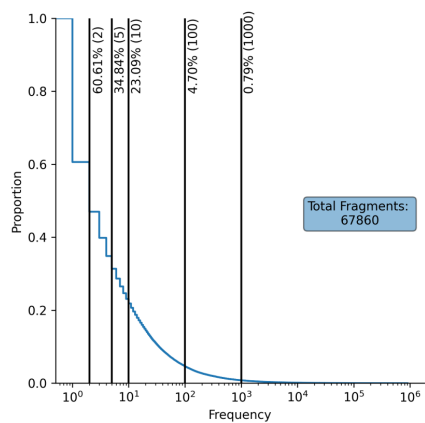


Figure 5: TDC MolGen task dataset (Huang et al., 2021) fragment frequencies.

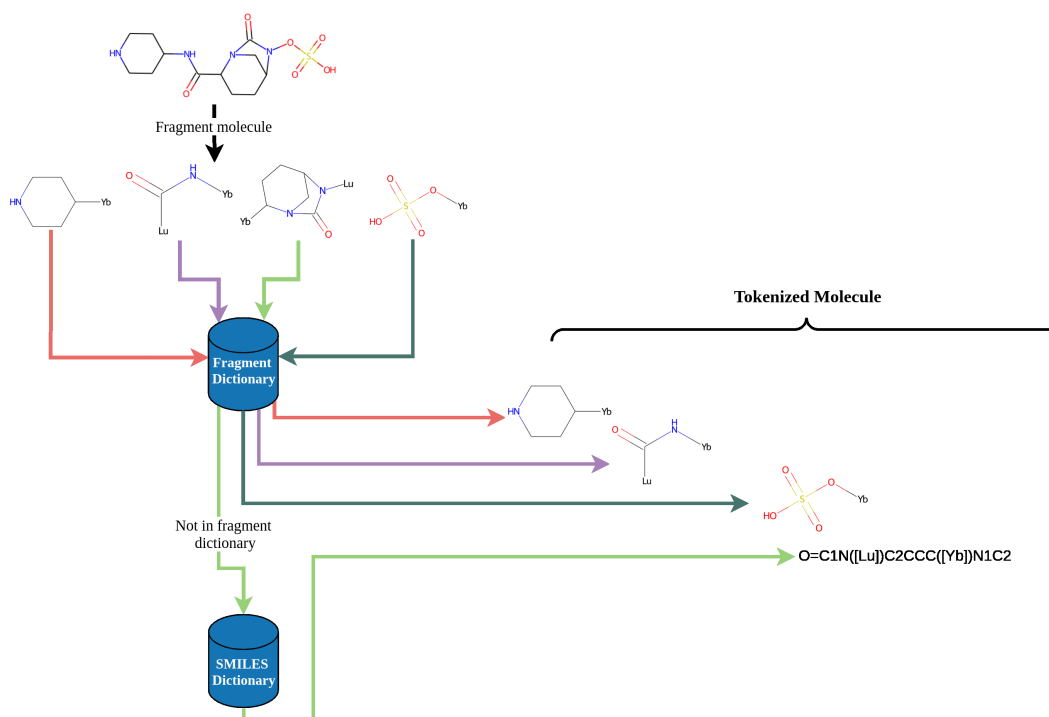


Figure 6: Fragment-SMILES hybrid tokenization strategy.

