# ChemMLLM: Chemical Multimodal Large Language Model

## **Anonymous ACL submission**

#### Abstract

Multimodal large language models (MLLMs) have made impressive progress in many applications in recent years. However, chemi-004 005 cal MLLMs that can handle cross-modal understanding and generation remain underex-007 plored. To fill this gap, in this paper, we propose ChemMLLM, a unified chemical multimodal large language model for molecule understanding and generation. Also, we design five multimodal tasks across text, molecular 011 SMILES strings, and image, and curate the datasets. We benchmark ChemMLLM against a range of general leading MLLMs and Chem-015 ical LLMs on these tasks. Experimental results show that ChemMLLM achieves supe-016 017 rior performance across all evaluated tasks. For example, in molecule image optimiza-019 tion task, ChemMLLM outperforms the best baseline (GPT-40) by 118.9% (4.27 vs 1.95 property improvement). The code is publicly available at https://anonymous.4open. science/r/ChemMLLM-0D98/.

## 1 Introduction

037

Multimodal large language models (MLLMs) have shown strong abilities in understanding and generating content across text, images, and audio (OpenAI, 2024; Sun et al., 2024; Team, 2024; Liu et al., 2024; Zhou et al., 2024; Xie et al., 2024; Wang et al., 2024), enabling more natural human–AI interaction. Chemistry is inherently multimodal, involving textual descriptions, structured formats like SMILES (Weininger, 1988)<sup>1</sup>, and molecular images. Recent works have demonstrated initial success in adapting MLLMs for chemical applications such as property prediction and reaction analysis (Cao et al., 2023; Zhang et al., 2024b; Luo et al., 2024; Li et al., 2025). However, these models largely focus on understanding tasks and treat images only as input. In contrast, molecular visuals are central to how chemists communicate and reason. Enabling image generation from chemical language or structure would greatly expand the expressiveness of MLLMs in chemistry (Kosenkov and Kosenkov, 2021). Yet, an integrated Chemical MLLM that supports both multimodal understanding and generation for chemistry remains lacking. 038

039

040

041

042

043

044

045

046

050

051

053

055

056

059

060

061

062

063

064

065

066

067

068

069

070

071

073

074

075

077

The challenges to build such a model are as follows: (1) Vision-based chemistry tasks and datasets remain underexplored. (2) Specificity of molecule image. Unlike natural images, molecule images are sparse, containing large areas of empty space, and are composed strictly of straight lines. General MLLMs result in unclear and distorted molecules (Figure 15). (3) Challenge for selecting an effective framework for chemical MLLM to seamlessly fuse diverse modalities, including discrete SMILES string/text, and continuous molecule images.

To address these issues, we propose ChemM-LLM, a chemical multimodal large language model that understands and generates molecules in a unified framework. Specifically, to handle three challenges above, (1) we identify five multimodal chemistry tasks with three modalities (text, SMILES, image), which contain both generation and comprehension tasks; (2) we finetune molecule imagelevel Vector Quantized Generative Adversarial Network (VQGAN) to bridge the gap between molecule images and natural images; (3) we introduce the "Image Tokenizer-LLM-Image Detokenizer" architecture into multimodal chemical tasks to fuse different modalities in early stages and enable models to generate images directly. Also, we design a two-stage training strategy and prove its effectiveness empirically.

**Contribution.** For ease of exposition, we summarize our main contribution as:

<sup>&</sup>lt;sup>1</sup>A SMILES (Simplified Molecular Input Line Entry System) string is a compact, text-based representation of a molecule's structure that encodes its atomic composition and connectivity in a linear format.



Figure 1: Overall architecture of ChemMLLM. (a) Image tokenizer and de-tokenizer. The image tokenizer employs CNN to extract spatial feature maps, where each  $n_z$ -dimensional spatial code is quantized into a discrete latent code via vector quantization (VQ). The resulting codebook indices serve as the final image tokens. Image de-tokenizer uses CNN to reconstruct image from discrete feature map. Then, a patch-based discriminator predicts whether the patch is fake (f) or real (r) (Section 3.1); (b) SMILES tokenizer and de-tokenizer. SMILES tokenization is consistent with text, and is mapped into a token sequence via text tokenizer; (c) ChemMLLM Training; (d) ChemMLLM Inference; (e) two-stage training paradigm for ChemMLLM (Section 3.3).

- A chemical multimodal LLM understanding and generating molecule in a unified framework. We propose and implement ChemMLLM, the first unified model to understand and generate molecules in text, SMILES, and image modality to the best of our knowledge (Section 3).
- A multimodal chemical dataset suite. We develop five new datasets to train and evaluate the chemical multimodal capability of MLLMs, which encompass a diverse spectrum of multimodal processing (Section 4).
- A variety of reliable evaluation. We benchmark the performance of different models on our proposed five tasks and ChemMLLM achieves dominating performance. For example, in molecule image optimization (image-to-image)

tasks, ChemMLLM outperforms the best baseline (GPT-40) by 118.9%, achieving a logP (optimizing property) increase of 4.27 compared to 1.95 (best baseline, GPT-40) (Section 5). 094

100

101

102

103

104

108

## 2 Related Work

Multimodal LLM (MLLM) With the advancement of large language models and multimodal learning, numerous high-performing multimodal large language models (MLLMs) have recently emerged. Notable MLLMs include GPT-40 (OpenAI, 2024), which extends GPT-4V (OpenAI, 2023) to understand and generate across different modalities, including text, image, and audio. Emu-3 (Wang et al., 2024) unifies vision understanding and generation via discrete token model-

Model	<b>BLEU-2</b> (†)	<b>BLEU-4</b> (†)	ROUGE-1 ( $\uparrow$ )	<b>ROUGE-2</b> $(\uparrow)$	ROUGE-L $(\uparrow)$	<b>METEOR</b> $(\uparrow)$
Qwen-VL-Chat InternVL-Chat-v1.5 LLaVA-v1.5-7B GPT-40	$\begin{array}{c} 0.09 \pm 0.001 \\ 0.04 \pm 0.0008 \\ 0.08 \pm 0.001 \\ \underline{0.16} \pm 0.003 \end{array}$	$\begin{array}{c} 0.01 \pm 0.0009 \\ 0.001 \pm 0.0002 \\ 0.004 \pm 0.0004 \\ 0.07 \pm 0.002 \end{array}$	$\begin{array}{c} 0.32 \pm 0.003 \\ \underline{0.38} \pm 0.003 \\ 0.33 \pm 0.002 \\ 0.28 \pm 0.005 \end{array}$	$\begin{array}{c} 0.07 \pm 0.001 \\ 0.07 \pm 0.001 \\ 0.07 \pm 0.001 \\ 0.13 \pm 0.003 \end{array}$	$\begin{array}{c} 0.22 \pm 0.002 \\ \underline{0.26} \pm 0.002 \\ 0.23 \pm 0.002 \\ 0.23 \pm 0.004 \end{array}$	$\begin{array}{c} 0.19 \pm 0.001 \\ 0.19 \pm 0.001 \\ 0.20 \pm 0.001 \\ 0.22 \pm 0.004 \end{array}$
ChemVLM-8B ChemMLLM (ours)	$\begin{array}{c} 0.15 \pm 0.004 \\ \textbf{0.33} \pm 0.005* \end{array}$	$\frac{0.08}{0.21} \pm 0.003$	0.28±0.006 <b>0.50</b> ±0.005*	$\frac{0.13}{0.31} \pm 0.003$	$0.23 \pm 0.005$ $0.43 \pm 0.005*$	$\frac{0.23}{0.43} \pm 0.005$

