SMORE-DRL: SCALABLE MULTI-OBJECTIVE RO BUST AND EFFICIENT DEEP REINFORCEMENT LEARNING FOR MOLECULAR OPTIMIZATION

Anonymous authors

Paper under double-blind review

Abstract

The adoption of machine learning 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 achieved through the application of deep reinforcement learning (DRL) frameworks. 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. Finetuning 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 robustness. 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.

005

008 009 010

011

013

014

015

016

017

018

019

021

022

025

026

027

028

029

031

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 10²⁰ and 10²⁰⁰
possible drug-like molecules (Brown, 2015).

044 Machine learning methods, including deep reinforcement learning (DRL), offer a promising solution to this problem (Al-Jumaily et al., 2023; Goel et al., 2021; Gottipati et al., 2021; 046 Kim et al., 2021; Pereira et al., 2021; Popova et al., 2018; Ståhl et al., 2019; Tang et al., 2023; Wang & Zhu, 2024; Yang et al., 2021). Molecular optimization is a process that 048 requires a model to perform minor modifications on a lead molecule to improve its drug-like 049 qualities while preserving structural similarity. As molecules with comparable structures are anticipated to exhibit similar behaviours, this approach aims at preventing the model 051 from producing unrealistic or undesirable molecules. This differs from molecular generation, where a model's task is to generate novel and diverse compounds from scratch (Ståhl et al., 052 2019). Additionally, present methodologies are hindered by lengthy training requirements and sub-optimal use of training data, resulting in impaired scalability and performance.

054 In this work, we present SMORE-DRL (Scalable Multi-Objective Robust and Efficient 055 Deep Reinforcement Learning), a gradient-alignment-based multi-objective DRL (MODRL) framework for molecular optimization. The key contributions of this research include: (1) 057 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 059 tokenization strategy is proposed to represent the a molecule as a hybrid of fragments and 060 SMILES, enabling efficient policy learning and effective representation of any new molecules; 061 and (3) a synergistic integration of gradient alignment, hybrid fragment-SMILES representa-062 tion, contrastive learning, and a transformer-encoder allows for scalability and generalization 063 capability superior to existing MODRL methods. Moreover, SMORE-DRL demonstrated 064 its ability to effectively scale and generalize its optimization process to new molecules after 065 fine-tuning. This is a particularly notable aspect of our work, as existing DRL methods lack 066 scalability and their generalization capacities are under-explored in current literature.

067 068 069

071

072

2 Related Work

The fundamental techniques of MODRL and aligned multi-task learning, closely related to this study, are reviewed below. In addition, see Appendix A.1 for a review of transformerencoder architectures and MLM.

073 074 075

2.1 Deep Reinforcement Learning

076 Reinforcement learning (RL) agents learn through a trial-and-error process guided by the 077 Markov decision process (MDP). See Appendix A.2 for an introduction to basic RL concepts. 078 When the RL task entails exploring a vast state or action space, as is often the case in drug 079 design, learning an exact optimal policy or value function can become computationally intractable. Thus, DRL is used to approximate policies or value functions (Arulkumaran 081 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 082 et al., 2021; Pereira et al., 2021; Popova et al., 2018; Ståhl et al., 2019; Tang et al., 2023; 083 Wang & Zhu, 2024; Yang et al., 2021). The actor model is responsible for learning a parameterized policy π_{θ_A} . This is guided by feedback known as temporal difference (TD) 085 error from the critic model, which evaluates the actor's actions based on the state. One 086 approach to this is by learning the advantage function $A^{\pi}(s,a) = Q^{\pi}(s,a) - V^{\pi}(s)$, which 087 measures the desirability of taking action a compared to alternative actions available from 088 state s (Graesser & Keng, 2019).

089 090

091

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 et al., 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). MODRL approaches have been developed for molecular design. Zhou et al. (2019) developed Molecule Deep Q-Networks (MolDQN), a multi-objective molecular 096 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 098 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 100 is used for MOO. Deep Fragment-based Multi-Parameter Optimization (DeepFMPO), in-101 troduced by Ståhl et al. (2019), is an actor-critic multi-objective method for molecular 102 optimization. In this work, a library of fragments is derived from a set of lead molecules 103 and fragments are encoded using a balanced binary tree such that similar molecules have 104 similar binary encoding. One modification step involves replacing a fragment in the lead 105 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 106 achieved or a reward of zero. If all objectives are met, the reward is doubled. A dynamic 107 reward mechanism is also implemented, where the model is penalized if it begins to underperform compared to previous epochs. Bolcato et al. (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. When recognizing the existence of other RL approaches
for molecular optimization, the three aforementioned MODRL approaches representatively
form the benchmarks to compare with our proposed framework.