A.5 Pretraining Diagrams

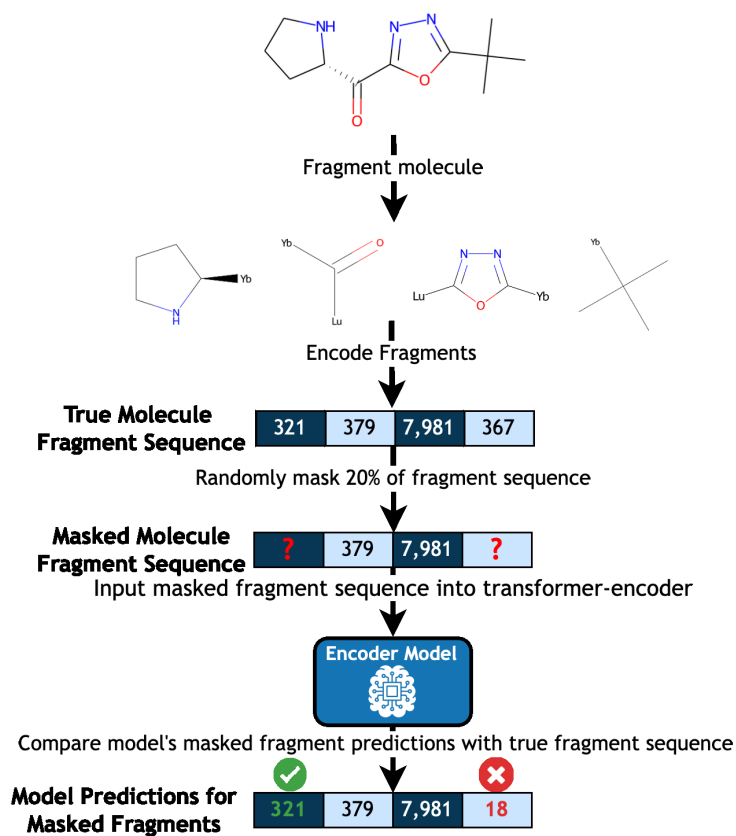


Figure 7: MLM training process for one molecule.