Table 1: Results on img2caption task (best: bold, second best: underlined, \*: significantly better (statistically)).

ing. Chameleon (Team, 2024) aligns modalities at 109 the token level for flexible multimodal generation, 110 while LlamaGen (Sun et al., 2024) treats images 111 as language-like sequences for scalable autoregres-112 sive image generation. Lumina-mGPT (Liu et al., 2024) trains a family of models to generate flex-114 ible photorealistic images from text descriptions 115 based on Chameleon. Transfusion (Zhou et al., 116 2024) and Show-o (Xie et al., 2024) combine diffu-117 sion and transformer for multimodal understanding 118 and generation. For vision understanding, well-119 known models include LLaVA (Liu et al., 2023), 120 BLIP-2 (Li et al., 2023), Qwen-VL-Chat (Bai 121 et al., 2023), InternVL-Chat (Chen et al., 2023), 122 MiniGPT-4 (Zhu et al., 2023), Gemini (Team et al., 123 2023), Flamingo (Alayrac et al., 2022) and Open-124 Flamingo (Awadalla et al., 2023). Although cur-125 rent MLLMs perform well across modalities, these 126 models struggle with chemical tasks due to a mis-127 alignment between general and domain-specific 128 knowledge. 129

Chemical LLM MLLMs have also exhibited 130 strong potential in addressing chemistry-related 131 tasks, particularly in bridging the modality gap be-132 tween textual descriptions and molecular represen-133 tations. Concretely, Instruct-Mol (Cao et al., 2023) 134 and MV-Mol (Luo et al., 2024) utilize LLaVA's ar-135 chitecture (Liu et al., 2023) and Q-former (Li et al., 2023) to align molecular structure and text modal-137 ity, respectively. ChemLLM (Zhang et al., 2024a) 138 uses a high-quality chemical dataset to fine-tune 139 InternLM2 (Cai et al., 2024). ChemVLM (Li et al., 140 2025) extends ChemLLM (Zhang et al., 2024a) to 141 understand images by adopting a projector-based 142 method to align vision information and text in-143 formation. UniMoT (Zhang et al., 2024b) uses 144 a molecule tokenizer to align the graph modality 145 146 molecule with text. Despite progress in chemical tasks, existing models can not generate molecu-147 lar images, limiting their utility in more intuitive, 148 visual forms of interaction. To fill this gap, We pro-149 pose ChemMLLM, a unified framework that under-150

stands and generates molecules in text, SMILES, and image formats.

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

183

184

185

186

187

188

189

190

191

# 3 Method

**Overview**. Our framework uses an image tokenizer to transfer images into discrete tokens, aligning with texts and molecule's SMILES strings at the token level, known as "Image Tokenizer-LLM-Image De-tokenizer" architecture (Team, 2024). First, Section 3.1 discusses the molecule image tokenizer. Then, Section 3.2 describes how to combine images, texts and SMILES strings. Finally, the training strategy is discussed in Section 3.3. The whole pipeline is shown in Figure 1. For ease of understanding, we list key mathematical notations in Table 10 in Appendix.

# 3.1 Mol-VQGAN for Molecule Image Generation

Following Chameleon (Team, 2024), to enable multimodal information alignment in the early stage, images need to be discretized into token sequences similar to text. To proceed it, we train Vector Quantized Generative Adversarial Network (VQ-GAN) (Esser et al., 2021) on molecule images, known as Mol-VQGAN. Mol-VQGAN compresses images into a discrete latent space, which aligns well with the sequential nature of text modeling in large language models. Specifically, Mol-VQGAN uses Vector Quantized Variational Auto-Encoder (VQVAE) (Van Den Oord et al., 2017) as the generator and patch-based discriminator (Isola et al., 2017). VQVAE compresses images into discrete spaces and reconstructs images from that space. Mol-VQGAN adds a discriminator and perceptual loss to VQVAE to keep good perceptual quality and improve the performance. Formally, Mol-VQGAN takes an image  $x \in \mathbb{R}^{H \times W \times 3}$  (H/W)are the height/width of the input image), and transfers it into discrete representations  $z_q$  by encoding  $\hat{z} \in \mathbb{R}^{h \times w \times n_z} = E(x) (h/w/n_z)$  are the height/width/channels of the feature map), and finding the closest codebook entry for each spatial code

		MW			LogP			TPSA		
Method	Pearson (†)	MSE $(\downarrow)$	MAE $(\downarrow)$	Pearson (†)	$\text{MSE}\left(\downarrow\right)$	MAE $(\downarrow)$	Pearson (†)	MSE $(\downarrow)$	MAE $(\downarrow)$	valid%( $\uparrow$ )
Qwen-VL-Chat InternVL-Chat-v1.5 LLaVA-v1.5-7B	$\begin{array}{c} 0.78 \pm \! 0.02 \\ 0.59 \pm \! 0.02 \\ 0.36 \pm \! 0.08 \end{array}$	$\frac{7073.3}{>1e+4} \pm 1595.5$ $>3e+4 \pm 5078.6$	$58.0 \pm 1.8$ $83.6 \pm 3.6$ $115.7 \pm 6.9$	$\begin{array}{c} 0.14 \pm \! 0.02 \\ 0.04 \! \pm \! 0.05 \\ - 0.003 \! \pm \! 0.05 \end{array}$	$\frac{4.4}{9.2\pm0.62}\pm0.62$ 5.1±1.0	$\frac{1.4}{2.3\pm0.08}\pm0.08$ 1.6±0.08	$\begin{array}{c} 0.19 \pm 0.02 \\ \underline{0.29} \pm 0.03 \\ 0.01 \pm 0.04 \end{array}$	$>1e+4 \pm 563.6$ $\underline{2158.8} \pm 432.6$ $>5e+4 \pm > 3e + 4$	$\begin{array}{c} 82.8 \pm 1.9 \\ \underline{28.7} \pm 1.5 \\ 99.1 \pm 11.9 \end{array}$	<u>99.3</u> % 55.0% 35.8%
ChemVLM-8B ChemMLLM (ours)	$\frac{0.84}{\textbf{0.97} \pm 0.004*} \pm 0.004$	9573.4 ±1992.9 789.7±162.2*	$\frac{56.9}{16.1 \pm 0.72*} \pm 4.5$	$\frac{0.38}{\textbf{0.92} \pm 0.01} \pm 0.01 ^{*}$	$\begin{array}{c} \textbf{4.9} \pm 0.58 \\ \textbf{0.70} {\pm 0.14} {*} \end{array}$	$\begin{array}{c} 1.6 \pm \! 0.08 \\ \textbf{0.52} \! \pm \! 0.02 * \end{array}$	$\begin{array}{c} 0.26 \pm \! 0.06 \\ \textbf{0.97} {\pm} 0.005^* \end{array}$	$5332.0 \pm 747.5 \\ \textbf{152.6} \pm 39.2^*$	$\begin{array}{c} 53.66 \pm 2.7 \\ \textbf{6.0} \pm 0.33 * \end{array}$	31.3% <b>99.6</b> %

Table 2: Results on img2property task: MW, LogP and TPSA (best, 2nd best, \*: significantly better (statistically)).