115 116

2.3 Aligned Multi-Task Learning

117 Two potential issues that arise when solving an MOO problem directly using gradient descent 118 are dominating and conflicting gradients. A dominating objective gradient is characterized 119 by the largest magnitude, which leads to a bias in the solution favouring the corresponding 120 task (Senushkin et al., 2023). When two objective gradients are conflicting, an increase 121 in the solution towards one objective decreases the solution for the conflicting objective. 122 Conflicting gradients are characterized by having a negative cosine similarity (Yu et al., 123 2020). To address these challenges in the context of multi-task learning, Senushkin et al. 124 (2023) propose aligned-multi-task learning (AMTL). Let $\mathcal{L}_k(\boldsymbol{\theta})$ represent the objective of 125 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: 126

- 127
- 128
- 129 130

140

141 142

143

144 145

146

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

To mitigate conflicting and dominating gradients, AMTL aligns the principal components 131 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 133 minimal singular value of the gradient matrices is identified and all other singular values 134 are adjusted to match it, resulting in aligned gradients. Subsequently, these aligned gra-135 dients are combined into a common gradient (Senushkin et al., 2023). Inspired by AMTL 136 which formulates a multi-task learning problem to an MOO problem, we integrate the same 137 gradient alignment technique to solve MODRL for drug design in this study. A detailed 138 description of AMTL-based MODRL can be found in Appendix A.3. 139

3 Methods

In this section, we discuss the data pre-processing, pretraining, and fine-tuning processes 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 et al. (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 et al., 2022).

154 While fragmentation is a good technique for reducing the chemical search space, it may 155 result in a vast token dictionary size. Figure 4 in Appendix A.4 is a fragment frequency 156 chart based on the 4-million molecule dataset, which shows that nearly 68,000 fragments 157 were extracted. Most of these fragments are rarely encountered in the dataset, with 95%158 appearing fewer than 100 times. Moreover, building a token dictionary solely from the 159 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 160 not present in the dictionary, especially if the fine-tuning dataset differs from the one used 161 for pretraining.

To reduce the dictionary size while still representing fragments that are absent or infre-163 quently encountered, we propose a novel hybrid tokenization strategy that uses both frag-164 ments and SMILES. Following the construction of a fragment dictionary from the pretraining 165 dataset, we append tokens for SMILES and exclude all fragments that appear less than twice. 166 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 repre-167 sented atom by atom. See Appendix A.4 for a diagram of the hybrid tokenization strategy. 168 As part of our ablation studies in Section 4.2, we demonstrate that further reducing the 169 token dictionary by retaining only the most frequently encountered fragments (resulting in 170 molecules being primarily represented by SMILES atoms) hinders training performance. 171

172

174

173 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.

- 181 182 183
- 3.2.1 Masked Language Model (MLM)

184 Given a training batch of molecules, they are fragmented and encoded into their token 185 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 187 dynamic approach introduced by Liu et al. (2019) in Robustly Optimized BERT Pretraining 188 Approach (RoBERTa). In each training batch, 20% of the tokens are randomly selected for 189 masking, with 90% of those selected tokens end up being masked. If a selected token 190 is part of a sequence of atom tokens representing a fragment that does not exist in the 191 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 192 reconstruct the original values of the masked tokens. See Appendix A.5 for a diagram of 193 the MLM process. 194

195

196 3.2.2 Contrastive Learning

197 We further refine the SMORE-DRL's contextual understanding of molecules by allowing it 198 to align its representations of similar molecules. This is particularly valuable during fine-199 tuning, where the model is tasked with optimizing lead molecules over multiple timesteps 200 while ensuring that the optimized molecules in the earlier timesteps retain chemical sim-201 ilarity to the original lead molecule. To accomplish this, we introduce a straightforward 202 contrastive learning technique that builds upon the MLM approach. Rather than directly 203 masking tokens in the fragment sequence as previously described, a "separation" token is 204 inserted at the end of the sequence, followed by an augmented version of that sequence. 205 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 represent-206 ing a fragment that is not in the token dictionary, the entire SMILES atom sequence for 207 that fragment is replaced with a randomly selected fragment token. The masking process 208 consists of keeping the original molecule sequence fully visible to the model, while masking 209 25% of the augmented sequence using the same technique described earlier. See Appendix 210 A.5 for a diagram of the contrastive learning process. 211

212

214

213 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.

216 3.3.1 AGENTS

222 223 224

236 237 238

249 250 251

253

218 Masker Model: The masker model, denoted as π_{θ_M} , is responsible for selecting which 219 tokens to mask from the lead molecule token sequence. At each timestep, it masks at least 220 one token and up to 70% of the token sequence. The model is designed to prefer masking 221 fragment tokens over SMILES atom tokens. The loss for the masker model is:

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

225 Actor Model: After the lead molecule token sequence is masked by the masker model, it 226 is passed to the actor model, denoted as π_{θ_A} . The actor model utilizes the same training 227 head that was used during the pretraining phase. Hence, its task is to replace the masked 228 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 230 that the resulting sequence represents an optimized yet chemically similar version of the 231 lead molecule. The actor model employs AMTL (Senushkin et al., 2023) to identify a common objective gradient, thereby avoiding conflicting and dominating gradients, which 232 ensures that all molecular properties are optimized equally. This process involves obtaining 233 a gradient matrix that collects all K objective gradients, represented as $G = \{g_1, \dots, g_K\}$, 234 where $\boldsymbol{g}_{k} = \nabla L_{k} \left(\theta_{A} \right)$ and 235

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

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's desirability as a potential drug candidate. The loss of the critic model is:

