

RTMol: Rethinking Molecule-text Alignment in a Round-trip View

Letian Chen^{*1,2}, Runhan Shi^{*2}, Gufeng Yu², Yang Yang^{2†}

¹Shanghai Innovation Institute

²School of Computer Science, Shanghai Jiao Tong University

{clt2001, han.run.jiangming, jm5820zz}@sjtu.edu.cn, yangyang@cs.sjtu.edu.cn

Abstract

Aligning molecular sequence representations (e.g., SMILES notations) with textual descriptions is critical for applications spanning drug discovery, materials design, and automated chemical literature analysis. Existing methodologies typically treat molecular captioning (molecule-to-text) and text-based molecular design (text-to-molecule) as separate tasks, relying on supervised fine-tuning or contrastive learning pipelines. These approaches face three key limitations: (i) conventional metrics like BLEU prioritize linguistic fluency over chemical accuracy, (ii) training datasets frequently contain chemically ambiguous narratives with incomplete specifications, and (iii) independent optimization of generation directions leads to bidirectional inconsistency. To address these issues, we propose RTMol, a bidirectional alignment framework that unifies molecular captioning and text-to-SMILES generation through self-supervised round-trip learning. The framework introduces novel round-trip evaluation metrics and enables unsupervised training for molecular captioning without requiring paired molecule-text corpora. Experiments demonstrate that RTMol enhances bidirectional alignment performance by up to 47% across various LLMs, establishing an effective paradigm for joint molecule-text understanding and generation.

Code — <https://github.com/clt20011110/RTMol>

Introduction

Understanding molecular sequence representations and enabling *de novo* molecular design constitute fundamental challenges in cheminformatics and computational chemistry (Mouchlis et al. 2021). With molecules expressible as linear strings through the Simplified Molecular Input Line Entry System (SMILES) (Weininger 1988), recent studies have employed large language models (LLMs) for molecular understanding and generation (Li et al. 2024a; Zhang et al. 2024b; Sadeghi et al. 2024). These models leverage vast biochemical knowledge embedded in scientific literature and hold promise for unifying representation learning across textual descriptions and molecular structural domains. (See Appendix A for more related work.)

^{*}These authors contributed equally.

[†]Corresponding author.

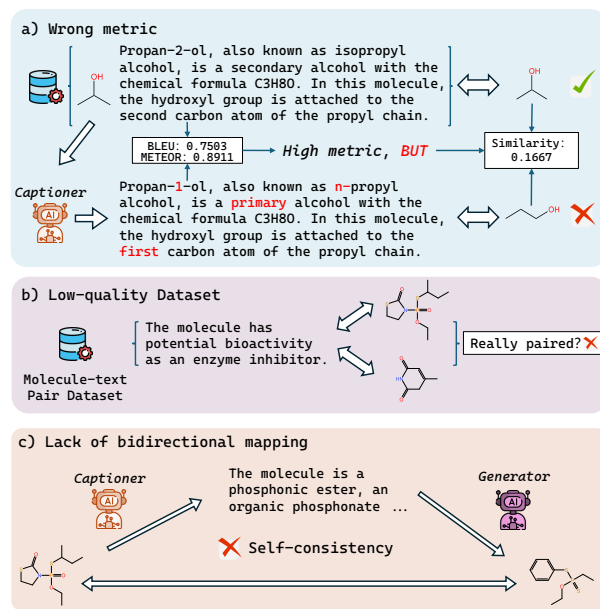


Figure 1: Limitations of current molecule-text alignment: (a) textual metrics ignore chemical fidelity; (b) captions are noisy or incomplete; and (c) separate modeling fails to enforce bidirectional consistency.

The critical challenge resides in establishing bidirectional alignment between molecular sequences and textual descriptions, a prerequisite for robust cross-modal reasoning. Current methodologies approach this through two decoupled tasks: i) Molecular captioning (molecule-to-text), and ii) Text-based molecular design (text-to-molecule). Although conceptually inverse, existing frameworks optimize these objectives independently through separate training regimes. Moreover, the prevailing reliance of current methodologies on supervised learning creates a dependency on high-quality paired datasets and robust evaluation metrics. As illustrated in Fig. 1, this paradigm introduces three fundamental limitations.

First, evaluation standards for molecular captioning are often misleading. Metrics such as BLEU (Papineni et al. 2002) and METEOR (Banerjee and Lavie 2005), which are commonly used to assess generated texts, primarily mea-

sure textual similarity through n -gram overlap. While these metrics reward fluent and keyword-rich captions, they often ignore chemical factuality. As a result, captions that achieve high BLEU/METEOR scores may still misrepresent molecular structures, containing incorrect or misleading chemical details. Recent analyses, such as those in molecular captioning benchmarks, demonstrate that high-scoring captions may violate chemical facts, indicating a weak alignment between textual quality and chemical correctness (Guo et al. 2023).

Second, the quality of existing datasets is problematic. Many publicly available molecule-text datasets suffer from ambiguous or incomplete descriptions (Edwards et al. 2022; Zhou et al. 2024). These descriptions are often generic, failing to uniquely identify the associated molecule or to describe key structural features. Using such datasets for supervised learning in molecular captioning tasks impairs the model’s ability to generalize reliably during inference.

Third, the lack of bidirectional alignment causes inconsistent understanding. Current training treats molecular captioning and text-based molecular design as separate tasks. Consequently, a model proficient in one often fails when its output is used as input for the other, leading to fragmented knowledge of molecule-text relationships (Guo et al. 2023). This inconsistency is a key barrier to unified modeling. Furthermore, achieving this bidirectional alignment is significantly hindered by the two previously mentioned limitations: the absence of robust evaluation metrics and large-scale, high-quality datasets.

To address these challenges, we propose a round-trip learning framework that unifies molecular captioning and text-based molecular design into a single training paradigm. The framework encourages round-trip consistency: the model first generates a caption from a molecule and then reconstructs the molecule solely from that caption. The similarity between the original and reconstructed molecules serves as a reward signal, directly optimizing for chemically faithful descriptions and bypassing reliance on noisy labels or purely textual metrics. This integrated approach promotes robust bidirectional alignment between molecular and textual representations. Our main contributions are:

- We propose a round-trip metric that evaluates captions by chemical fidelity rather than textual overlap, overcoming the limitations of conventional metrics.
- We present RTMol, a round-trip reinforcement learning framework that jointly aligns molecular captioning and text-based molecule generation.
- Our framework enables self-supervised training for molecular captioning, reducing reliance on noisy or incomplete annotations.
- Experiments demonstrate consistent improvements in both round-trip captioning and text-to-molecule tasks across various LLM backbones.

Method

Problem Formulation

Let \mathcal{M} denote the space of molecules and \mathcal{T} the space of text descriptions. We consider two conditional probability

distributions:

- $p_\theta(y|x)$, which models the probability of generating a textual description $y \in \mathcal{T}$ given a molecule $x \in \mathcal{M}$.
- $q_\phi(x'|y)$, which models the probability of generating a molecule $x' \in \mathcal{M}$ given a text description $y \in \mathcal{T}$.

The round-trip process involves first sampling text $y \sim p_\theta(y|x)$ given the original molecule x , and then sampling a reconstructed molecule $x' \sim q_\phi(x'|y)$.

To encourage faithful alignment between the molecular and textual modalities, we aim to minimize the expected discrepancy between the original molecule x and its reconstruction x' . This discrepancy is quantified by a distance function $d : \mathcal{M} \times \mathcal{M} \rightarrow \mathbb{R}_{\geq 0}$, which measures the structural difference between two molecules.