		Hbd			Hba			Rb			QED	
Method	Pearson (†)	$\text{MSE}\left(\downarrow\right)$	MAE $(\downarrow)$	Pearson (†)	$\text{MSE}\left(\downarrow\right)$	$\text{MAE}\left(\downarrow\right)$	Pearson (†)	$\text{MSE}\left(\downarrow\right)$	$\text{MAE}\left(\downarrow\right)$	Pearson (†)	$MSE(\downarrow)$	MAE $(\downarrow)$
Qwen-VL-Chat InternVL-Chat-v1.5 LLaVA-v1.5-7B	$\begin{array}{c} \text{-0.02} \pm 0.008 \\ 0.03 \pm 0.05 \\ 0.004 \pm 0.05 \end{array}$	$\begin{array}{r} 4.8 \pm 0.75 \\ 9.9 \pm 0.89 \\ 53.3 \pm 27.3 \end{array}$	1.4±0.05 2.2±0.09 3.7±0.33	$\substack{0.05 \pm 0.01 \\ 0.22 \pm 0.04 \\ 0.04 \pm 0.05}$	$30.5\pm1.2$ $10.2\pm0.85$ $23.7\pm5.4$	$\begin{array}{c} 4.9{\pm}0.08\\ \underline{2.4}{\pm}0.08\\ \overline{2.9}{\pm}0.19\end{array}$	$\begin{array}{c} \underline{0.19}{\pm 0.1} \\ 0.04{\pm 0.04} \\ 0.03 {\pm 0.04} \end{array}$	$145.4 \pm 30.8$ $43.6 \pm 4.3$ $\underline{39.9} \pm 10.4$	$6.8 \pm 0.31$ $4.8 \pm 0.19$ $3.8 \pm 0.26$	$\begin{array}{c} 0{\pm}0.0\\ \underline{0.003}{\pm}0.03\\ {-}0.11{\pm}0.08\end{array}$	$4.0\pm 3.2$ >1e+5 $\pm$ > 1e + 5	0.4±0.08 24.5±18.6
ChemVLM-8B ChemMLLM (ours)	$\begin{array}{c} \underline{0.49} \pm 0.09 \\ \textbf{0.96} \pm 0.004 * \end{array}$	$\begin{array}{c} \underline{4.2} \pm 0.85 \\ \textbf{0.18} \pm 0.02 * \end{array}$	$\begin{array}{c} \underline{1.3} \pm 0.08 \\ \textbf{0.13} \pm 0.01 * \end{array}$	$\begin{array}{c} \underline{0.32} \pm 0.07 \\ \textbf{0.94} \pm 0.007 * \end{array}$	$\begin{array}{c} 27.2 \pm 1.95 \\ \textbf{0.79} {\pm} 0.12^* \end{array}$	$\begin{array}{c} \textbf{4.5} \pm 0.14 \\ \textbf{0.44} {\pm 0.02} * \end{array}$	$\begin{array}{c} 0.10 \pm \! 0.06 \\ \textbf{0.94} {\pm} 0.01 {*} \end{array}$	$\begin{array}{c} 45.8 \pm \! 5.02 \\ \textbf{1.6} {\pm} 0.33^* \end{array}$	$\begin{array}{c} {\rm 5.5} \pm 0.22 \\ {\rm 0.59} \pm 0.03^* \end{array}$	$\begin{array}{c} \textbf{-0.003} \pm 0.06 \\ \textbf{0.91} \pm 0.006 \ast \end{array}$	$\begin{array}{c} \underline{0.24} \pm 0.02 \\ \textbf{0.008} \pm 0.0004 * \end{array}$	$\begin{array}{c} \underline{0.37} \pm 0.01 \\ \textbf{0.06} \pm 0.002 * \end{array}$

Table 3: Results on img2property task: Hbd, Hba, Rb, and QED (best, 2nd best, \*: significantly better (statistically)).

 $\hat{z}_{ij} \in \mathbb{R}^{n_z}$  (also known as vector quantization (VQ) process, denoted  $\mathbf{q}(\cdot)$ ):

192

193

194

197

199

200

$$z_q = \mathbf{q}(\hat{z}) = \left(\arg\min_{z_k \in Z} \|\hat{z}_{ij} - z_k\|\right) \in \mathbb{R}^{h \times w \times n_z},$$
(1)

where Z is codebook, a set of learnable vectors, and  $z_k \in Z = \{z_i\}_{i=1}^n \subset \mathbb{R}^{n_z}$ . Then, Mol-VQGAN reconstructs image  $\hat{x} \in \mathbb{R}^{H \times W \times 3}$  from  $z_q$ :

$$\hat{x} = G(z_q) = G(\mathbf{q}(E(x))). \tag{2}$$

The discretization and reconstruction process is illustrated in Figure 1(a).

The VQGAN's objective function is:

$$\min_{E,G,Z} \max_{D} \left[ \mathcal{L}_{vqvae}(E,G,Z) + \lambda_1 \mathcal{L}_{perceptual}(E,G,Z) + \lambda_2 \mathcal{L}_{GAN}(\{E,G,Z\},D) \right],$$
(3)

where (1) the first term  $\mathcal{L}_{vqvae}(E, G, Z) = ||x - U||^2$  $G(\hat{z} - sg(\hat{z} - z_q))||_2^2 + ||sg[\hat{z}] - z_q||_2^2 + ||sg[z_q] - \hat{z}||_2^2$ 204 is VQVAE loss, where  $\mathcal{L}_{rec} = ||x - G(\hat{z} - sg(\hat{z} - sg(\hat{$  $|z_q\rangle||_2^2$  is reconstruction loss. Since the quantization process is non-differentiable, a stop-gradient 207 operation  $sq[\cdot]$  is used so that the forward pass operates on quantized vectors, whereas the backward pass leverages continuous vectors for gradient com-210 putation (Van Den Oord et al., 2017); (2) the second 211 term  $\mathcal{L}_{perceptual}(E, G, Z) = ||P(x) - P(G(\hat{z} - C))||P(x)||P(x) - P(G(\hat{z} - C))||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||P(x)||$ 212  $sg(\hat{z} - z_q))|_2$  is perceptual loss, P denotes a 213 perceptual model like Learned Perceptual Image 214 Patch Similarity (LPIPS) (Zhang et al., 2018), 215 216 which is used to extract the high-level semantic features; (3) the third term  $\mathcal{L}_{GAN}(\{E, G, Z\}, D) =$ 217  $\log D(x) + \log(1 - D(\hat{z} - sg(\hat{z} - z_a)))$  is GAN loss, 218 D is a patch-based discriminator (Isola et al., 2017) 219 aiming to differentiate original and reconstructed 220

images.  $\lambda_1$  is a hyperparameter and  $\lambda_2$  is an adaptive weight computed dynamically to balance the weight of  $\mathcal{L}_{GAN}$  and stabilize training, following VQGAN (Esser et al., 2021).