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

3.3.2 Reward System

254 To assess a molecule's potential as a drug candidate, we use the following three properties: 255 (1) logarithm of partition coefficient (ClogP), which impacts a drug's administration, ab-256 sorption, transport and excretion, (2) synthetic accessibility score (SAS), which measures 257 the difficulty of synthesizing a molecule, and (3) binding affinity score (BAS) to LPA1, 258 which quantifies the binding capability of a drug to a target protein (Brown, 2015; Ertl & 259 Schuffenhauer, 2009; Li et al., 2019). However, the number and type of properties can be 260 tailored to any specific optimization task. RDKit is used for ClogP and SAS calculations, and QuickVina2-GPU-2.1 (Tang et al., 2024) was used to calculate BAS. Lysophosphatidic 261 acid receptor 1 (LPA1/LPAR1), a bioactive lipid mediator primarily derived from membrane 262 phospholipids, is chosen as the target protein for BAS. LPA Receptors (LPARs) have been 263 found to be over-expressed in multiple types of cancer, with LPA1 specifically expressed in 264 ovarian cancer, breast cancer, liver cancer, gastric cancer, pancreatic cancer, lung cancer, 265 glioblastoma and osteosarcoma. LPA1 promotes metastasis and tumor motility, making it a 266 natural choice for targeting in efforts to inhibit cancer spread and cell movement (Lin et al., 267 2021).268

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



the molecule, and in the final timestep, BAS is calculated and a full reward is given. 306 307 The following three criteria are also examined for each epoch: (1) validity: the ratio of 308 chemically valid optimized molecules, checked using RDKit, (2) novelty: the ratio of opti-309 mized 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 310 et al., 2022). In each timestep, if a molecule is valid, unique, and novel, it is assigned its 311

reward, else it is assigned a reward of -1 for each objective, for a total reward of -3 if it is 312 the final timestep. If all properties are achieved by a molecule, it is provided with extra 313 reinforcement by doubling the final reward. 314

Thus, the output of the reward function is $\mathbf{r}^{K} = [r^{1}, r^{2}, \dots, r^{K}]$, where K = 3 for the final 315 timestep and K = 2 for all intermediate timesteps. The values of the reward system are 316 listed in Table 1 of the next section. 317

319 4

318

320

- **EXPERIMENTS & RESULTS**
- 4.1 Pretraining 321
- The encoder model was pretrained for six consecutive epochs on a combination of MLM 323 and contrastive learning tasks, where the same training head was used throughout learning.

While the main experiment employs the 2-frequency token dictionary (68,000 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.

336 337

338

328 329

330

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

339 We compare SMORE-DRL's optimization performance with three other MODRL optimiza-340 tion frameworks: (1) DeepFMPOv3D, (2) DeepFMPO and (3) MolDQN. A primary goal 341 of this paper is to present the scalability of SMORE-DRL. However, it is not feasible to 342 conduct large-scale optimization using thousands of lead molecules to compare with the 343 other models, as these benchmarks lack the efficiency for scalability. As described in their 344 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 345 dataset was utilized. Scalability and generalizability of SMORE-DRL are demonstrated in 346 the following sections. Challenging property values were selected for the optimization task, 347 as noted in Table 1. The results of this comparative study represent the mean scores of 348 three separate runs for all models. 349

350 351

352

353

355 356 357 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)

358 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 70 359 epochs, during which all lead molecules were optimized over 4 timesteps per epoch. The 360 molecules optimized in the final timestep of the last epoch are used for our comparisons. 361 DeepFMPOv3D and DeepFMPO trained for 1,000 epochs, optimizing random batches of 362 512 unique molecules per epoch. DeepFMPOv3D performed optimization over 4 timesteps 363 and DeepFMPO used 8 timesteps. In the final epoch, all lead molecules were optimized, 364 and results from this epoch are used for our comparisons. MolDQN was trained for 6,000 365 epochs. For the first 5,000, a random molecule from the dataset is selected and optimized 366 for 20 timesteps. The final 1,000 epochs focus on optimizing each lead molecule, with the 367 resulting optimized molecules utilized for comparison. To calculate property rewards, all 368 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. 369 For MolDQN, a partial reward excluding BAS is assigned at each intermediate timestep, 370 while a full reward is assigned in the final timestep. 371

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 8 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, 379 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	-	$\sim 6 \text{ hrs}$	~4.5 hrs	$\sim 6 \text{ hrs}$	$\sim 5.5 \text{ hrs}$
Validity	-	$99.80\% (\pm 0.00)$	80.00% (±2.16)	$77.33\% (\pm 3.09)$	100% (±0.00)
Novelty	-	$98.52\% (\pm 0.05)$	80.01% (±1.85)	82.99% (± 0.27)	100% (±0.00)
Uniqueness	-	$93.63\% (\pm 1.47)$	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%	$71.83\%~(\pm 2.29)$	10.96% (±0.60)	22.86% (± 3.82)	$0.70\% (\pm 0.29)$
BAS (LPA1)	8.65%	30.85% (± 2.77)	$6.60\% (\pm 0.17)$	$7.86\% (\pm 1.43)$	$1.73\% (\pm 0.66)$
All Properties	0%	$23.54\%~(\pm 1.56)$	0.03% (±0.05)	0.92% (±0.42)	0% (±0.00)