Formally, the goal is to learn optimal parameters θ^* and ϕ^* such that the expected reconstruction loss over the joint distribution of (x, y, x') is minimized. The sampling procedure follows $x \sim p(x)$, $y \sim p_\theta(y|x)$, $x' \sim q_\phi(x'|y)$, where $p(x)$ denotes the empirical distribution of molecules in the dataset. The optimization objective is defined as:

$$(\theta^*, \phi^*) = \arg \min_{\theta, \phi} \mathbb{E}_{x \sim p(x), y \sim p_\theta(y|x), x' \sim q_\phi(x'|y)} [d(x, x')]. \quad (1)$$

This objective aims to learn a pair of conditional distributions that effectively align the molecular space \mathcal{M} and the textual space \mathcal{T} , by minimizing the loss incurred during the round-trip process from a molecule to text and back to a reconstructed molecule.

Round-trip Metric

To evaluate the quality of molecule-text alignment, we introduce the round-trip metric R . The metric measures how well the model preserves information after translating the molecule into text and reconstructing it back into molecular form. The metric is defined as,

$$R(\theta, \phi) = \mathbb{E}_{x \sim p(x), y \sim p_\theta(y|x), x' \sim q_\phi(x'|y)} [\mathbb{1}\{d(x, x') = 0\}], \quad (2)$$

where $\mathbb{1}$ is a boolean indicator function. A higher R means a higher expectation that the reconstructed molecule x' is the same as the original molecule x , indicating better round-trip fidelity and stronger alignment between the two modalities.

Round-trip Optimization

We now show that minimizing Eq. 1’s objective is equivalent to maximizing a variational lower bound on the molecular-textual mutual information.

Let $X \in \mathcal{M}$ denote a random molecule drawn from $p(x)$, and let $Y \in \mathcal{T}$ be the text description corresponding to X , i.e. $Y | X \sim p_\theta(y | X)$. Viewing X and Y as two random variables, the mutual information between X and Y is

$$I(X, Y; \theta) = \mathbb{E}_{p(x) p_\theta(y|x)} [\log p_\theta(y|x) - \log p(y)]. \quad (3)$$

Using the Barber-Agakov variational decomposition (Barber and Agakov 2004), Eq. 3 becomes

$$\begin{aligned} I(X, Y; \theta) &= H(X) + \mathbb{E}_{p(x) p_\theta(y|x)} [\log q_\phi(x|y)] \\ &\quad + \text{KL}(p(x|y) \parallel q_\phi(x|y)) \\ &\geq H(X) + \mathbb{E}_{p(x) p_\theta(y|x)} [\log q_\phi(x|y)], \end{aligned} \quad (4)$$

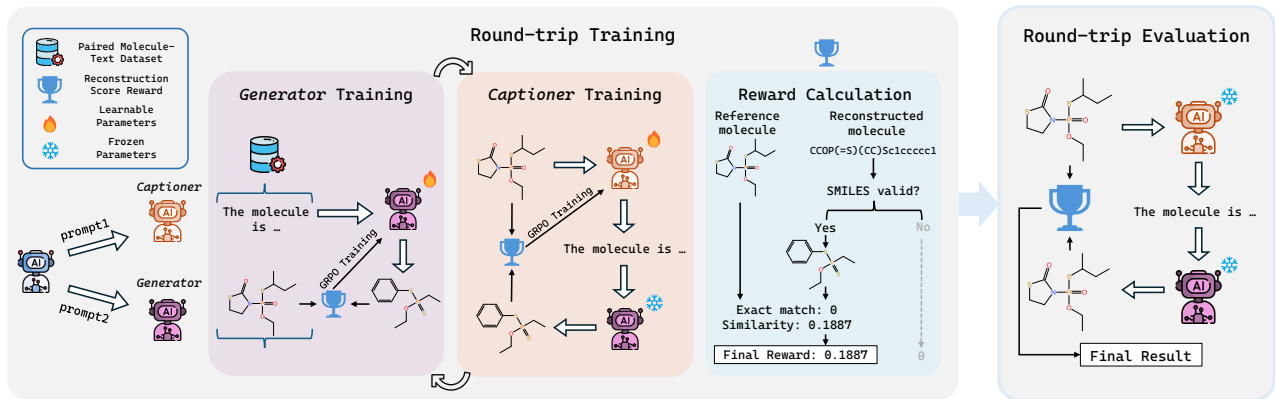


Figure 2: Overview of RTMol. A single LLM serves as both the Captioner and Generator for molecule-text alignment, with their training alternating cyclically in a complementary manner to reinforce each other’s performance.

where $H(X)$ is the entropy of X and $\text{KL}(\cdot||\cdot)$ denotes the Kullback-Leibler divergence. If q_ϕ is L -Lipschitz with respect to the distance function d defined above, then for some $\alpha > 0$, we have

$$\log q_\phi(x|y) \geq \mathbb{E}_{q_\phi(x'|y)} [-\alpha d(x, x') + C], \quad (5)$$

where C is a constant (proof in Appendix B). Substituting Eq. 5 into Eq. 4 yields

$$I(X, Y; \theta) \geq H(X) - \alpha \mathbb{E}_{p(x) p_\theta(y|x) q_\phi(x'|y)} [d(x, x')] + C, \quad (6)$$

which concludes the proof. Higher round-trip metrics correspond to smaller distances between original and reconstructed molecules, creating a theoretically grounded metric through its direct relationship with the distance function.

Model Architecture

Our framework, as shown in Fig. 2, uses a single large language model (LLM) that operates in two complementary roles:

- **Generator** ($q_\phi(x|y)$): converts textual descriptions into molecules.
- **Captioner** ($p_\theta(y|x)$): generates textual descriptions from molecules.

Unlike traditional methods that train these tasks independently, we couple them in a unified training process to establish consistent bidirectional alignment between molecules and text.

Generator: Grounding in Chemical Knowledge. The Generator translates natural language descriptions into valid molecular SMILES. It is trained using supervised fine-tuning on paired molecule-text data and optimized to maximize the reconstruction score S . The score evaluates three aspects:

- Similarity $S_{\text{sim}}(\cdot, \cdot)$: the similarity of molecular fingerprints between the original and reconstructed molecules.
- Exact match $S_{\text{exact}}(\cdot, \cdot)$: whether the reconstructed SMILES exactly matches the original.
- Validity $S_{\text{valid}}(\cdot)$: whether the generated SMILES string represents a chemically valid molecule.

Formally, the reconstruction score is defined as,

$$S(x, x') = \begin{cases} 0, & \text{if } S_{\text{valid}}(x') = 0, \\ S_{\text{sim}}(x, x') + S_{\text{exact}}(x, x'), & \text{otherwise,} \end{cases} \quad (7)$$

where

$$\begin{aligned} S_{\text{valid}}(x') &= \mathbf{1}\{x' \text{ is valid}\}, \\ S_{\text{exact}}(x, x') &= \mathbf{1}\{x = x'\}, \\ S_{\text{sim}}(x, x') &= T_{\text{MACCS}}(x, x') + T_{\text{RDKit}}(x, x') \\ &\quad + T_{\text{Morgan}}(x, x'). \end{aligned} \quad (8)$$

Here each $T_f(x, x')$ denotes the Tanimoto coefficient (Bajusz, Racz, and Haberger 2015) computed using a specific chemical fingerprint representation:

$$T_f(x, x') = \frac{|f(x) \cap f(x')|}{|f(x) \cup f(x')|}, \quad (9)$$

where $f \in \{\text{MACCS}, \text{RDKit}, \text{Morgan}\} : \mathcal{M} \rightarrow 2^{\{0, 1, \dots, n\}}$ is a fingerprint function that maps a molecule to a set of integer feature identifiers. These functions represent molecular structures through MACCS keys (Durant et al. 2002), RDKit (RDK) path-based fingerprints (Landrum 2024), and Morgan circular fingerprints (Rogers and Hahn 2010), respectively. We combine three molecular fingerprints to capture complementary chemical features, including substructures, physicochemical properties, and topological scaffolds. This multi-fingerprint similarity metric grounds the Generator’s understanding in established chemical representations, enabling it to serve as a reliable evaluator.