## 3.2 ChemMLLM

Chemical molecules are typically encoded in the format of Simplified Molecular Input Line Entry System (SMILES) (Weininger, 1988), a compact ASCII string, serving as a distinct modality in computational chemistry (Anderson et al., 1987). In ChemMLLM, SMILES tokenization is the same as the text. Specifically, SMILES is mapped into a token sequence via Chaemelon (Team, 2024) text tokenizer trained based on Byte Pair Encoding (BPE) algorithm (Sennrich et al., 2015). The process is shown in Figure 1(b).

After training the Mol-VQGAN, ChemMLLM uses it as image tokenizer to align image and text at token level, unifying the training and inference for image and text by maximizing the standard nexttoken prediction cross-entropy loss:

$$\mathcal{L}_{LLM} = \sum_{i=1}^{L} \log p_{\theta}(s_i | s_1, ..., s_{i-1}) + \lambda \sum_{k} (\log \sum_{j=1}^{V} \exp(z_{k,j}))^2,$$
(4)

where  $\lambda$  is a hyper-parameter; in the first term, Lis the length of total sequence tokens;  $s_i \in S = \{S_I, S_T\}$ .  $S_I$  is the image tokens sequence tokenized by the image tokenizer, and  $S_T$  is the text/SMILES token sequence tokenized by the text tokenizer. The second term is z-loss, a regularization term to mitigate the problem of logit shift in the final softmax and stabilize training (Chowdhery 241

242

243

245

246

247

248

249

250

221

		MW			LogP			TPSA		
Method	Pearson (†)	MSE $(\downarrow)$	MAE $(\downarrow)$	Pearson (†)	MSE $(\downarrow)$	MAE $(\downarrow)$	Pearson (†)	MSE $(\downarrow)$	MAE $(\downarrow)$	valid%( $\uparrow$ )
InternVL-Chat-v1.5	$0.04{\pm}0.05$	$>1e+5\pm>1e+5$	$210.5 {\pm} 29.4$	$-0.02 \pm 0.11$	$23.4 \pm 9.3$	$2.96{\pm}0.31$	$0.04{\pm}0.11$	$> 1e + 4 \pm 6205.6$	$58.9 \pm 8.1$	75.0%
LLaVA-v1.5-7B	$0.07 {\pm} 0.27$	$> 2e+6\pm > 9e+5$	$894.9 {\pm} 285.0$	$-0.52 \pm 0.16$	$1104.4 {\pm} 522.3$	$18.0 {\pm} 6.7$	$0.48 \pm 0.27$	$>1e+5\pm>1e+5$	$146.3 {\pm} 88.5$	8.5%
GPT-40	$0.77 \pm 0.04$	7633.2±1954.0	55.7±4.7*	$0.48 {\pm} 0.07$	$5.7 \pm 0.87$	$1.7 \pm 0.11$	<u>0.69</u> ±0.04	$1209.2 \pm 211.1$	<b>24.3</b> ±1.7	<u>99.0</u> %
ChemLLM-7B-Chat	$-0.25 {\pm} 0.08$	$>7e+4\pm>1e+4$	$244.8{\pm}16.2$	$0.0005 {\pm} 0.05$	$10.0{\pm}1.2$	$2.7 {\pm} 0.19$	$-0.35 {\pm} 0.09$	$7666.0 \pm 937.2$	$80.6 {\pm} 4.3$	35.0%
ChemVLM-8B	$0.14 {\pm} 0.15$	$> 2e+5\pm > 1e+5$	$197.39 {\pm} 30.9$	$0.08 {\pm} 0.08$	<u>9.0</u> ±1.5	$2.2 \pm 0.14$	$0.22 {\pm} 0.22$	$> 2e+4\pm > 2e+4$	$62.5 \pm 11.0$	100.0%
ChemMLLM (ours)	$0.71 \pm 0.05$	$> \underline{3e+4} \pm > 1e+4$	$119.5 \pm 12.1$	$0.42 \pm 0.05$	$13.2 \pm 9.0$	$1.8 \pm 0.26$	$\textbf{0.71}{\pm}0.05$	<b>1191.6</b> ±183.6	$26.5 \pm 1.91$	65.5%

Table 4: Results on property2img task: MW, LogP and TPSA (best, 2nd best, \*: significantly better (statistically)).

		Hbd			Hba			Rb			QED	
Method	Pearson (†)	$\text{MSE}\left(\downarrow\right)$	MAE $(\downarrow)$	Pearson (†)	MSE $(\downarrow)$	$\text{MAE}\left(\downarrow\right)$	Pearson (↑)	MSE $(\downarrow)$	MAE $(\downarrow)$	Pearson (†)	MSE $(\downarrow)$	MAE $(\downarrow)$
InternVL-Chat-v1.5 LLaVA-v1.5-7B GPT-4o	-0.1±0.08 0.0±0.0 <b>0.58</b> ±0.07	7.3±3.8 <u>1.8</u> ±0.35	$1.5 \pm 0.18$ <u><math>0.86 \pm 0.07</math></u>	$\begin{array}{c} 0.09{\pm}0.08\\ 0.44{\pm}0.27\\ \underline{0.57}{\pm}0.07\end{array}$	$55.3{\pm}31.1 \\ 529.7{\pm}495.4 \\ \underline{6.2}{\pm}1.0$	3.4±0.52 9.3±5.0 <b>1.7</b> ±0.12	$\begin{array}{c} 0.06{\pm}0.08\\ \text{-}0.24{\pm}0.26\\ \underline{0.45}{\pm}0.07\end{array}$	$\begin{array}{c} 44.9{\pm}15.5\\ 5789.7{\pm}3515.2\\ \textbf{11.9}{\pm}2.1\end{array}$	$\begin{array}{c} 4.3{\pm}0.41\\ 37.5{\pm}15.7\\ \textbf{2.4}{\pm}0.17\end{array}$	0.07±0.08 0.11±0.24 <b>0.59</b> ±0.04*	$\begin{array}{c} 0.09{\pm}0.009\\ 0.10{\pm}0.03\\ \textbf{0.03}{\pm}0.003\end{array}$	$\begin{array}{c} 0.24{\pm}0.01\\ 0.25{\pm}0.04\\ \textbf{0.15}{\pm}0.008\end{array}$
ChemLLM-7B-Chat ChemVLM-8B ChemMLLM (ours)	$\begin{array}{c} \text{-}0.23{\pm}0.07\\ \underline{0.47}{\pm}0.30\\ 0.45{\pm}0.08 \end{array}$	$\begin{array}{c} 3.4{\pm}0.82\\ 6.0{\pm}0.99\\ \textbf{1.5}{\pm}0.32\end{array}$	$\begin{array}{c} 1.5{\pm}0.14\\ 1.7{\pm}0.12\\ \textbf{0.83}{\pm}0.08\end{array}$	-0.22±0.06 0.02±0.13 <b>0.66</b> ±0.06	30.3±4.3 75.6±58.4 <b>5.8</b> ±0.84	$\begin{array}{c} 4.9{\pm}0.31\\ 3.49{\pm}0.56\\ \underline{1.8}{\pm}0.13\end{array}$	-0.22±0.07 0.44±0.29 <b>0.62</b> ±0.05*	$46.9 \pm 6.5$ $25.8 \pm 4.3$ $17.3 \pm 8.1$	$5.9{\pm}0.43 \\ 3.7{\pm}0.24 \\ \underline{2.4}{\pm}0.29$	$\begin{array}{c} -0.07{\pm}0.03\\ 0.10{\pm}0.06\\ \underline{0.34}{\pm}0.08\end{array}$	$\begin{array}{c} 0.07{\pm}0.009\\ 0.08{\pm}0.007\\ \underline{0.08}{\pm}0.009 \end{array}$	$\begin{array}{c} 0.24{\pm}0.01\\ 0.24{\pm}0.01\\ \underline{0.24}{\pm}0.01\end{array}$