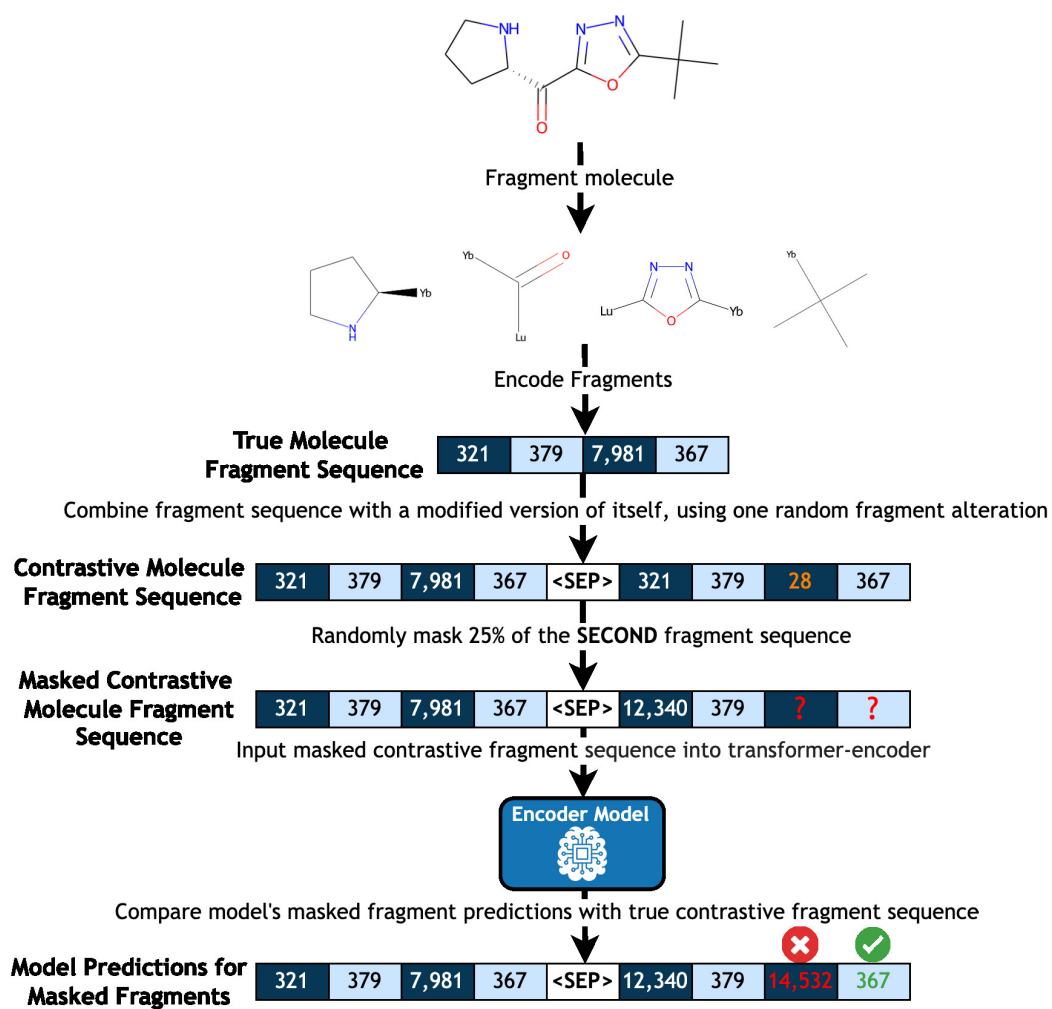


Figure 8: Contrastive learning training process for one molecule.

A.6 Pretraining Results

Table 5: Pretraining results using the 2-Frequency token dictionary on a 4-million molecule dataset from TDC (Huang et al., 2021).

Epoch	Testing Loss	Testing Accuracy	Compute Time
Epoch 1: MLM	1.02	0.72	~8 hrs
Epoch 2: MLM	0.93	0.75	~8 hrs
Epoch 3: Contrastive Learning	0.37	0.90	~20 hrs
Epoch 4: MLM	0.87	0.78	~8 hrs
Epoch 5: Contrastive Learning	1.70	0.91	~20 hrs
Epoch 6: MLM	0.87	0.79	~7 hrs

Table 6: Pretraining results using the 100-Frequency token dictionary on a 4-million molecule dataset from TDC (Huang et al., 2021).

Epoch	Testing Loss	Testing Accuracy	Compute Time
Epoch 1: MLM	0.90	0.70	~3 hrs
Epoch 2: MLM	0.83	0.72	~3 hrs
Epoch 3: Contrastive Learning	0.29	0.89	~14 hrs
Epoch 4: MLM	0.80	0.73	~3 hrs
Epoch 5: Contrastive Learning	0.26	0.90	~15 hrs
Epoch 6: MLM	0.77	0.74	~3 hrs

Table 7: Pretraining results using the 1000-Frequency token dictionary on a 4-million molecule dataset from TDC (Huang et al., 2021).

Epoch	Testing Loss	Testing Accuracy	Compute Time
Epoch 1: MLM	0.88	0.70	~3 hrs
Epoch 2: MLM	0.81	0.72	~3 hrs
Epoch 3: Contrastive Learning	0.24	0.90	~15 hrs
Epoch 4: MLM	0.78	0.73	~3 hrs
Epoch 5: Contrastive Learning	0.22	0.91	~15 hrs
Epoch 6: MLM	0.76	0.74	~3 hrs

A.7 Property-Wise Density Plots For the Comparative and Scalability Studies

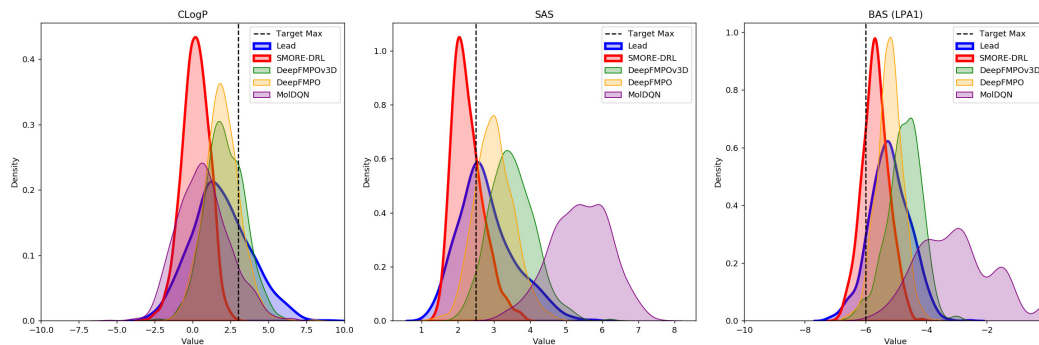


Figure 9: Property-wise comparisons between the lead molecules (blue) and the molecules optimized in final epoch by SMORE-DRL (red), DeepFMPOv3D (Bolcato, Heid, and Boström, 2022) (green), DeepFMPO (Stähl et al., 2019) (yellow), and MoIDQN (Zhou et al., 2019) (purple). All objectives are to be minimized and the targeted maximums are indicated by the black dashed line.