Captioner: Learning via Round-trip Process. The Captioner learns to describe a molecule with natural language. Critically, its training is **unsupervised** concerning textual labels: instead of matching reference captions, the Captioner learns through a round-trip interaction with the Generator. For a given molecule x , the Captioner first samples a description $y \sim p_\theta(y|x)$, which is immediately passed to the Generator to reconstruct a molecule $x' \sim q_\phi(x'|y)$. The Captioner’s objective is to maximize the reconstruction score defined in

Eq. 7 between x and x' , directly rewarding captions that preserve chemically relevant information. This self-consistency training bypasses the limitations of ambiguous or low-quality text labels, driving the model to produce descriptions that are both precise and informative for reconstruction.

Training Strategy

Choice of Optimization Algorithm. We optimize our model using reinforcement learning (RL) due to the non-differentiable nature of our objective. The final reconstruction score S defined in Eq. 7 involves components such as chemical validity checks and exact SMILES string matches, which cannot be optimized with standard backpropagation. We therefore frame the task as an RL problem, where the Captioner acts as a policy network. It generates textual descriptions (actions) to maximize an expected reward signal. For the optimization itself, we employ the Group Relative Policy Optimization (GRPO) (see Appendix C), which is well-suited for this rule-based reward setup.

Choice of Reward Function. Directly using the final round-trip metric $R(\theta, \phi)$ defined in Eq. 2 as the sole reward is ineffective, particularly in the early stages of training. The initial, untrained Captioner backbone frequently generates syntactically invalid SMILES strings and produces captions that result in a reconstruction score of zero. Such a sparse and uninformative reward signal fails to provide a meaningful gradient for optimization. As shown in Eq. 7, we include validity checking to enforce the Captioner to generate valid SMILES strings and a similarity reward to make the reward landscape smoother, providing a continuous optimization gradient even when an exact molecular match is not achieved.

Coupled Training. The Captioner and Generator are trained in a closed-loop process. Crucially, the Generator acts as an evaluator for the Captioner’s output, creating a dependency: the Captioner’s improvement is directly dependent on the Generator’s feedback. While the Generator’s training remains independent of the Captioner’s outputs, the two models are trained in parallel to prevent degradation of the Generator’s core text-to-molecule capability. This ensures the Captioner consistently learns from a stable, proficient evaluator, which is essential for refining molecular-textual alignment. The training pseudo-code is provided in Appendix D.

Experiments

Experimental Setup

Dataset. Our primary experiments use the ChEBI-20 dataset (Edwards, Zhai, and Ji 2021) (33,010 molecule-text pairs), with a standard 8:1:1 random split for training, validation, and testing. For fair comparison, main evaluations follow prior work by using a 100-sample benchmark test set (Zhao et al. 2024b). We additionally evaluate on two external datasets, L+M (Edwards et al. 2024) and Mol-Instruction (Fang et al. 2024). While Mol-Instruction shows high overlap with ChEBI-20 (88.9%), L+M suffers from underspecified textual descriptions. To prevent data leakage, we remove overlapping samples, yielding filtered versions L+M-F (180,178 pairs) and Mol-Instruct-F (239,659 pairs).

Baselines. We benchmark our model against two categories of baselines: 1) General-purpose LLMs: GPT-4o (OpenAI, Achiam et al. 2024), Gemini-2.5-Flash (Comanici, Bieber et al. 2025), Qwen-3-8B (Yang, Li et al. 2025), and DeepSeek-V3 (DeepSeek-AI, Liu et al. 2024); and 2) Domain-specific models: ChemT5-0.2B (Christofidelis et al. 2023), ChemDFM-8B (Zhao et al. 2024b), and ether0 (Narayanan et al. 2025). For the RTMol framework, we implement three variants using Qwen-3-8B (Qwen3), ChemT5-0.2B (ChemT5), and ChemDFM-8B (ChemDFM) as base models.

Training Details and Metrics. LLMs are trained following the process described in Section Method. We evaluate model performance using (1) Chemical metrics focus on assessing the quality of molecular outputs. We use the same metrics defined in Eq. 8. (2) Textual metrics assess the quality of generated natural language descriptions, and we can also apply them to molecular outputs. We report BLEU and METEOR, two widely-used generation metrics that measure n -gram overlap and semantic similarity between the generated and reference texts. Additional implementation details and prompts are provided in Appendix E.

Round-trip Molecular Captioning

Table 1 presents round-trip evaluation results. General-purpose LLMs like GPT-4o and Gemini-2.5-Flash perform poorly on chemical structure recovery metrics (MACCS, RDK, Morgan), indicating that strong general reasoning capabilities fail to ensure reliable molecular understanding in generative tasks.

Crucially, integrating RTMol consistently enhances performance across all backbones. With ChemDFM, exact match improves by roughly 47%, validity by 9%, and Morgan similarity by 31%. Similar gains are observed for ChemT5 and Qwen, demonstrating that RTMol is model-agnostic and broadly effective. These improvements extend beyond chemical metrics to textual alignment (BLEU, METEOR), reflecting enhanced bidirectional molecule-text understanding.

Among the three backbones evaluated, ChemDFM achieves the largest performance gains, followed by Qwen3, with ChemT5 showing the least improvement. This trend can be attributed to the characteristics of each model. Qwen3, being a general-purpose language model, benefits from the introduction of chemical rewards and round-trip training but remains limited by its lack of domain-specific knowledge, resulting in comparatively lower overall performance. ChemT5, in contrast, is domain-specific but significantly smaller (0.2B parameters). Its limited capacity leads to overfitting on the available datasets and leaves less room for improvement under our framework. ChemDFM combines a large model size with extensive pre-training on chemical data, making it well-suited to our round-trip optimization. As a result, it not only achieves the highest baseline performance but also exhibits the largest gains when enhanced with our method.

These findings highlight the value of round-trip evaluation as a more faithful measure of caption quality than traditional text-only metrics. Cases in Fig. 3 further illustrate how RTMol promotes chemically accurate and semantically

Model	Exact(%) \uparrow	Validity(%) \uparrow	MACCS \uparrow	RDK \uparrow	Morgan \uparrow	BLEU \uparrow	METEOR \uparrow
<i>Baselines</i>		<i>Chemical metrics</i>			<i>Textual metrics</i>		
GPT-4o	3.0	72.0	0.541	0.384	0.209	0.425	0.548
DeepSeek-V3	3.0	85.0	0.704	0.521	0.343	0.575	0.694
Gemini-2.5-Flash	17.0	40.0	0.374	0.343	0.293	0.597	0.565
ether0	4.0	71.0	0.385	0.291	0.198	0.280	0.405
<i>Round-trip training</i>							
Qwen3	7.0	60.0	0.409	0.326	0.234	0.392	0.569
Qwen3 + RTMol	<u>9.0</u>	<u>92.0</u>	<u>0.580</u>	<u>0.409</u>	<u>0.274</u>	<u>0.465</u>	<u>0.581</u>
ChemT5	12.0	86.0	0.691	0.605	0.482	<u>0.595</u>	0.695
ChemT5 + RTMol	<u>14.0</u>	<u>87.0</u>	<u>0.699</u>	<u>0.613</u>	<u>0.486</u>	0.590	<u>0.701</u>
ChemDFM	19.0	90.0	0.669	0.579	0.457	0.603	0.734
ChemDFM + RTMol	<u>28.0</u>	<u>98.0</u>	<u>0.826</u>	<u>0.729</u>	<u>0.597</u>	<u>0.722</u>	<u>0.812</u>

Table 1: Benchmark results of different models in round-trip evaluation using generated descriptions. The best results are highlighted in underline for each model.