Table 5: Results on property2img Task: Hbd, Hba, Rb, and QED (best, 2nd best, \*: significantly better).

et al., 2023), where  $z_{k,j}$  denotes the logit in the last layer, V denotes the size of the vocabulary.

Specifically, ChemMLLM adopts Chameleon VQGAN as the image tokenizer/de-tokenizer and Chameleon-7B as the language model. The image tokenizer takes images of 256×256 resolution as input. Simultaneously, text and SMILES information passes through the text tokenizer to be converted to text tokens. The text, SMILES and image tokens are then concatenated to form a unified token sequence to feed into the LLM during training and inference.

#### 3.3 Training

252

253

257

260

261

262

263 264

265

266

268

269

273

274

275

277

281

ChemMLLM's training can be divided into two stages: (i) Mol-VQGAN training and (ii) ChemM-LLM supervised fine-tuning (SFT) training, as shown in Figure 1(e).

(i) Mol-VQGAN training. The original Chameleon VQGAN is only trained on the natural image dataset and can not discrete and reconstruct molecule images well. So, the first stage focuses on improving VQGAN's performance in encoding and decoding molecule images. Concretely, we use the well-trained VQGAN (trained on natural images) as the initialization and then fine-tune it on molecule image datasets.

(ii) ChemMLLM Supervised Fine-Tuning Training. In the second stage, we freeze the Mol-278 VQGAN and only finetune the language model on 5 downstream tasks. We utilize Lumina-mGPT (Liu et al., 2024) as training framework to train our ChemMLLM and uses Chameleon-7B as the base model. The weight related to the image tokens in the last layer will first be initialized as zero during finetuning. LLM uses the output of MolVQGAN as finetuning data, *i.e.*, the data is first pre-tokenized by Mol-VQGAN and text tokenizer into token sequences and then fed into LLM.

Task	Input	Output	Source	# train/test
molecule image captioning (img2caption)	image +text	text	chebi-20 (Edwards et al., 2022) mol-instruct (Fang et al., 2023)	70K/3K
molecule image property prediction (img2property)	image +text	text	PubChem (Kim et al., 2021)	95K/5K
image-to-SMILES conversion (img2smiles)	image +text	SMILES	PubChem (Kim et al., 2021)	95K/5K
controllable multi-objective molecule image design (property2img)	text	image	PubChem (Kim et al., 2021)	95K/5K
molecule image optimization (img2img)	image +text	image	TDC (Huang et al., 2021)	157K/17K

Table 6: Tasks and datasets.

#### 4 **Tasks and Data Curation**

In this paper, we design five vision-based chemistry research tasks, defined as follows.

(1) Molecule image captioning (img2caption) is an image-to-text task, where the models are expected to generate a caption concerning the source, functionality, structure feature and usage for each molecule image. This image-to-caption task requires models to translate molecule images into natural language descriptions, which is a process mirroring how chemists annotate experimental data. Examples for this task are shown in Table 11.

(2) Molecule image property prediction (img2property) is an image-to-text task, where models are expected to generate the value of seven different important properties for each molecule image, including molecule weight (MW), Partition Coefficient (P) of a solute between octanol and water (LogP), Topological Polar Surface Area (TPSA), Hydrogen Bond Donor (Hbd), Hydrogen Bond Acceptor (Hba), Rotatable Bond (Rb), and Quantitative Estimate of Drug-likeness (QED). More details for the properties can be found in

290

291

292

293

303

304

305

306

307

308

309

310

Model	domain	architecture	txt2txt	img2txt	txt2img	img2img
Qwen-VL-Chat	general	text tokenizer, vision encoder	$\checkmark$	$\checkmark$	×	×
InternVL-Chat-v1.5	general	text tokenizer, vision encoder	$\checkmark$	$\checkmark$	×	×
LLaVA-v1.5-7B	general	text tokenizer, vision encoder	$\checkmark$	$\checkmark$	×	×
GPT-40	general	close-sourced	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
ChemLLM-7B-Chat	chemistry	text tokenizer	$\checkmark$	×	×	×
ChemVLM-8B	chemistry	text tokenizer, vision encoder	$\checkmark$	$\checkmark$	×	×
ChemMLLM (ours)	chemistry	text tokenizer, vision tokenizer/de-tokenizer	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$

Table 7: Architectures and capabilities of MLLMs and Chemical LLMs approaches.

Appendix G. This image-to-property prediction task evaluates a model's ability to infer key chemical properties directly from 2D molecular images, enabling researchers to extract actionable insights from molecular images without specialized software, which could accelerate high-throughput screening in drug/material design (Lu et al., 2021). Table 12 shows some examples.

(3) Image-to-SMILES conversion (img2smiles)
is a fundamental chemistry task, where models are
expected to recognize the SMILES in each molecular image. The image-to-SMILES translation task
challenges models to convert 2D molecular images
into SMILES strings, requiring precise recognition
of atoms, bonds, rings, and stereochemistry. Examples for this task are shown in Table 13.

(4) Controllable multi-objective molecule image
design (property2img) is the inverse problem of
molecule image property prediction and is a text-toimage task, where models are expected to generate
the image of a molecule conditioned on target properties. It is the core of molecule design (Du et al.,
2022). The challenge lies in simultaneously optimizing multiple property constraints while maintaining chemical validity. Examples for this task
are shown in Table 14.

(5) **Molecule image optimization (img2img)** is an image-to-image task, where models take a molecular structure with less desirable molecular properties (*e.g.*, LogP) as input and generate a similar molecular structure with more desirable properties while preserving desired chemical properties. It imitates the process of lead optimization, a fundamental problem in drug discovery (Huang et al., 2021; Fu et al., 2020). Examples for this task are shown in Table 15.