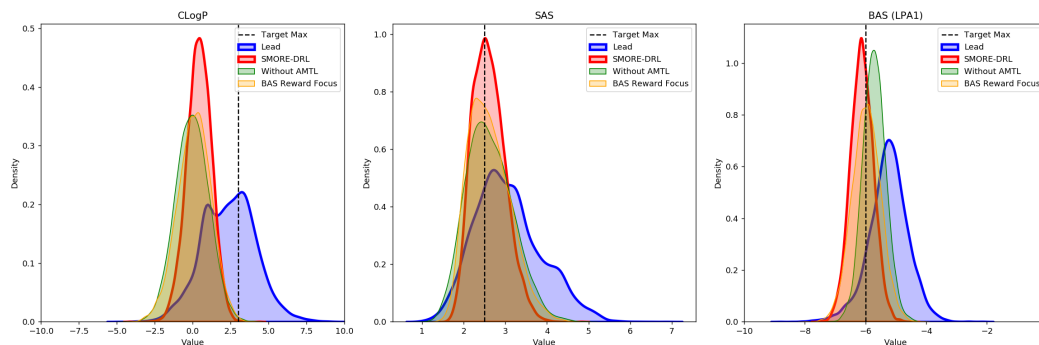


Figure 10: Property-wise comparisons between the lead molecules (blue) and the molecules optimized in final epoch by SMORE-DRL (red), SMORE-DRL without AMTL (green), and SMORE-DRL with a reward emphasis on BAS (yellow). All objectives are to be minimized and the targeted maximums are indicated by the black dashed line.

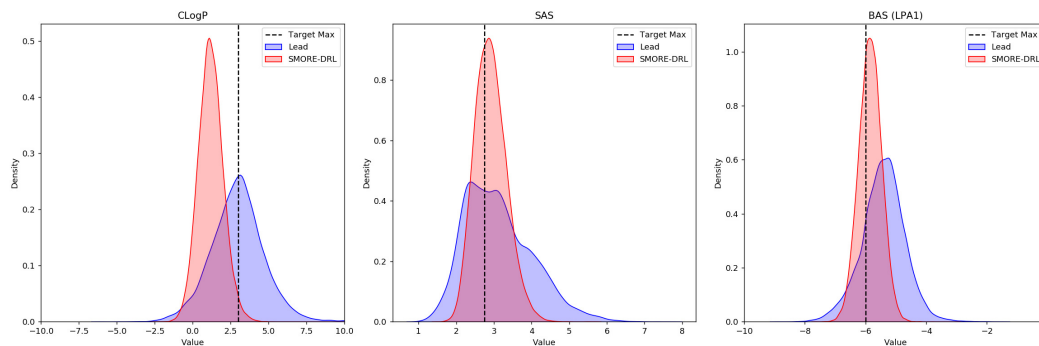


Figure 11: Generalization results – property-wise comparisons between the test lead molecules (blue) and the molecules optimized by SMORE-DRL (red). All objectives are to be minimized and the targeted maximums are indicated by the black dashed line.