Model	Exact(%) \uparrow	Validity(%) \uparrow	MACCS \uparrow	RDK \uparrow	Morgan \uparrow	BLEU \uparrow	METEOR \uparrow
<i>Baselines</i>		<i>Chemical metrics</i>			<i>Textual metrics</i>		
GPT-4o	5.0	75.0	0.593	0.432	0.277	0.427	0.555
DeepSeek-V3	18.0	86.0	0.764	0.618	0.468	0.530	0.674
Gemini-2.5-Flash	16.0	69.0	0.635	0.551	0.441	0.605	0.662
ether0	5.0	28.0	0.165	0.126	0.096	0.121	0.234
<i>Round-trip training</i>							
Qwen3	2.0	53.0	0.290	0.180	0.106	0.187	0.432
Qwen3 + RTMol	<u>3.0</u>	<u>88.0</u>	<u>0.469</u>	<u>0.281</u>	<u>0.142</u>	<u>0.283</u>	<u>0.465</u>
ChemT5	12.0	89.0	0.753	0.636	0.508	<u>0.519</u>	0.691
ChemT5 + RTMol	<u>13.0</u>	<u>90.0</u>	<u>0.765</u>	<u>0.648</u>	<u>0.524</u>	0.518	<u>0.694</u>
ChemDFM	24.0	96.0	0.864	0.755	0.620	0.797	0.903
ChemDFM + RTMol	<u>54.0</u>	<u>99.0</u>	<u>0.936</u>	<u>0.889</u>	<u>0.800</u>	<u>0.903</u>	<u>0.924</u>

Table 2: Benchmark results of different models for the text-based molecular design task using reference descriptions. The best results are highlighted in underline for each model.

informative descriptions.

Reference Text-based Molecular Design

We further evaluate all models using reference molecular descriptions from the ChEBI-20 dataset (Table 2). Compared to round-trip generation, performance generally improves, reflecting the benefit of high-quality textual input. For example, ChemDFM+RTMol achieves an 18% increase in RDK similarity. The only exception is ChemT5, which shows marginal change. Despite these gains, most LLMs still struggle to fully recover molecular SMILES from text. Chemical and textual scores remain modest, with validity below 90% across baselines. This highlights the persistent challenge of achieving precise molecular understanding and alignment, even when reference descriptions are clean and informative.

Applying our proposed RTMol consistently improves results across all backbones. Qwen3 benefits noticeably, while

ChemDFM shows the strongest improvement, with exact match and Morgan similarity rising by 125% and 29%, respectively. These results show that round-trip training complements LLMs, confirming the effectiveness of our self-supervised alignment strategy. Finally, we note that such evaluations depend on the quality and coverage of reference captions. While ChEBI-20 offers reliable descriptions, other datasets remain noisy or underspecified. We analyze the impact of description quality in the next section and present cross-dataset ablation studies to demonstrate the robustness of RTMol with less informative descriptions.

Addressing Noisy Datasets through Unsupervised Captioning

Current molecular captioning approaches rely heavily on supervised learning with paired datasets, making performance vulnerable to annotation quality. While curated datasets like

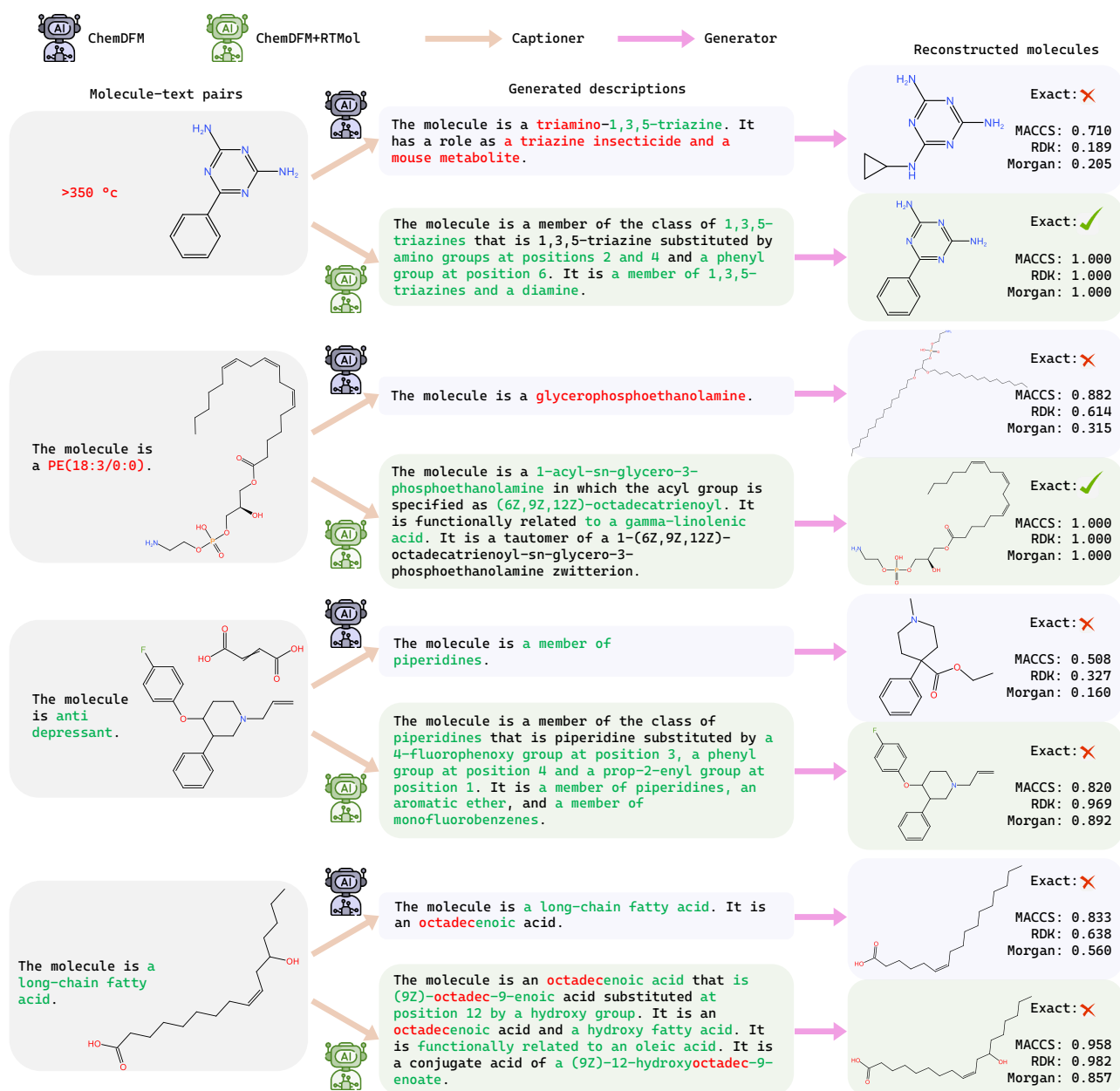


Figure 3: Cases of round-trip evaluation from the L+M-F and Mol-Instruct-F datasets. Mistakes and corrections are highlighted in red and green, respectively.