## 4.1 Data Curation

338

341

342

343

345

347

349

352

We employ RDKit (Landrum et al., 2006) to convert the original SMILES strings into molecular images across all five tasks. We primarily follow the methodology of SketchMol (Wang et al., 2025) and ChemVLM (Li et al., 2025) for data curation and diversity natural language templates synthesis. The input/output modalities, raw data sources, and sizes of training/test sets for all tasks are shown in Table 6. Further details on data curation are provided in Appendix A. 353

354

355

356

357

358

359

360

361

362

363

364

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

386

387

389

## 5 Experiment

## 5.1 Experimental Setup

**Baseline Methods** cover both general multimodal LLM and chemical LLM. For general-domain multimodal LLM, we chose Qwen-VL-Chat (Bai et al., 2023), InternVL-Chat-v1.5 (Chen et al., 2023), LLaVA-v1.5-7B (Liu et al., 2023) and GPT-40 (OpenAI, 2024). For chemical LLM, we chose ChemLLM-7B-Chat (Zhang et al., 2024a) and ChemVLM-8B (Li et al., 2025). We compare their capabilities in Table 7. Please refer to Appendix D for more descriptions.

**Evaluation metrics** and **implementation details** are elaborated in Appendix F and H, respectively. The code is publicly available at https://anonymous.4open.science/r/ ChemMLLM-0D98/.

# 5.2 Result

Molecule image captioning (img2caption). We compare our model with various multimodal LLMs (MLLMs) including Qwen-VL-Chat (Bai et al., 2023), InternVL-Chat-v1.5 (Chen et al., 2023), LLaVA-v1.5-7B (Liu et al., 2023), GPT-4o (OpenAI, 2024), ChemVLM-8B (Li et al., 2025). The evaluation results are shown in Table 1. Our model exhibits strong performance on this task, outperforming all competing MLLM models on all six metrics. An example is shown in Figure 2. Our model generates captions that closely match the ground truth, while Qwen-VL-Chat includes fewer semantically informative details.

Molecule	ima	ge pro	operty	pr	ediction	390
(img2prope	rty).	GPT-40	can	not	predict	391



Figure 2: An example on img2caption task, comparison between Qwen-VL-Chat and our ChemMLLM.

properties from image directly, so in this task we do not compare with GPT-40. As shown in Table 2 and 3, our model consistently outperforms competing methods across all seven molecular properties, yielding the highest Pearson correlation coefficients alongside the lowest MSE and MAE values. An example is shown in Figure 3. Among the properties predicted, our model has 5 accurate values and 2 close values while Qwen-VL-Chat has 2 close values and 5 inaccurate values.

392

393

399

400

401

Model	Avg Sim (†)	<b>Tani@1.0</b> (†)	valid $\%(\uparrow)$	
Qwen-VL-Chat	$0.08 \pm 0.006$	$0.0\pm0.0$	8.2%	
InternVL-Chat-v1.5	$0.09 \pm 0.003$	$0.0 \pm 0.0$	20.7%	
LLaVA-v1.5-7B	$0.05 \pm 0.004$	$0.0 \pm 0.0$	11.1%	
GPT-40	$0.29 {\pm} 0.005$	$0.01 \pm 0.004$	74.5%	
ChemVLM-8B	$\underline{0.55} \pm 0.009$	$0.15 \pm 0.01$	<u>85.2</u> %	
ChemMLLM (ours)	0.75±0.009*	0.49±0.01*	97.1%	

Table 8: Results on img2smiles Task. Tanimoto similarities are written as Avg Sim, and Tanimoto@1.0 written as Tani@1.0 (**best**, 2nd best, \*: significantly better).



Figure 3: A comparison of answers on img2property task on Qwen-VL-Chat and our ChemMLLM. Accurate answers are highlighted in bottle-green, close answers are highlighted in light-green and inaccurate answers are highlighted in red. **Image-to-SMILES conversion (img2smiles).** The evaluation results are shown in Table 8. ChemMLLM performances best in both Tanimoto similarity and Tanimoto@1 metrics. For Tanimoto similarity, ChemMLLM (0.75) surpasses domainspecific model ChemVLM (0.55) by 36.4%. An example is shown in Figure 4. Our model recognizes SMILES from images successfully, while GPT-40 predicts wrong SMILES with low Tanimoto similarity.

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437



Figure 4: A comparison of GPT-4o's and ChemM-LLM's answers on img2smiles task.

Controllable multi-objective molecule image design (property2img). Since the GPT-40 API does not perform as well on this task as its web interface, and other MLLMs lack the capability to generate images, we treat this task as a purely text-based problem when evaluating other MLLMs. Specifically, our model is used to generate molecular images, while other MLLMs are tasked with directly generating the corresponding SMILES strings in text form. Given this setting, we also include ChemLLM-Chat-7B (Zhang et al., 2024a), a domain-specific chemical language model, in our evaluation. Furthermore, we exclude Owen-VL-Chat (Bai et al., 2023) from comparison on this task, as it fails to generate valid SMILES strings on all test samples. The evaluation results are shown in Table 4 and 5. Our model achieves top-2 or better performance on 90% of the evaluation metrics. Several examples are shown in Figure 5. Our model generates the image molecules with desired properties directly. For more result examples for this task, please refer to Figure 12.

**Molecule image optimization (img2img).** Since the GPT-40 API does not achieve the same level of performance on this task as its web-based interface and other MLLMs lack the ability to generate



Figure 5: Examples on property2img task on our ChemMLLM. Accurate answers are highlighted in bottle-green, close answers are highlighted in light-green.

images, we formulate the evaluation as an image-to-438 text task for these models. Specifically, our model 439 generates molecular images directly given opti-440 mized molecule images, while other MLLMs are 441 442 required to generate SMILES strings. As GPT-40 cannot directly generate SMILES with optimized 443 LogP from molecular images, we instead treat its 444 evaluation as a text-to-text task, by providing the in-445 put SMILES in textual form rather than as images. 446 As shown in Table 9, ChemMLLM achieves the 447 highest increase in LogP, outperforms GPT-40 by 448 118.9%. Several examples are shown in Figure 6. 449 Our model generates molecule images with higher 450 LogP directly. For more result examples for this 451 452 task, please refer to Figure 13.

Model	Increased LogP $(\uparrow)$	Diversity $(\uparrow)$	Novelty (†)	valid $\%(\uparrow)$
Qwen-VL-Chat InternVL-Chat-v1.5 LLaVA-v1.5-7B GPT-40	$\begin{array}{c} -2.0 \pm 0.11 \\ -0.77 \pm 0.17 \\ -0.86 \pm 0.59 \\ 1.95 \pm 0.08 \end{array}$	$\begin{array}{c} \underline{0.95}{\pm}0.01\\ 0.90{\pm}0.004\\ \textbf{0.96}{\pm}0.005\\ 0.86{\pm}0.002 \end{array}$	$1.0\pm0.0$ $1.0\pm0.0$ $1.0\pm0.0$ $1.0\pm0.0$	4.0% 48.0% 37.5% <b>99.0</b> %
ChemVLM-8B ChemMLLM (ours)	$\frac{0.45}{4.2\pm0.14}$	$0.87 {\pm} 0.002$ $0.88 {\pm} 0.001$	0.97±0.01 <b>1.0</b> ±0.0	<u>92.5</u> % 91.0%

Table 9: Results on img2img task (best, 2nd best).