409 410 411

412

413

378

380

381

394

395

396 397 398

399

400

4.2.2SCALABILITY OF SMORE-DRL

To examine scalability, SMORE-DRL was trained for 70 epochs using 10,000 molecules-5,000 414 from the DrugBank database (Wishart et al., 2018) and 5,000 from the Collection of Open 415 Natural Products (COCONUT) database (Sorokina et al., 2021). COCONUT molecules 416 were included to attempt to test the model's robustness, as they are different from those 417 typically encountered during pretraining. All lead molecules were optimized over 4 timesteps 418 per epoch, with those from the final timestep of the last epoch used for comparisons. While 419 the baseline model was run five times, ablation studies were also conducted. These included: 420 (1) without the use of AMTL, (2) with a weight emphasis on the BAS reward (0.25 for 421 ClogP, 0.25 for SAS, and 0.5 for BAS), (3) with the freezing of encoder weights for all 422 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 423 weights and pretraining on the 100-frequency and 1000-frequency dictionaries significantly 424 impair the model's learning progress. As such, these experiments were limited to a single 425 run of 25 epochs and excluded from further analysis. To evaluate the impact of omitting 426 AMTL and placing greater emphasis on the BAS reward, these model variations were run 427 three times for 70 epochs, while the baseline model ran five timess. Presented results are 428 based on run averages. 429

As displayed in Figure 3, incorporating AMTL improves training stability and enhances the 430 overall quality of optimized molecules. Findings in Table 3 support this by demonstrating 431 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 9
(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	-	$\sim 60 \text{ hrs}$	$\sim 60 \text{ hrs}$	$\sim 60 \text{ hrs}$
Validity	-	$99.54\% (\pm 0.08)$	99.33% (±0.34)	$99.64\%~(\pm 0.10)$
Novelty	-	$99.98\% (\pm 0.01)$	99.97% (± 0.02)	99.99% (± 0.00)
Uniqueness	-	$89.31\% (\pm 1.59)$	67.30% (±3.88)	90.49% (±1.91)
ClogP	62.06%	$94.57\% (\pm 1.35)$	95.04% (±1.90)	95.83% (± 0.50)
SAS	25.64%	$57.08\%~(\pm 0.92)$	51.67% (±2.17)	$56.74\% (\pm 0.70)$
BAS (LPA1)	9.79%	$53.87\%~(\pm 2.05)$	31.71% (±1.30)	$46.79\% (\pm 2.90)$
All Properties	0%	$32.22\%~(\pm 1.15)$	17.20% (±1.03)	29.20% (±1.92)

486 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 unex-plored. 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 av-erage 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). Propertywise distributions are seen in Figure 10 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	-	$\sim 1.25 \text{ hrs}$
Validity	-	$99.44\% \ (\pm 0.14)$
Novelty	-	99.81% (±0.10)
Uniqueness	-	89.50% (± 0.94)
ClogP	51.22%	86.01% (± 2.53)
SAS	39.15%	$51.55\% (\pm 1.49)$
BAS (LPA1)	16.90%	38.58% (±1.62)
All Properties	0%	19.11% (± 0.55)

A more detailed discussion of results is found in Appendix A.10.

5 Conclusion

In this work, we present SMORE-DRL, a scalable gradient-alignment-based MODRL frame-work for molecular optimization. A novel hybrid fragment-SMILES representation to depict molecules enables SMORE-DRL to selects 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 molecular optimizer, significantly outperforming other MODRL models while demonstrating scalability. Furthermore, when evaluated on new molecules post fine-tuning, SMORE-DRL effectively generalized its opti-mization 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 ver-sion, but with greater similarity to lead compounds. The implementation of SMORE-DRL is available at https://anonymous.4open.science/r/SMORE-DRL-F38B.