ChEBI-20 provide accurate molecule-text alignments, others like L+M-F and Mol-Instruct-F contain substantial noise, ambiguity, and generic descriptions. This quality gap severely impacts downstream performance: Table 3 (“reference” rows) shows that captions from these noisy datasets yield poor reconstruction for ChemDFM+RTMol (e.g., L+M-F exact match $<0.1\%$, Morgan similarity <0.15), indicating insufficient structural information.

Our unsupervised approach overcomes this limitation by generating captions without relying on reference texts. As shown in Table 3 (“generated” rows), it yields substantial

improvements on L+M-F, with exact match rate increasing by over 2.5 and fingerprint similarities rising by more than 0.45 across all metrics. Similar gains are observed on both chemical and textual metrics on both datasets, demonstrating its effectiveness and robustness across diverse datasets.

Case studies (Fig. 3) demonstrate this qualitatively. Reference captions are often vague or even incorrect. With RTMol training, the model can produce accurate and detailed descriptions of input molecules and successfully reconstruct the original molecules (first and second cases). In other scenarios, it generates generally accurate descriptions that, while

Dataset	Exact(%) \uparrow	Validity(%) \uparrow	MACCS \uparrow	RDKit \uparrow	Morgan \uparrow	BLEU \uparrow	METEOR \uparrow
L+M-F	<i>Chemical metrics</i>			<i>Textual metrics</i>			
reference	<0.1	97.4	0.378	0.231	0.142	0.150	0.211
generated	<u>2.7</u>	<u>99.9</u>	<u>0.892</u>	<u>0.792</u>	<u>0.624</u>	<u>0.708</u>	<u>0.750</u>
Mol-Instruct-F							
reference	<u>4.3</u>	97.8	0.604	0.423	0.196	0.288	0.471
generated	2.6	<u>99.9</u>	<u>0.676</u>	<u>0.486</u>	<u>0.245</u>	<u>0.358</u>	<u>0.520</u>

Table 3: Molecular design task using reference and generated descriptions on L+M-F and Mol-Instruct-F datasets with ChemDFM+RTMol. The best results are highlighted in underline for each dataset.

Model	Dataset	Exact(%) \uparrow	Validity(%) \uparrow	MACCS \uparrow	RDKit \uparrow	Morgan \uparrow	BLEU \uparrow	METEOR \uparrow
ChemDFM		<i>Chemical metrics</i>			<i>Textual metrics</i>			
without RT	ChEBI-20	<u>30.0</u>	94.0	0.791	0.702	0.571	0.712	0.801
text. with RT	ChEBI-20	26.0	96.0	0.814	0.720	0.583	<u>0.735</u>	<u>0.823</u>
chem. with RT (ours)	ChEBI-20	28.0	<u>98.0</u>	<u>0.826</u>	<u>0.729</u>	<u>0.597</u>	0.722	0.812
without RT	L+M-F	1.2	96.6	0.839	0.741	0.577	0.669	<u>0.721</u>
with RT (ours)	L+M-F	<u>1.3</u>	<u>99.8</u>	<u>0.860</u>	<u>0.778</u>	<u>0.620</u>	<u>0.674</u>	<u>0.720</u>
without RT	Mol-Instruct-F	<u>2.1</u>	82.7	0.513	0.360	0.162	<u>0.348</u>	<u>0.517</u>
with RT (ours)	Mol-Instruct-F	1.3	<u>99.5</u>	<u>0.615</u>	<u>0.435</u>	<u>0.202</u>	0.328	0.498

Table 4: Ablation in round-trip evaluation using generated descriptions. The best results are highlighted in underline for each dataset.

not exact, still enable reconstruction of molecules with high structural similarity to the input. In the third case, our method accurately reconstructs the mixture’s principal active component, which is the key ingredient in antidepressant formulations. In the fourth case, we recover a long-chain fatty acid with precise substituent positions and close physicochemical properties. Both cases achieve similarity above 0.8 across all three descriptors, with the RDKit fingerprint yielding particularly high values (>0.95).

Unlike prior works, our method also provides a diagnostic application: the round-trip metric identifies high-confidence captions within noisy datasets (L+M-F, Mol-Instruct-F), enabling dataset curation (see Appendix F). This demonstrates our approach’s dual role in both *improving generation* and *evaluating dataset quality*.

Ablation Studies

To assess the contribution of each component in RTMol, we perform ablation studies under two configurations: (1) Reward signal: Replace the chemically grounded rewards (similarities) with purely textual metrics (BLEU and METEOR). (2) Round-trip objective: Remove the round-trip consistency loss and train only with unidirectional supervision (text-to-molecule and molecule-to-text). We then compare performance across multiple datasets to assess the impact of this objective. Results in Table 4 reveal two key findings:

Importance of chemical rewards. Substituting chemical similarities with textual metrics consistently degrades chemical fidelity (e.g., Morgan similarity drops about 2.3% on ChEBI-20). While BLEU and METEOR scores remain com-

parable, chemically aware rewards better capture structural correctness, underscoring their necessity for molecule-text alignment.

Effect of round-trip consistency. Removing the round-trip objective leads to lower validity and similarity across various datasets. For instance, validity decreases by approximately 4.1%, 3.2%, and 16.9% on ChEBI-20, L+M-F, and Mol-Instruct-F, respectively. Similarly, Morgan similarity drops by about 4.4%, 6.9%, and 19.8% on these datasets. This confirms that enforcing cycle consistency improves bidirectional alignment and yields more accurate reconstructions.

Overall, both chemically grounded rewards and round-trip consistency are indispensable components for achieving robust, generalizable alignment between molecules and text.

Conclusion

We present RTMol, a reinforcement learning framework that unifies molecular captioning and text-based molecular design through round-trip consistency. Unlike prior approaches that train these tasks separately and rely on text-centric evaluation, RTMol directly optimizes for chemical fidelity by rewarding accurate reconstruction. This formulation addresses noisy captions, misleading metrics, and fragmented bidirectional alignment. Experiments on multiple datasets demonstrate that RTMol consistently improves round-trip scores while generating higher-quality molecule-text pairs, demonstrating its adaptability to diverse pre-trained language models. In future work, we plan to extend RTMol to multi-modal chemical data (e.g., 3D structures, spectra) and downstream tasks such as drug discovery and reaction planning.

Acknowledgment

This work was supported by the National Natural Science Foundation of China (No. 62272300).