A.8 Visualizations of SMORE-DRL's Molecular Optimization During Scalability Experiments

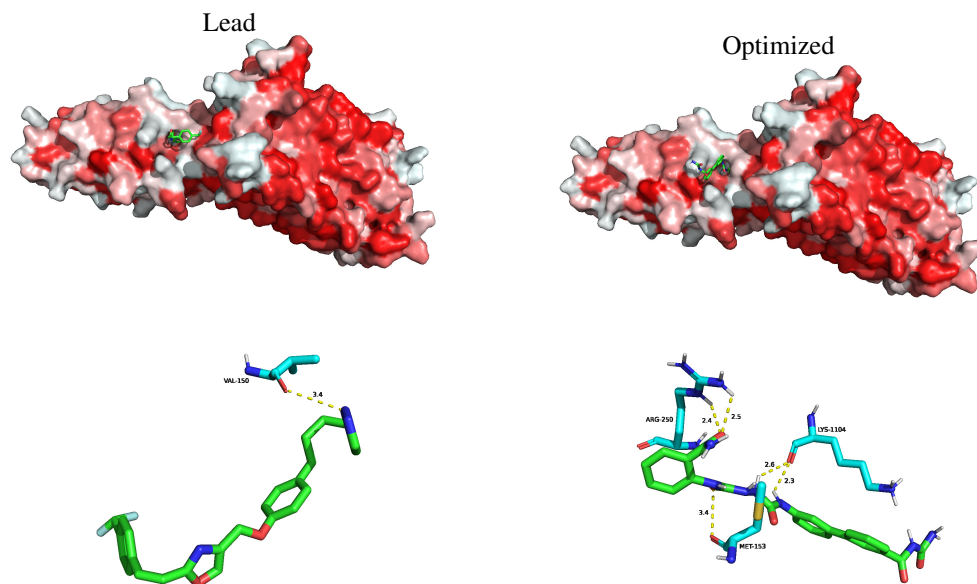


Figure 12: Binding visualization of a lead molecule (FC(F)(F)c1ccc(C=Cc2nc(C0c3ccc(CCCn4ccnn4)cc3)co2)cc1, ClogP = 6.05, SAS = 2.73, BAS = -4.7) and SMORE-DRL's optimized version (NC(=O)NC(=O)c1ccc(-c2cccc(NC(=O)NNC(=O)Nc3ccccc3C(N)=O)c2)c1, ClogP = 2.11, SAS = 2.31, BAS = -9.0) from the scalability experiments.