540 REFERENCES

551

561

562

563

571

572

573

578

579

580

- Aws Al-Jumaily, Muhetaer Mukaidaisi, Andrew Vu, Alain Tchagang, and Yifeng Li. Examining multi-objective deep reinforcement learning frameworks for molecular design. *Biosystems*, 232:104989, 2023.
- Kai Arulkumaran, Marc Deisenroth, Miles Brundage, and Anil Bharath. A brief survey of deep reinforcement learning. *IEEE Signal Processing Magazine*, 34, 08 2017. doi: 10.1109/MSP.2017.2743240.
- Giovanni Bolcato, Esther Heid, and Jonas Boström. On the value of using 3D shape and electrostatic similarities in deep generative methods. *Journal of Chemical Information and Modeling*, 62(6):1388–1398, 2022.
- 552 Nathan Brown. In silico Medicinal Chemistry: Computational Methods to Support Drug
 553 Design. Royal Society of Chemistry, 2015.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), 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), pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
 - Peter Ertl and Ansgar Schuffenhauer. Estimation of synthetic accessibility score of druglike molecules based on molecular complexity and fragment contributions. *Journal of Cheminformatics*, 1(1):1–11, 2009.
- Manan Goel, Shampa Raghunathan, Siddhartha Laghuvarapu, and U Deva Priyakumar.
 MoleGuLAR: Molecule generation using reinforcement learning with alternating rewards.
 Journal of Chemical Information and Modeling, 61(12):5815–5826, 2021.
- Sai Krishna Gottipati, Yashaswi Pathak, Boris Sattarov, Rohan Nuttall, Mohammad Amini, Matthew E Taylor, Sarath Chandar, et al. Towered actor critic for handling multiple action types in reinforcement learning for drug discovery. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 142–150, 2021.
 - Laura Graesser and Wah Loon Keng. Foundations of Deep Reinforcement Learning: Theory and Practice in Python. Addison-Wesley Professional, 2019.
- 574 Kexin Huang, Tianfan Fu, Wenhao Gao, Yue Zhao, Yusuf Roohani, Jure Leskovec, Connor W. Coley, Cao Xiao, Jimeng Sun, and Marinka Zitnik. Therapeutics Data Commons:
 576 Machine learning datasets and tasks for drug discovery and development. In *NeurIPS Datasets and Benchmarks*, 2021.
 - Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Hierarchical generation of molecular graphs using structural motifs. In *International Conference on Machine Learning*, pp. 4839–4848, 2020.
- Jintae Kim, Sera Park, Dongbo Min, and Wankyu Kim. Comprehensive survey of recent drug discovery using deep learning. *International Journal of Molecular Sciences*, 22(18), 2021.
- Yanjun Li, Mohammad A Rezaei, Chenglong Li, and Xiaolin Li. DeepAtom: A framework
 for protein-ligand binding affinity prediction. In 2019 IEEE International Conference on
 Bioinformatics and Biomedicine (BIBM), pp. 303–310. IEEE, 2019.
- Yu-Hsuan Lin, Yueh-Chien Lin, and Chien-Chin Chen. Lysophosphatidic acid receptor antagonists and cancer: the current trends, clinical implications, and trials. *Cells*, 10(7): 1629, 2021.
- Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent for multi-task learning. Advances in Neural Information Processing Systems, 34: 18878–18890, 2021.

594	Chunming Liu, Xin Z	Xu, and Dewen Hu.	Multiobjectiv	ve reinforce	ment learning	: A compre	e-
595	hensive overview.	IEEE Transactions	on Systems,	Man, and	Cybernetics:	Systems, 4	15
596	(3):385-398, 2015.		Ŭ,		Ŭ		
597							

- 598 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy,
 599 Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A robustly optimized
 600 bert pretraining approach, 2019.
- Hannes H Loeffler, Jiazhen He, Alessandro Tibo, Jon Paul Janet, Alexey Voronov, Lewis H
 Mervin, and Ola Engkvist. Reinvent 4: Modern AI-driven generative molecule design.
 Journal of Cheminformatics, 16(1):20, 2024.
- Muhetaer Mukaidaisi, Andrew Vu, Karl Grantham, Alain Tchagang, and Yifeng Li. Multiobjective drug design based on graph-fragment molecular representation and deep evolutionary learning. *Frontiers in Pharmacology*, 13:920747, 2022.
- Thanh Thi Nguyen, Ngoc Duy Nguyen, Peter Vamplew, Saeid Nahavandi, Richard Dazeley,
 and Chee Peng Lim. A multi-objective deep reinforcement learning framework. *Engineer- ing Applications of Artificial Intelligence*, 96:103915, 2020.
- Tiago Pereira, Maryam Abbasi, Bernardete Ribeiro, and Joel P Arrais. Diversity oriented deep reinforcement learning for targeted molecule generation. Journal of Cheminformatics, 13(1):1–17, 2021.
- Mariya Popova, Olexandr Isayev, and Alexander Tropsha. Deep reinforcement learning for
 de novo drug design. *Science Advances*, 4(7):eaap7885, 2018.
- 617
 618
 619
 620
 619
 620
 619
 620
 610
 611
 612
 612
 613
 614
 615
 615
 616
 617
 618
 619
 619
 619
 610
 610
 611
 612
 612
 613
 614
 615
 615
 616
 617
 617
 618
 619
 619
 619
 610
 610
 611
 612
 612
 614
 615
 615
 616
 617
 616
 617
 618
 619
 619
 610
 610
 610
 610
 611
 612
 612
 614
 615
 614
 615
 616
 617
 618
 619
 619
 610
 610
 610
 610
 610
 610
 610
 610
 611
 612
 612
 614
 614
 615
 614
 615
 614
 614
 615
 614
 614
 615
 614
 615
 614
 615
 614
 614
 614
 614
 614
 614
 614
 614
 614
 614
 614
 614
 614
 614
 614
 614
 614
 614
 614
 614
 614
 614
 614
 614
 614
- S. Shreyashree, Pramod Sunagar, S. Rajarajeswari, and Anita Kanavalli. A literature review on bidirectional encoder representations from transformers. In S. Smys, Valentina Emilia Balas, and Ram Palanisamy (eds.), *Inventive Computation and Information Technologies*, pp. 305–320, Singapore, 2022. Springer Nature Singapore.
- Maria Sorokina, Peter Merseburger, Kohulan Rajan, Mehmet Aziz Yirik, and Christoph
 Steinbeck. COCONUT online: Collection of open natural products database. Journal of Cheminformatics, 13(1):2, 2021.
- Niclas Ståhl, Göran Falkman, Alexander Karlsson, Gunnar Mathiason, and Jonas Boström.
 Deep reinforcement learning for multiparameter optimization in de novo drug design. Journal of Chemical Information and Modeling, 59(7):3166–3176, 2019.
- Duxin Sun, Wei Gao, Hongxiang Hu, and Simon Zhou. Why 90% of clinical drug development fails and how to improve it? Acta Pharmaceutica Sinica B, 12(7):3049–3062, 2022.
- Richard S Sutton and Andrew G Barto. Reinforcement Learning: An Introduction. MIT press, 2018.
- Huidong Tang, Chen Li, Shuai Jiang, Huachong Yu, Sayaka Kamei, Yoshihiro Yamanishi,
 and Yasuhiko Morimoto. EarlGAN: An enhanced actor-critic reinforcement learning
 agent-driven gan for de novo drug design. *Pattern Recognition Letters*, 175:45–51, 2023.
- Shidi Tang, Ji Ding, Xiangyu Zhu, Zheng Wang, Haitao Zhao, and Jiansheng Wu. Vina-GPU 2.1: Towards further optimizing docking speed and precision of autodock vina and its derivatives. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, pp. 1–13, 2024.
- Maha Thafar, Mona Alshahrani, Somayah Albaradei, Takashi Gojobori, Magbubah Essack, and Xin Gao. Affinity2Vec: Drug-target binding affinity prediction through representation learning, graph mining, and machine learning. *Scientific Reports*, 12, 03 2022.