References

- Bajusz, D.; Rácz, A.; and Háberger, K. 2015. Why is Tanimoto index an appropriate choice for fingerprint-based similarity calculations? *Journal of Cheminformatics*, 7: 20.
- Banerjee, S.; and Lavie, A. 2005. METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. In Goldstein, J.; Lavie, A.; Lin, C.-Y.; and Voss, C., eds., *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, 65–72. Ann Arbor, Michigan: Association for Computational Linguistics.
- Barber, D.; and Agakov, F. 2004. The im algorithm: a variational approach to information maximization. *Advances in neural information processing systems*, 16(320): 201.
- Brislin, R. W. 1970. Back-translation for cross-cultural research. *Journal of cross-cultural psychology*, 1(3): 185–216.
- Christofidellis, D.; Giannone, G.; Born, J.; Winther, O.; Laino, T.; and Manica, M. 2023. Unifying Molecular and Textual Representations via Multi-task Language Modelling. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, 6140–6157. PMLR.
- Comanici, G.; Bieber, E.; et al. 2025. Gemini 2.5: Pushing the Frontier with Advanced Reasoning, Multimodality, Long Context, and Next Generation Agentic Capabilities. *arXiv preprint arXiv:2507.06261*.
- DeepSeek-AI; Liu, A.; et al. 2024. DeepSeek-V3 Technical Report. *arXiv preprint arXiv:2412.19437*.
- Durant, J. L.; Leland, B. A.; Henry, D. R.; and Nourse, J. G. 2002. Reoptimization of MDL Keys for Use in Drug Discovery. *Journal of Chemical Information and Computer Sciences*, 42(6): 1273–1280.
- Edwards, C.; Lai, T.; Ros, K.; Honke, G.; Cho, K.; and Ji, H. 2022. Translation between Molecules and Natural Language. In *2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022*, 375–413. Association for Computational Linguistics (ACL).
- Edwards, C.; Wang, Q.; Zhao, L.; and Ji, H. 2024. L+M-24: Building a Dataset for Language+Molecules @ ACL 2024. In *Proceedings of the 1st Workshop on Language + Molecules (L+M 2024)*, 1–9. Bangkok, Thailand: Association for Computational Linguistics.
- Edwards, C.; Zhai, C.; and Ji, H. 2021. Text2mol: Cross-modal molecule retrieval with natural language queries. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 595–607.
- Fang, H.; and Xie, P. 2022. An end-to-end contrastive self-supervised learning framework for language understanding. *Transactions of the Association for Computational Linguistics*, 10: 1324–1340.
- Fang, Y.; Liang, X.; Zhang, N.; Liu, K.; Huang, R.; Chen, Z.; Fan, X.; and Chen, H. 2024. Mol-Instructions: A Large-Scale Biomolecular Instruction Dataset for Large Language Models. In *The Twelfth International Conference on Learning Representations*.
- Godard, C.; Mac Aodha, O.; and Brostow, G. J. 2017. Unsupervised monocular depth estimation with left-right consistency. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 270–279.
- Guo, T.; Nan, B.; Liang, Z.; Guo, Z.; Chawla, N.; Wiest, O.; Zhang, X.; et al. 2023. What can large language models do in chemistry? a comprehensive benchmark on eight tasks. *Advances in Neural Information Processing Systems*, 36: 59662–59688.
- He, D.; Xia, Y.; Qin, T.; Wang, L.; Yu, N.; Liu, T.-Y.; and Ma, W.-Y. 2016. Dual learning for machine translation. *Advances in neural information processing systems*, 29.
- Landrum, G. 2024. RDKit: Open-source cheminformatics. <https://zenodo.org/records/12782092>.
- Li, J.; Liu, W.; Ding, Z.; Fan, W.; Li, Y.; and Li, Q. 2025. Large language models are in-context molecule learners. *IEEE Transactions on Knowledge and Data Engineering*.
- Li, J.; Liu, Y.; Liu, W.; Le, J.; Zhang, D.; Fan, W.; Zhou, D.; Li, Y.; and Li, Q. 2024a. MolReFlect: Towards In-Context Fine-grained Alignments between Molecules and Texts. *arXiv preprint arXiv:2411.14721*.
- Li, S.; Liu, Z.; Luo, Y.; Wang, X.; He, X.; Kawaguchi, K.; Chua, T.-S.; and Tian, Q. 2024b. Towards 3D Molecule-Text Interpretation in Language Models. In *The Twelfth International Conference on Learning Representations*.
- Liu, Q.; Xu, J.; Jiang, R.; and Wong, W. H. 2021. Density estimation using deep generative neural networks. *Proceedings of the National Academy of Sciences*, 118(15): e2101344118.
- Liu, S.; Nie, W.; Wang, C.; Lu, J.; Qiao, Z.; Liu, L.; Tang, J.; Xiao, C.; and Anandkumar, A. 2023a. Multi-modal molecule structure–text model for text-based retrieval and editing. *Nature Machine Intelligence*, 5(12): 1447–1457.
- Liu, S.; Zhang, D.; Tu, Z.; Dai, H.; and Liu, P. 2024. Evaluating Molecule Synthesizability via Retrosynthetic Planning and Reaction Prediction. *arXiv preprint arXiv:2411.08306*.
- Liu, Z.; Li, S.; Luo, Y.; Fei, H.; Cao, Y.; Kawaguchi, K.; Wang, X.; and Chua, T.-S. 2023b. MolCA: Molecular Graph-Language Modeling with Cross-Modal Projector and Uni-Modal Adapter. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 15623–15638.
- Mouchlis, V. D.; Afantitis, A.; Serra, A.; Fratello, M.; Padiamantis, A. G.; Aidinis, V.; Lynch, I.; Greco, D.; and Melagraki, G. 2021. Advances in De Novo Drug Design: From Conventional to Machine Learning Methods. *International Journal of Molecular Sciences*, 22.
- Narayanan, S. M.; Braza, J. D.; Griffiths, R.-R.; Bou, A.; Wellawatte, G.; Ramos, M. C.; Mitchener, L.; Rodrigues, S. G.; and White, A. D. 2025. Training a Scientific Reasoning Model for Chemistry. *arXiv preprint arXiv:2506.17238*.

OpenAI; Achiam, J.; et al. 2024. GPT-4 Technical Report. *arXiv preprint arXiv:2303.08774*.

Papineni, K.; Roukos, S.; Ward, T.; and Zhu, W.-J. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics, ACL '02*, 311–318. USA: Association for Computational Linguistics.

Pei, Q.; Wu, L.; Gao, K.; Zhu, J.; and Yan, R. 2024. 3D-MolT5: Leveraging Discrete Structural Information for Molecule-Text Modeling. In *International Conference on Learning Representations*.

Rogers, D.; and Hahn, M. 2010. Extended-Connectivity Fingerprints. *Journal of Chemical Information and Modeling*, 50(5): 742–754.

Sadeghi, S. S.; Bui, A.; Forooghi, A.; Lu, J.; and Ngom, A. 2024. Can large language models understand molecules? *BMC Bioinformatics*, 25.

Su, B.; Du, D.; Yang, Z.; Zhou, Y.; Li, J.; Rao, A.; Sun, H.; Lu, Z.; and Wen, J.-R. 2022. A molecular multimodal foundation model associating molecule graphs with natural language. *arXiv preprint arXiv:2209.05481*.

Weininger, D. 1988. SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules. *Journal of chemical information and computer sciences*, 28(1): 31–36.

Yang, A.; Li, A.; et al. 2025. Qwen3 Technical Report. *arXiv preprint arXiv:2505.09388*.

Yang, Y.; Shi, R.; Li, Z.; Jiang, S.; Lu, B.-L.; Yang, Y.; and Zhao, H. 2024. BatGPT-Chem: A Foundation Large Model For Retrosynthesis Prediction. *arXiv preprint arXiv:2408.10285*.

Zhang, D.; Liu, W.; Tan, Q.; Chen, J.; Yan, H.; Yan, Y.; Li, J.; Huang, W.; Yue, X.; Ouyang, W.; et al. 2024a. Chemllm: A chemical large language model. *arXiv preprint arXiv:2402.06852*.