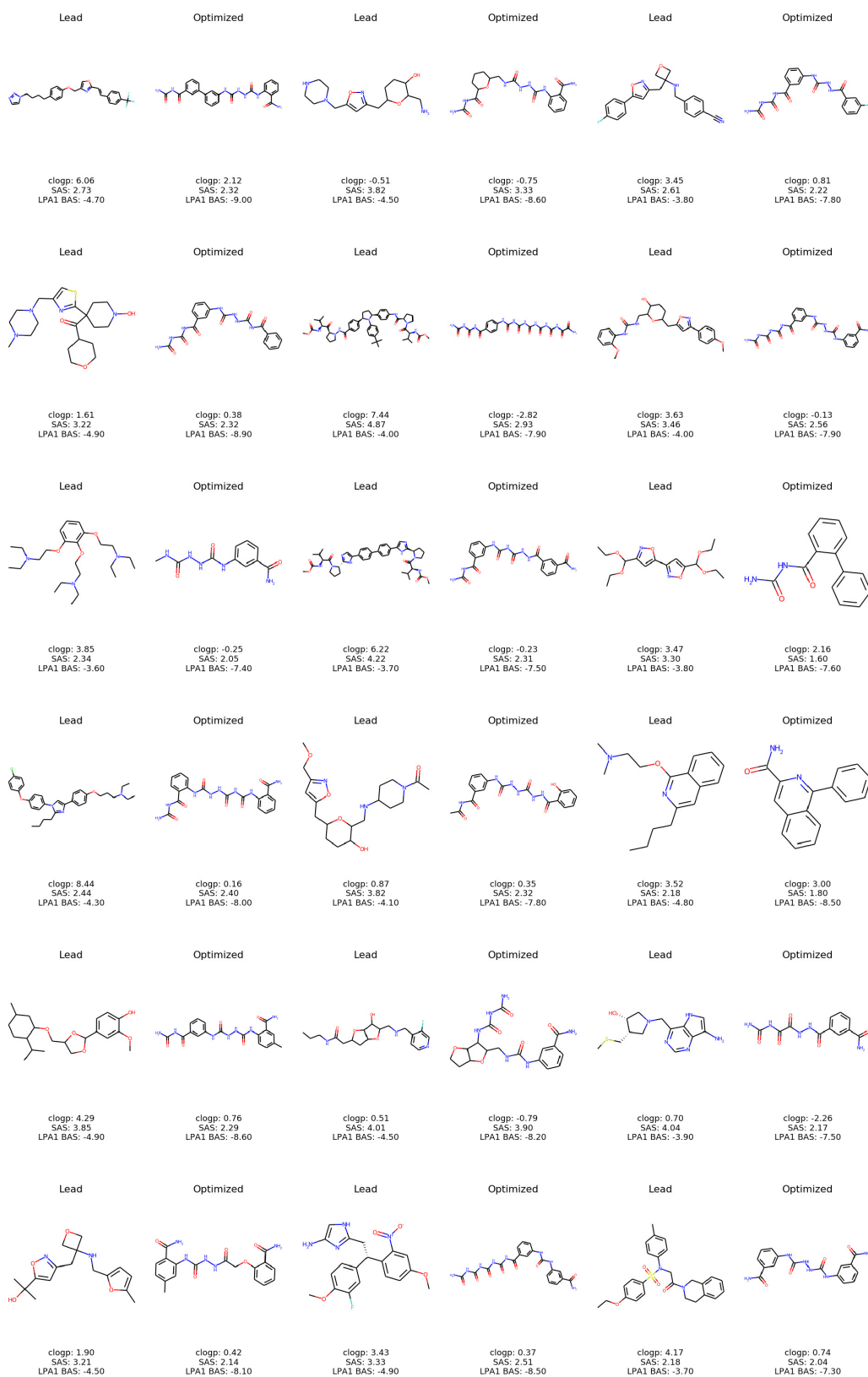


Figure 13: Lead molecules optimized by SMORE-DRL from the scalability experiments.

A.9 Visualizations of SMORE-DRL's Molecular Optimization During Generalization Experiments

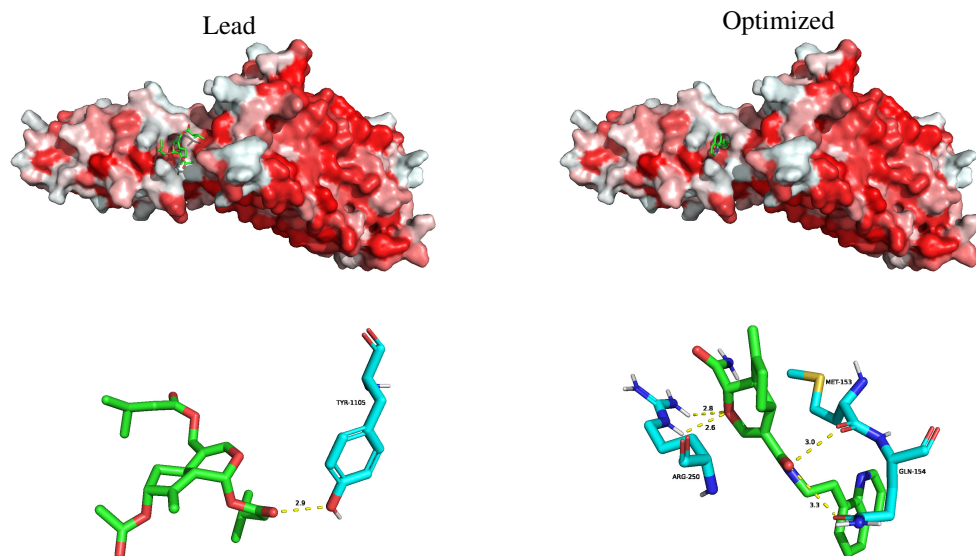


Figure 14: Binding visualization of a lead molecule (C=C1C(OC(C)=O)CC2C(COC(=O)CC(C)C)=COC(OC(=O)CC(C)C)C12, ClogP = 3.52, SAS = 4.33, BAS = -3.8) and SMORE-DRL's optimized version (C=C1CCC2C(C(=O)NCCc3cccc4ccnc34)=COC(C(N)=O)C12, ClogP = 2.24, SAS = 3.88, BAS = -8.2) from the generalization experiments.