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N.
 Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Neural In*formation Processing Systems, 2017.
- Jing Wang and Fei Zhu. ExSelfRL: An exploration-inspired self-supervised reinforcement learning approach to molecular generation. *Expert Systems with Applications*, 260:125410, 2024. ISSN 0957-4174.
- Alexander Wettig, Tianyu Gao, Zexuan Zhong, and Danqi Chen. Should you mask 15% in masked language modeling? In Andreas Vlachos and Isabelle Augenstein (eds.), Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, pp. 2985–3000, Dubrovnik, Croatia, May 2023. Association for Computational Linguistics.
- ⁶⁶⁰ David S. Wishart, Yannick D. Feunang, An C. Guo, Elvis J. Lo, Ana Marcu, Jason R. Grant, Tanvir Sajed, Daniel Johnson, Carin Li, Zinat Sayeeda, Nazanin Assempour, Ithayavani Iynkkaran, Yifeng Liu, Adam Maciejewski, Nicola Gale, Alex Wilson, Lucy Chin, Ryan Cummings, Diana Le, Allison Pon, Craig Knox, and Michael Wilson. DrugBank 5.0: A major update to the DrugBank database for 2018. *Nucleic Acids Research*, 46(Database): D1074–D1082, 2018.
- Soojung Yang, Doyeong Hwang, Seul Lee, Seongok Ryu, and Sung Ju Hwang. Hit and
 lead discovery with explorative RL and fragment-based molecule generation. Advances in *Neural Information Processing Systems*, 34, 2021.
- Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea
 Finn. Gradient surgery for multi-task learning. Advances in Neural Information Processing
 Systems, 33:5824–5836, 2020.
- Kiao-Chen Zhang, Cheng-Kun Wu, Jia-Cai Yi, Xiang-Xiang Zeng, Can-Qun Yang, Ai-Ping Lu, Ting-Jun Hou, and Dong-Sheng Cao. Pushing the boundaries of molecular property prediction for drug discovery with multitask learning BERT enhanced by SMILES enumeration. *Research*, 2022:0004, 2022.
- Zhenpeng Zhou, Steven Kearnes, Li Li, Richard N Zare, and Patrick Riley. Optimization of molecules via deep reinforcement learning. *Scientific Reports*, 9(1):1–10, 2019.

651

702 A APPENDIX

704 A.1 TRANSFORMER-ENCODER 705

Transformer-encoder models are made for pretraining on unlabeled data in a bidirectional 706 fashion (Vaswani et al., 2017; Devlin et al., 2019; Shreyashree et al., 2022). To extract 707 features, an embedding layer transforms the input fragment tokens $x = (x_1, x_2, \ldots, x_n)$ 708 into learnable embedding vectors $w = (w_1, w_2, \ldots, w_n)$, with the addition of a sinusoidal 709 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 710 711 712 multi-head self-attention layer, it is linearly transformed into the following h = 1, 2, ..., H713 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 714 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: $\mathbf{O}_h = \operatorname{softmax} \left(\frac{\mathbf{Q}_h \mathbf{K}_h^T}{\sqrt{d_k}} \right) \mathbf{V}_h$, 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 715 716 717 718 fed into a linear layer. Finally, during pretraining, the attention output is processed by a 719 feed-forward network, referred to as the "pretraining head." This head is typically replaced 720 with task-specific head during the fine-tuning stage (Zhang et al., 2022). 721

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 \tilde{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 \tilde{x} to recover x:

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

733 A.2 Reinforcement Learning

729 730 731

734 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 735 action spaces, \mathcal{P} is the state transition probability distribution $\mathcal{P}(s_{t+1}|s_t, a_t), \mathcal{R}$ is the reward 736 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 738 reward function for t (Al-Jumaily et al., 2023; Graesser & Keng, 2019). The goal of an RL 739 agent is to learn a policy distribution $\pi(a_t|s_t)$ that maximizes long-term cumulative rewards 740 through exploration of the environment over multiple timesteps. This is accomplished by 741 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 & Barto, 2018). To assess the value of states and actions 742 with respect to expected long-term returns, two functions are formulated: $V^{\pi}(s)$, which 743 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: $Q^{\pi}(s,a) = \mathbb{E}_{s_0=s,a_0=a,\tau\sim\pi} \left[\sum_{t=0}^T \gamma^t r_t \right]$ 744 745 746 (Graesser & Keng, 2019). 747