Zhang, J.; Bian, Y.; Chen, Y.; and Yao, Q. 2024b. Unimot: Unified molecule-text language model with discrete token representation. *arXiv preprint arXiv:2408.00863*.

Zhang, S.; Li, H.; Chen, L.; Zhao, Z.; Lin, X.; Zhu, Z.; Chen, B.; Chen, X.; and Yu, K. 2025. Reasoning-Driven Retrosynthesis Prediction with Large Language Models via Reinforcement Learning. *arXiv preprint arXiv:2507.17448*.

Zhao, Z.; Chen, B.; Li, J.; Chen, L.; Wen, L.; Wang, P.; Zhu, Z.; Zhang, D.; Li, Y.; Dai, Z.; et al. 2024a. ChemDFM-X: towards large multimodal model for chemistry. *Science China Information Sciences*, 67(12): 220109.

Zhao, Z.; Ma, D.; Chen, L.; Sun, L.; Li, Z.; Xia, Y.; Chen, B.; Xu, H.; Zhu, Z.; Zhu, S.; et al. 2024b. ChemDFM: A Large Language Foundation Model for Chemistry. *arXiv preprint arXiv:2401.14818*.

Zhou, P.; Wang, J.; Li, C.; Wang, Z.; Liu, Y.; Sun, S.; Lin, J.; Wang, L.; and Zeng, X. 2024. Instruction multi-constraint molecular generation using a teacher-student large language model. *BMC Biology*, 23.

Zhou, T.; Krahenbuhl, P.; Aubry, M.; Huang, Q.; and Efros, A. A. 2016. Learning dense correspondence via 3d-guided

cycle consistency. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 117–126.

Zhu, J.-Y.; Park, T.; Isola, P.; and Efros, A. A. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, 2223–2232.

Appendix

A. Related Work

Molecule-text Modeling

The concept of molecule-text modeling is first brought up by Text2Mol (Edwards, Zhai, and Ji 2021), which introduces the ChEBI-20 dataset with pairs of molecule SMILES representations and their textual captions. MolT5 (Edwards et al. 2022) further introduces molecule-caption translation tasks to evaluate the quality of molecule-text alignment. Recent advancements in molecule-text modeling are largely driven by Large Language Models (LLMs) as they demonstrate great potential in language-related applications within scientific domains (Guo et al. 2023). LLMs such as ChemDFM (Zhao et al. 2024b), ChemLLM (Zhang et al. 2024a), and BatGPT-Chem (Yang et al. 2024) are pre-trained on chemical literature and SMILES data. These models are typically adapted for downstream tasks, such as molecule captioning and text-to-molecule generation, through supervised fine-tuning (SFT) on paired datasets. In-context learning techniques are also leveraged to help LLMs better understand molecules (Li et al. 2025, 2024a). To overcome the limitations of purely text-based representations, many approaches incorporate richer structural information. MoMu (Su et al. 2022), MoleculeSTM (Liu et al. 2023a), and MolCA (Liu et al. 2023b) combines 1D SMILES and 2D graph representations for molecule-to-text generation. 3D-MoLM (Li et al. 2024b) and 3D-MolT5 (Pei et al. 2024) leverage the 3D conformation of molecules to help LLM better understand molecules. ChemDFM-X (Zhao et al. 2024a) further combines five modalities in the field of chemistry to enhance the ability of molecule understanding. Separately, ether0 (Narayanan et al. 2025) leverages reinforcement learning methods to perform chemical reasoning tasks using verifiable rewards. However, the application of reinforcement learning to molecule-text alignment remains less explored. Our work aims to fill this gap by proposing a brand-new round-trip framework to systematically solve this problem.

Round-trip Strategy

Round-trip process, also known as cycle consistency, is a fundamental concept with broad applications across various scientific and engineering domains. Historically, this principle is employed by human translators for back-translation to verify and improve translation quality (Brislin 1970). It has since been adapted into machine learning in dual learning for machine translation (He et al. 2016) and to enable end-to-end self-supervised learning (Fang and Xie 2022). The applicability of cycle consistency also extends into computer vision, where it serves as a powerful constraint in tasks such as semantic alignment (Zhou et al. 2016), unsupervised depth

estimation (Godard, Mac Aodha, and Brostow 2017), and unpaired image-to-image translation (Zhu et al. 2017). This principle is also leveraged in chemistry, where a “retrosynthesis-reaction prediction” loop is utilized to evaluate molecular synthesizability and verify retrosynthesis pathways (Liu et al. 2024; Zhang et al. 2025). Beyond these empirical applications, the effectiveness of the round-trip process has been theoretically grounded in the area of density estimation (Liu et al. 2021).

B. Proof on Round-trip Optimization

Assumptions. We make two mild regularity assumptions.

1. **Lipschitz continuity.** Let $\ell_\phi(x | y) = \log q_\phi(x | y)$. There exists a constant $L > 0$ such that

$$|\ell_\phi(x | y) - \ell_\phi(x' | y)| \leq L d(x, x') \quad \forall x, x', y, \quad (10)$$

i.e. $\ell_\phi(\cdot | y)$ is L -Lipschitz w.r.t. the distance $d(\cdot, \cdot)$ defined in the main text.

2. **Random reconstruction.** For a given y , let $x' \sim q_\phi(x' | y)$ be a reconstructed molecule sampled from the Generator.

Inequality derivation. For any x' , the Lipschitz condition in Eq. 10 implies

$$\ell_\phi(x | y) = \log q_\phi(x | y) \geq \ell_\phi(x' | y) - L d(x, x'). \quad (11)$$

Taking the expectation over the variable $x' \sim q_\phi(x' | y)$ on both sides of Eq. 11 yields

$$\log q_\phi(x | y) \geq \mathbb{E}_{q_\phi(x' | y)} [-L d(x, x') + \ell_\phi(x' | y)]. \quad (12)$$

Setting $\alpha := L$ and the constant $C := \mathbb{E}_{q_\phi(x' | y)} [\ell_\phi(x' | y)]$ (which depends only on y) immediately gives the desired result

$$\log q_\phi(x | y) \geq \mathbb{E}_{q_\phi(x' | y)} [-\alpha d(x, x') + C]. \quad (13)$$

Interpretation. The inequality states that, under Lipschitz continuity, the log-likelihood assigned by the Generator to any molecule x cannot decay faster than linearly with its molecular distance from *any* sample x' drawn from the same Generator.

C. Reinforcement Learning Using Group Relative Policy Optimization

In this work, we use Group Relative Policy Optimization (GRPO) as our optimization strategy.

Given a question x from the dataset, we sample G completions $y_1, \dots, y_G \sim \pi(\cdot | x)$. Each is assigned a reward r_1, \dots, r_G and a corresponding advantage:

$$A_i = \frac{r_i - \text{mean}\{r_1, \dots, r_G\}}{\text{std}\{r_1, \dots, r_G\}}. \quad (14)$$

Given a single problem x and a group of completions $\{y_i\}$, the per-group objective is:

$$J(\theta, x, y_1, \dots, y_G) = \sum_{i=1}^G \frac{1}{|y_i|} \sum_{t=1}^{|y_i|} \left\{ \text{clip} \left(\frac{\pi_\theta(y_{i,t} | x, y_{i,<t})}{\pi_{\theta_{\text{old}}}(y_{i,t} | x, y_{i,<t})}, A_i, \epsilon \right) - \beta \hat{D}_{\text{KL}}[\pi_{\text{ref}} || \pi_r; x, y_i, z_t] \right\}, \quad (15)$$

where π_θ is the policy being optimized, $\pi_{\theta_{\text{old}}}$ is the policy from which we sampled rollouts, π_{ref} is a reference policy, and clip is the standard PPO clip function:

$$\text{clip}(r, A, \epsilon) = \min\{r \cdot A, \max\{\min\{r, 1 + \epsilon\}, 1 - \epsilon\} \cdot A\}. \quad (16)$$

D. Training Process

The pseudo-code for the training process of RTMol is shown in Algorithm 1.