### 5.3 Ablation Study

453

454

455

456

We conduct an ablation study on property2img task (see Appendix I) to assess the effects of Mol-VQGAN training and data augmentation (espe-



Figure 6: Examples of ChemMLLM on img2img task.

cially image rotation). Results show that both components significantly improve the correlation between generated images and molecular properties, with their combination yielding the best overall performance. This highlights the importance of high-quality visual representations in enhancing multimodal chemical tasks.

457

458

459

460

461

462

463

464

## 6 Conclusion

This paper has proposed ChemMLLM, a chemi-465 cal multimodal large language model that handles 466 molecule comprehension and generation across 467 three modalities (text, SMILES string, molecule 468 image). By jointly modeling text, SMILES strings, 469 and molecule images, ChemMLLM enables seam-470 less cross-modal comprehension and generation, 471 outperforming state-of-the-art MLLMs and special-472 ized chemical LLMs across a range of tasks. Also, 473 we design five cross-modal chemistry tasks and cu-474 rate datasets, providing a valuable resource for mul-475 timodal AI in chemistry. The experimental results 476 demonstrate ChemMLLM's strong performance, 477 highlighting its potential for real-world drug and 478 material discovery. 479

# 530 531 532 533 534 535 536 537 538 539 540 541 542 543 544 545 546 547 548 549 550 551 552 553 554 555 556 557 558 559 561 562 563 564 565 566 567 568 569 570 571 572

573

574

575

576

577

578

579

580

581

582

583

584

# Limitations

480

497

511

512

513

514

515

516

517

518

519

521

522

524

525

528

529

Despite its promising capabilities, our work has 481 several limitations that point to important direc-482 tions for future research: Currently, ChemMLLM 483 incorporates only three modalities - text, SMILES 484 strings, and 2D molecule images. Real-world 485 chemical data often includes richer modalities such 486 as 3D molecular structures, quantum mechanical 487 properties, or spectroscopic data. Incorporating 488 these would significantly enhance the model's abil-489 ity to capture complex molecular behaviors and 490 interactions. Also, our evaluation is primarily fo-491 cused on proof-of-concept chemistry tasks. Fur-492 ther studies are needed to validate the model's per-493 formance in real-world applications such as drug 494 discovery, materials design, or chemical synthesis 495 planning. 496

## Ethics Statement

The development and application of chemical AI 498 models, such as ChemMLLM, raise important eth-499 ical considerations. We ensure that all datasets 500 used in this work are sourced from publicly avail-501 able, non-sensitive chemical data and do not in-502 volve personal or private information. The model 503 is designed for scientific research purposes, including drug discovery and materials science, with the aim of advancing chemical understanding and innovation. We advocate for responsible use of AI in chemistry and encourage transparency, reproducibility, and fairness in future deployments of such technologies. 510

## References

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. *arXiv preprint arXiv:2204.14198*.
- Eric Anderson, Gilman D Veith, and David Weininger. 1987. *SMILES, a line notation and computerized interpreter for chemical structures.* US Environmental Protection Agency, Environmental Research Laboratory.
- Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. 2023. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*.
- Jinze Bai, Shuai Bai, Shusheng Yang, Sinan Wang, Zhitang Tan, Penghui Wang, Junyang Chen, Jun Zhou,

and Jingren Zhou. 2023. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv* preprint arXiv:2308.12966.

- Steven Bird, Ewan Klein, and Edward Loper. 2009. Natural language processing with Python: analyzing text with the natural language toolkit. " O'Reilly Media, Inc.".
- Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, et al. 2024. Internlm2 technical report. *arXiv preprint arXiv:2403.17297*.
- He Cao, Zijing Liu, Xingyu Lu, Yuan Yao, and Yu Li. 2023. Instructmol: Multi-modal integration for building a versatile and reliable molecular assistant in drug discovery. *arXiv preprint arXiv:2311.16208*.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. 2023. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238*.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.
- Yuanqi Du, Tianfan Fu, Jimeng Sun, and Shengchao Liu. 2022. Molgensurvey: A systematic survey in machine learning models for molecule design. *arXiv preprint arXiv:2203.14500*.
- Carl Edwards, Tuan Lai, Kevin Ros, Garrett Honke, Kyunghyun Cho, and Heng Ji. 2022. Translation between molecules and natural language. *arXiv preprint arXiv:2204.11817*.
- Patrick Esser, Robin Rombach, and Bjorn Ommer. 2021. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12873–12883.
- Yin Fang, Xiaozhuan Liang, Ningyu Zhang, Kangwei Liu, Rui Huang, Zhuo Chen, Xiaohui Fan, and Huajun Chen. 2023. Mol-instructions: A large-scale biomolecular instruction dataset for large language models. *arXiv preprint arXiv:2306.08018*.
- Tianfan Fu, Cao Xiao, and Jimeng Sun. 2020. CORE: Automatic molecule optimization using copy and refine strategy. *AAAI*.
- Kexin Huang, Tianfan Fu, Wenhao Gao, Yue Zhao, Yusuf Roohani, Jure Leskovec, Connor W Coley, Cao Xiao, Jimeng Sun, and Marinka Zitnik. 2021. Therapeutics data commons: machine learning datasets and tasks for therapeutics. *NeurIPS Track Datasets and Benchmarks*.