748 When the RL task entails exploring a vast state or action space, as is often the case in 749 drug design, learning an exact optimal policy or value function can become computationally 750 intractable. Thus, deep reinforcement learning (DRL) is used to approximate policies or 751 value functions (Arulkumaran et al., 2017). The actor-critic framework approximates both 752 and has been leveraged by various drug development frameworks (Al-Jumaily et al., 2023; 753 Goel et al., 2021; Gottipati et al., 2021; Pereira et al., 2021; Popova et al., 2018; Ståhl et al., 2019; Tang et al., 2023; Wang & Zhu, 2024; Yang et al., 2021). The actor model is 754 responsible for learning a parameterized policy π_{θ_A} , guided by feedback, known as temporal 755 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 & 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:

$$A^{\pi}(s_{t}, a_{t}) = Q^{\pi}(s_{t}, a_{t}) - V^{\pi}(s_{t})$$

$$\approx r_{t} + \gamma r_{t+1} + \gamma^{2} r_{t+2} + \dots + \gamma^{n} r_{t+n} + \gamma^{n+1} \hat{V}^{\pi}(s_{t+n+1}) - \hat{V}^{\pi}(s_{t}).$$
(8)

Thus, the value function is parameterized as $V^{\pi}_{\theta_C}(s)$ and is updated using loss function: 766

$$L_{\text{val}}(\theta_{C}) = \frac{1}{T} \sum_{t=0}^{T} \left(r_{t} + \hat{V}_{\theta_{C}}^{\pi}(s_{t+1}) - V_{\theta_{C}}^{\pi}(s_{t}) \right)^{2},$$

while the loss function for the actor is given by:

762 763

764

767

768 769 770

776

791

795 796

797

798

799

800

801 802

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

(9)

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.

782 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 783 gradient conflicts (Liu et al., 2021). However, the two metrics cannot offer a comprehensive 784 assessment if taken in isolation. One of the key components of AMTL is the proposal of the 785 condition number, a stability criterion that can indicate the presence of both challenges. 786 The value of the condition number is the ratio of the maximum and minimum singular 787 values of the corresponding matrix. Minimizing the condition number of the linear system 788 of gradients, a linear combination of gradients for all objectives, mitigates dominance and 789 conflicts within this system. If we apply singular value decomposition (SVD), we can have 790

$$G = U\Sigma V^{\mathrm{T}},\tag{11}$$

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

$$\boldsymbol{G}^{\mathrm{T}}\boldsymbol{G} = \boldsymbol{V}\boldsymbol{\Sigma}\boldsymbol{U}^{\mathrm{T}}\boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathrm{T}} = \boldsymbol{V}\boldsymbol{\Sigma}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathrm{T}} = \boldsymbol{V}\boldsymbol{\Lambda}\boldsymbol{V}^{\mathrm{T}}, \qquad (12)$$

where $\mathbf{\Lambda} = \operatorname{diag}(\lambda_1, \lambda_2, \dots, \lambda_K)$ 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 eigendecomposition 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:

$$\widehat{\boldsymbol{G}} = \boldsymbol{U}\widehat{\boldsymbol{\Sigma}}\boldsymbol{V}^{\mathrm{T}} = \boldsymbol{U}\boldsymbol{\sigma}\boldsymbol{I}\boldsymbol{V}^{\mathrm{T}} = \boldsymbol{\sigma}\boldsymbol{U}\boldsymbol{V}^{\mathrm{T}} = \boldsymbol{\sigma}\boldsymbol{G}\boldsymbol{V}\boldsymbol{\Sigma}^{-1}\boldsymbol{V}^{\mathrm{T}},\tag{13}$$

where $\sigma = \sqrt{\lambda_K}$ and $U = GV\Sigma^{-1}$ because of Equation (11), and \hat{G} is the aligned gradient matrix. A linear combination of the aligned objective-specific gradient vectors using the objective importance would be $\hat{G}\omega = \sum_{k=1}^{K} \omega_k \hat{g}_k$. 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 finetuning 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 foreach minibatch do foreach k = 1, 2, ..., K do $\mathbf{4}$ Compute loss $\mathcal{L}_k(\phi)$; Compute gradient $\boldsymbol{g}_{k} = \nabla_{\phi} \mathcal{L}_{k}(\phi);$ end Get the gradient matrix $G = \{g_1, ..., g_K\}$; // playing objective-specific gradient vectors as columns in GCompute task space Gram matrix $M \leftarrow G^{\mathrm{T}}G$; Get eigen-values and eigen-vectors $(\lambda, V) \leftarrow eigen(M)$; // eigen-decomposition such that $m{M} = m{V} m{\Lambda} m{V}^{ extsf{T}}$ where $m{\Lambda} = ext{diag}(m{\lambda})$ $\boldsymbol{\Sigma}^{-1} \leftarrow \operatorname{diag}\left(\sqrt{\frac{1}{\lambda_1}}, ..., \sqrt{\frac{1}{\lambda_K}}\right);$ Balance transformation $\boldsymbol{B} \leftarrow \sqrt{\lambda_n} \boldsymbol{V} \boldsymbol{\Sigma}^{-1} \boldsymbol{V}^T$; $\mathbf{12}$ Get new aligned gradient matrix $\widehat{G} = GB$; Updated gradient $\nabla \phi = \widehat{G}\omega$; Update policy parameter $\phi = \phi - \eta \nabla \phi$; $\mathbf{14}$ end $_{16}$ end 17 Return policy π_{ϕ} ;

864 A.4 Fragments-SMILES Hybrid Tokenization Strategy Figures



Figure 4: TDC MolGen task dataset (Huang et al., 2021) fragment frequencies.