Algorithm 1: Unified Molecule-Text Round-trip Training

Input: Molecule-text pair dataset $\mathcal{D} = \{(x_i, y_i)\}$; Large language model \mathcal{L}_Θ with parameters Θ ;

Prompt_{captioner}, Prompt_{generator}; batch size B ; step number k ; rollout number n

Output: Updated model parameters Θ^*

```

1: while not converged do
2:    $\mathcal{G}_\Theta \leftarrow \mathcal{L}_\Theta$  with Promptgenerator
3:    $\mathcal{C}_\Theta \leftarrow \mathcal{L}_\Theta$  with Promptcaptioner
4:   for  $t = 1$  to  $k$  do
5:     Sample batch  $\{(x_j, y_j)\}_{j=1}^B$  from  $\mathcal{D}$ 
6:     for  $j = 1$  to  $B$  do
7:       for  $s = 1$  to  $n$  do
8:          $x_j^{(s)} \sim \mathcal{G}_\Theta(y_j)$  //  $x_j^{(s)} \sim q_\phi(x' | y_j)$ 
9:          $r_j^{(s)} \leftarrow S(x_j, x_j^{(s)})$ 
10:      end for
11:    end for
12:     $\Theta \leftarrow$  Update  $\Theta$  using GRPO with rewards  $\{r_j^{(s)}\}$ 
13:  end for
14:   $\Theta_0 \leftarrow \Theta$ 
15:  for  $t = 1$  to  $k$  do
16:    Sample batch  $\{x_j\}_{j=1}^B$  from  $\mathcal{D}$ 
17:    for  $j = 1$  to  $B$  do
18:       $y'_j \sim \mathcal{C}_\Theta(x_j)$  //  $y'_j \sim p_\theta(y | x_j)$ 
19:      for  $s = 1$  to  $n$  do
20:         $x_j^{(s)} \sim \mathcal{G}_{\Theta_0}(y'_j)$ 
21:         $r_j^{(s)} \leftarrow S(x_j, x_j^{(s)})$ 
22:      end for
23:    end for
24:     $\Theta \leftarrow$  Update  $\Theta$  using GRPO with rewards  $\{r_j^{(s)}\}$ 
25:     $\Theta_0 \leftarrow \Theta$ 
26:  end for
27: end while
28:  $\Theta^* \leftarrow \Theta$ 
29: return  $\Theta^*$ 

```

E. Implementation Details

Hyper-parameters used in RTMol during training are shown in Table 5. We implement our framework on top of the VERL library¹ for GRPO training. The experiments are conducted on 8 NVIDIA A800 GPUs, requiring approximately 3, 14, and 50 hours to train ChemT5-0.2B, ChemDFM-8B, and Qwen-3-32B, respectively.

¹<https://github.com/volcengine/verl>

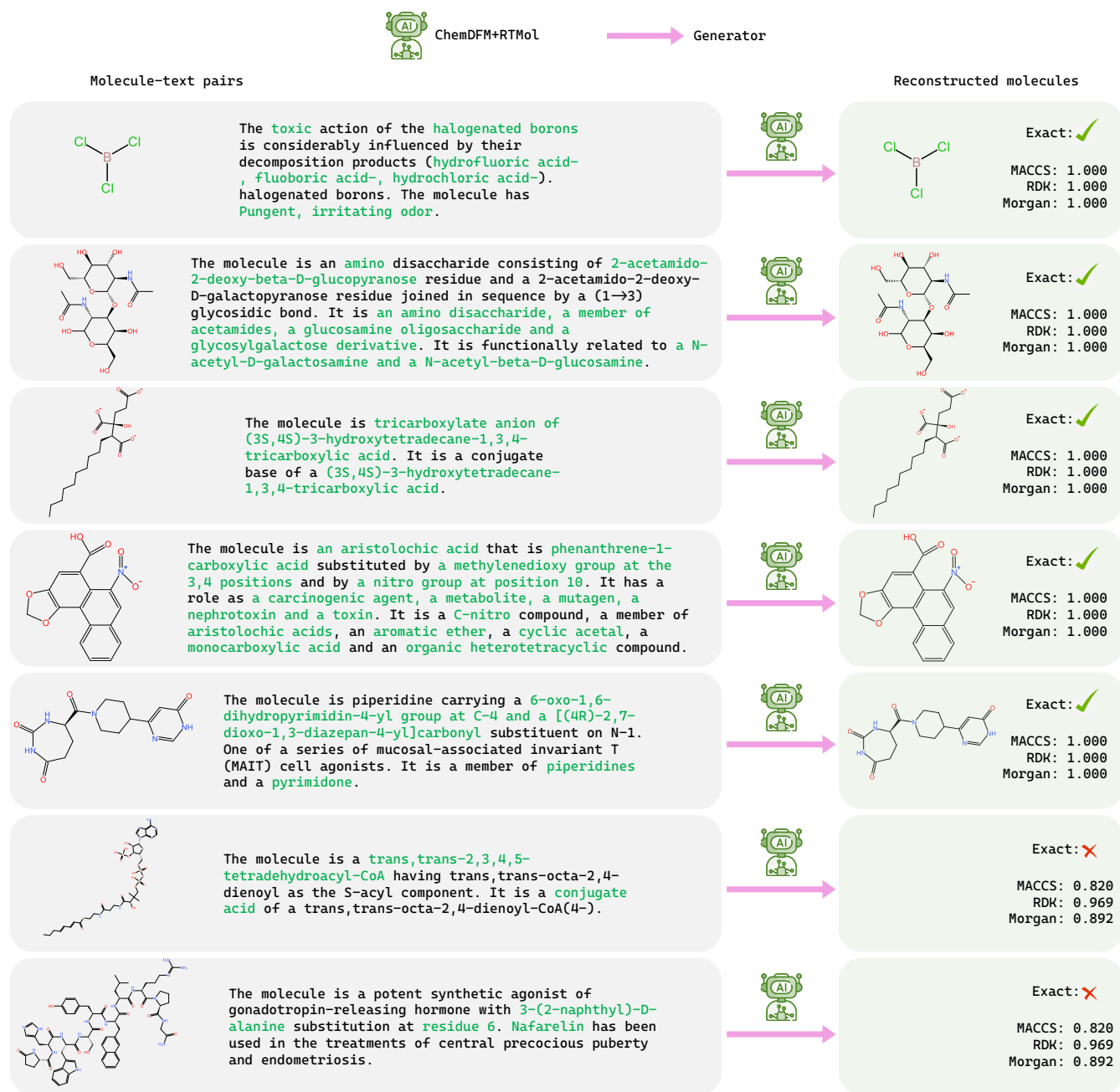


Figure 4: Examples of filtered high-quality molecule-text pairs from the L+M-F and Mol-Instruct-F datasets.

F. Examples of Generated Molecule Descriptions

Figure 4 demonstrates pairs that ChemDFM+RTMol achieves high round-trip scores.

Parameters	Value	Parameters	Value
Train batch size	128	KL loss coefficient	1e-3
PPO mini batch size	64	Learning rate	1e-6
Rollout number	32	Max steps	100

Table 5: Parameters for GRPO training.