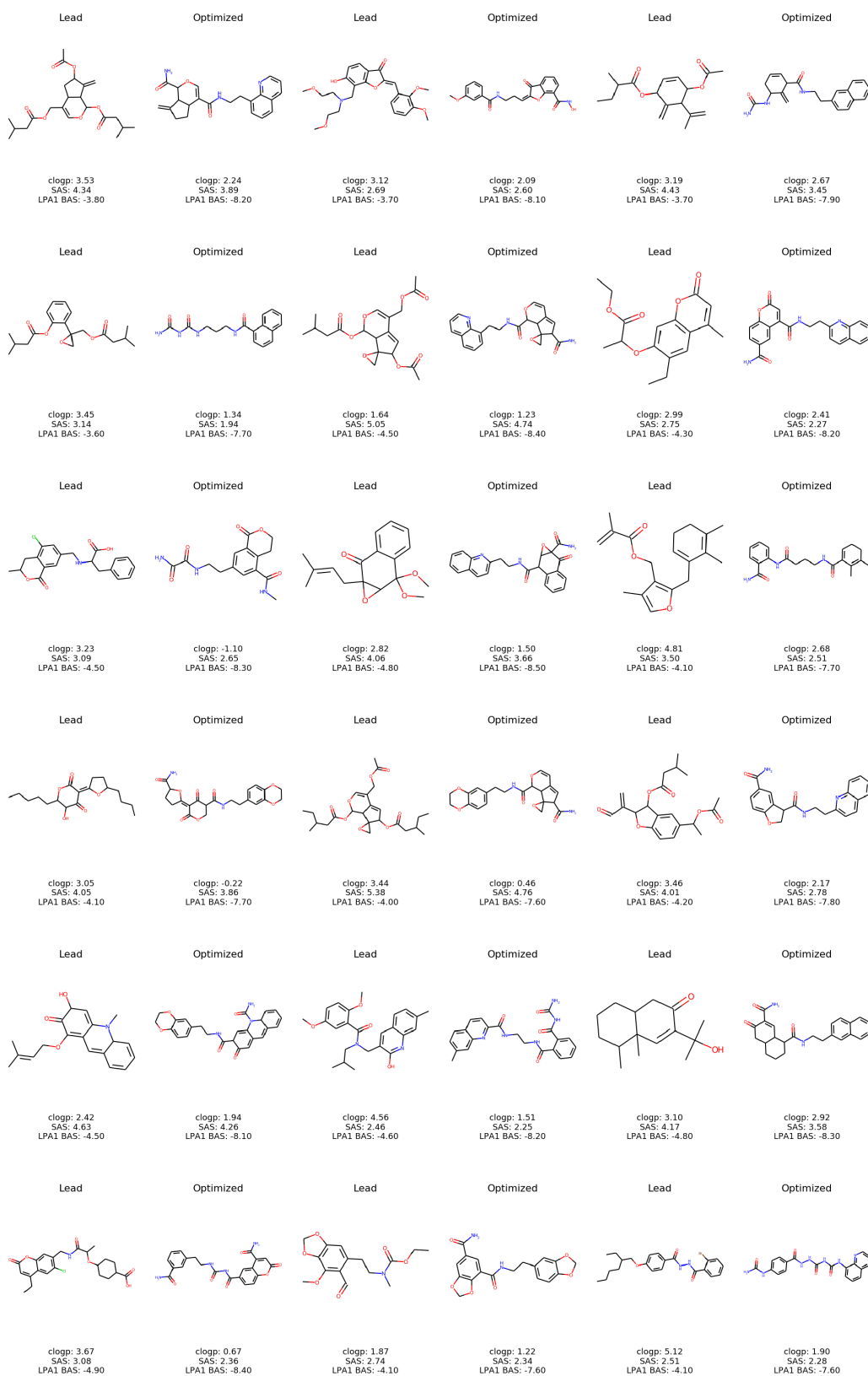


Figure 15: Lead molecules optimized by SMORE-DRL from the generalization experiments.