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	$\sim 8 \text{ hrs}$
Epoch 2: MLM	0.93	0.75	$\sim 8 \text{ hrs}$
Epoch 3: Contrastive Learning	0.37	0.90	$\sim 20 \text{ hrs}$
Epoch 4: MLM	0.87	0.78	$\sim 8 \text{ hrs}$
Epoch 5: Contrastive Learning	1.70	0.91	$\sim 20 \text{ hrs}$
Epoch 6: MLM	0.87	0.79	$\sim 7 \text{ hrs}$

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	$\sim 3 \text{ hrs}$
Epoch 2: MLM	0.83	0.72	$\sim 3 \text{ hrs}$
Epoch 3: Contrastive Learning	0.29	0.89	$\sim 14 \text{ hrs}$
Epoch 4: MLM	0.80	0.73	$\sim 3 \text{ hrs}$
Epoch 5: Contrastive Learning	0.26	0.90	$\sim 15 \text{ hrs}$
Epoch 6: MLM	0.77	0.74	$\sim 3 \text{ hrs}$

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	$\sim 3 \text{ hrs}$
Epoch 2: MLM	0.81	0.72	$\sim 3 \text{ hrs}$
Epoch 3: Contrastive Learning	0.24	0.90	$\sim 15 \text{ hrs}$
Epoch 4: MLM	0.78	0.73	$\sim 3 \text{ hrs}$
Epoch 5: Contrastive Learning	0.22	0.91	$\sim 15 \text{ hrs}$
Epoch 6: MLM	0.76	0.74	$\sim 3 \text{ hrs}$

PROPERTY-WISE DENSITY PLOTS FOR THE COMPARATIVE AND SCALABILITY A.7 STUDIES



Figure 8: Property-wise comparisons between the lead molecules (blue) and the molecules optimized in final epoch by SMORE-DRL (red), DeepFMPOv3D (Bolcato et al., 2022) (green), DeepFMPO (Ståhl et al., 2019) (yellow), and MolDQN (Zhou et al., 2019) (purple). All objectives are to be minimized and the targeted maximums are indicated by the black dashed line.



Figure 9: 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 10: 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.





Figure 12: Lead molecules optimized by SMORE-DRL from the scalability experiments.





Figure 14: Lead molecules optimized by SMORE-DRL from the generalization experiments.

1350 A.10 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 compara-1359 tive study, it significantly outperformed all other models. In the scalability study, SMORE-1360 DRL performed efficiently, optimizing a set of lead molecules that did not achieve all proper-1361 ties so that one third of produced molecules satisfied all properties. Additionally, SMORE-1362 DRL's robustness allowed it to successfully generalize its optimization approach to unseen 1363 molecules. With just two modification steps, it improved the lead molecules from 0% to 19%1364 achieving all target properties. The inclusion of AMTL has proven to be a vital component 1365 of SMORE-DRL, enhancing training stability and improving the overall performance. 1366

As discussed, the primary objective of molecular optimization is developing a novel molecule 1367 similar to a lead molecule, aiming to have both molecules exhibit comparable qualities. As 1368 such, the progression of SMORE-DRL's optimized molecules were analyzed by comparing 1369 their similarity to lead molecules and their corresponding rewards across all timesteps. Fig-1370 ure 15 depicts the average similarity and reward for each of the four optimization timesteps 1371 performed on 1,000 molecules during the scalability study. To measure similarity, we utilize 1372 the method described in DeepFMPO (Ståhl et al., 2019), which employs a combination 1373 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 1375 0.5 and 0.7 is considered medium similarity (Loeffler et al., 2024). While SMORE-DRL 1376 does not achieve high similarity, it still presents strong results. As seen in Figure 15, there is an inverse correlation between average similarity and average sum of rewards across all 1377 objectives, where as similarity decreases, reward increases. This represents a trade-off: re-1378 stricting the optimization process to minimal modifications of a lead molecule may result in 1379 high similarity, but will likely restrict exploration and hinder the development of superior 1380 candidates. Nonetheless, the next iteration of SMORE-DRL should balance exploration 1381 with maintaining similarity to lead molecules, aiming to generate high-quality compounds 1382 without sacrificing similarity. One possible approach involves incorporation of a dynamic 1383 similarity component into the reward function, allowing for exploration in the initial training 1384 epochs while penalizing molecules with low similarity to lead molecules in the later epochs. 1385



Figure 15: 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.

1400 1401

1386 1387

1388

1390

1392

1393

1394

1395

1396 1397

1398

1399

1402