585 Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Yizhen Luo, Kai Yang, Massimo Hong, Xing Yi Liu, Zikun Nie, Hao Zhou, and Zaiqing Nie. 2024. Learn-Efros. 2017. Image-to-image translation with condi-641 tional adversarial networks. In Proceedings of the ing multi-view molecular representations with struc-642 IEEE conference on computer vision and pattern tured and unstructured knowledge. In Proceedings of 643 589 recognition, pages 1125–1134. the 30th ACM SIGKDD Conference on Knowledge 644 Discovery and Data Mining, pages 2082–2093. 645 Dhiraj Kalamkar, Dheevatsa Mudigere, Naveen Mellempudi, Dipankar Das, Kunal Banerjee, Sasikanth Avan-OpenAI. 2023. Gpt-4v(ision) system card. Accessed: 646 cha, Dharma Teja Vooturi, Nataraj Jammalamadaka, 2024-07-20. 647 Jianyu Huang, Hector Yuen, et al. 2019. A study of 594 bfloat16 for deep learning training. arXiv preprint OpenAI. 2024. Gpt-4o: Our most advanced ai model. 648 595 arXiv:1905.12322. Accessed: 2024-07-20. 649 Sunghwan Kim, Jie Chen, Tiejun Cheng, Asta Gindu-Yujie Qian, Jiang Guo, Zhengkai Tu, Zhening Li, Con-650 lyte, Jia He, Siqian He, Qingliang Li, Benjamin A nor W. Coley, and Regina Barzilay. 2023. MolScribe: 651 Shoemaker, Paul A Thiessen, Bo Yu, et al. 2021. Pub-Robust molecular structure recognition with image-652 chem in 2021: new data content and improved web to-graph generation. Journal of Chemical Informa-653 interfaces. Nucleic acids research, 49(D1):D1388tion and Modeling. 654 D1395. Diederik P Kingma and Jimmy Ba. 2014. Adam: A Rico Sennrich, Barry Haddow, and Alexandra Birch. 655 method for stochastic optimization. arXiv preprint 2015. Neural machine translation of rare words with 656 604 arXiv:1412.6980. subword units. arXiv preprint arXiv:1508.07909. 657 Yana Kosenkov and Dmitri Kosenkov. 2021. Computer OpenEye Scientific Software. 2023. OpenEye Toolkits 658 vision in chemistry: Automatic titration. Documentation. 659 Greg Landrum et al. 2006. Rdkit: Open-source chemin-Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, 660 formatics. Bingyue Peng, Ping Luo, and Zehuan Yuan. 661 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2024. Autoregressive model beats diffusion: Llama 662 2023. Blip-2: Bootstrapping language-image prefor scalable image generation. arXiv preprint 663 arXiv:2406.06525. training with frozen image encoders and large lan-664 guage models. In International conference on machine learning, pages 19730-19742. PMLR. 613 Chameleon Team. 2024. Chameleon: Mixed-modal 665 early-fusion foundation models. arXiv preprint 666 Junxian Li, Di Zhang, Xunzhi Wang, Zeying Hao, Jingdi arXiv:2405.09818. 667 Lei, Qian Tan, Cai Zhou, Wei Liu, Yaotian Yang, Xinrui Xiong, et al. 2025. Chemvlm: Exploring the Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-668 power of multimodal large language models in chem-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan 669 istry area. In Proceedings of the AAAI Conference Schalkwyk, Andrew M Dai, Anja Hauth, Katie 670 on Artificial Intelligence, volume 39, pages 415-423. Millican, et al. 2023. Gemini: a family of 671 highly capable multimodal models. arXiv preprint 672 Chin-Yew Lin. 2004. Rouge: A package for automatic arXiv:2312.11805. 673 evaluation of summaries. In Text summarization branches out, pages 74-81. Aaron Van Den Oord, Oriol Vinyals, et al. 2017. Neural 674 Dongyang Liu, Shitian Zhao, Le Zhuo, Weifeng Lin, discrete representation learning. Advances in neural 675 Yu Qiao, Hongsheng Li, and Peng Gao. 2024. information processing systems, 30. 676 Lumina-mgpt: Illuminate flexible photorealistic textto-image generation with multimodal generative pre-Xinlong Wang, Xiaosong Zhang, Zhengxiong Luo, 677 training. arXiv preprint arXiv:2408.02657. Quan Sun, Yufeng Cui, Jinsheng Wang, Fan Zhang, 678 Yueze Wang, Zhen Li, Qiying Yu, et al. 2024. Emu3: 679 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Next-token prediction is all you need. arXiv preprint 680 Lee. 2023. Visual instruction tuning. Advances in arXiv:2409.18869. 681 neural information processing systems, 36:34892-34916. Zixu Wang, Yangyang Chen, Pengsen Ma, Zhou Yu, 682 Jianmin Wang, Yuansheng Liu, Xiucai Ye, Tetsuya 683 Ilva Loshchilov and Frank Hutter. 2017. Decou-Sakurai, and Xiangxiang Zeng. 2025. Image-based 684 pled weight decay regularization. arXiv preprint generation for molecule design with sketchmol. Na-685 arXiv:1711.05101. ture Machine Intelligence, pages 1–12. 686 Yingzhou Lu, Chiung-Ting Wu, Sarah J Parker, Lulu David Weininger. 1988. Smiles, a chemical language Chen, Georgia Saylor, Jennifer E Van Eyk, David M 687 Herrington, and Yue Wang. 2021. COT: an efficient and information system. 1. introduction to methodol-688 python tool for detecting marker genes among many ogy and encoding rules. Journal of chemical infor-689 subtypes. *bioRxiv*, pages 2021–01. mation and computer sciences, 28(1):31-36. 690

586

593

596

598

599

607

611

612

614

615

616

618

619

621

622

623

631

632

635

Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. 2024. Show-o: One single transformer to unify multimodal understanding and generation. *arXiv* preprint arXiv:2408.12528.

694

697

702

704

705

706

707 708

710

711

712

713

714

715

716

717

718

719 720

721

722

725

- Di Zhang, Wei Liu, Qian Tan, Jingdan Chen, Hang Yan, Yuliang Yan, Jiatong Li, Weiran Huang, Xiangyu Yue, Wanli Ouyang, et al. 2024a. Chemllm: A chemical large language model. *arXiv preprint arXiv:2402.06852*.
- Juzheng Zhang, Yatao Bian, Yongqiang Chen, and Quanming Yao. 2024b. Unimot: Unified moleculetext language model with discrete token representation. *arXiv preprint arXiv:2408.00863*.
- Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. 2018. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595.
- Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo, Chien-Chin Huang, Min Xu, Less Wright, Hamid Shojanazeri, Myle Ott, Sam Shleifer, et al. 2023. Pytorch fsdp: experiences on scaling fully sharded data parallel. *arXiv preprint arXiv:2304.11277*.
- Chunting Zhou, Lili Yu, Arun Babu, Kushal Tirumala, Michihiro Yasunaga, Leonid Shamis, Jacob Kahn, Xuezhe Ma, Luke Zettlemoyer, and Omer Levy. 2024. Transfusion: Predict the next token and diffuse images with one multi-modal model. *arXiv preprint arXiv:2408.11039*.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*.

# **A** Data Curation Details

(1) **img2caption**: The dataset used for this task is sourced from chebi-20 (Edwards et al., 2022)and mol-instruct (Fang et al., 2023). The original datasets contain SMILES-caption pairs. We utilize RDKit (Landrum et al., 2006) to transfer the SMILES into 256×256 image to form imagecaption pairs. For mol-instruct, we filter captions shorter than 150 words to screen clearer descriptions. We only use the test set of chebi-20 as the test set and the partitioning of the dataset is the same as chebi-20 (Edwards et al., 2022). 726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

773

774

(2) **img2property**: The dataset used for this task is sourced from PubChem (Kim et al., 2021). For one compound, there are three fileds related to it's SMILES, *i.e.*, "PUBCHEM\_SMILES" which means The Simplified Molecular Linear Input Specification (SMILES) for compounds, it is a string used to represent the chemical structure, "PUBCHEM\_OPENEYE\_CAN\_SMILES" which means Canonical SMILES generated using the OpenEye tool (Software, 2023) and "PUB-CHEM\_OPENEYE\_ISO\_SMILES" which means The isomer SMILES generated using the OpenEye tool. The processing steps are as follows:

- 1. Extract all the three fields from PubChem and remove the duplicate SMILES.
- 2. Use the Draw.MolToImage() function in RD-Kit (Landrum et al., 2006) to transfer SMILES into images.
- 3. Sample 100,000 SMILES for this work.
- 4. Choose 7 important properties as the prediction objective, *i.e.*, MW, LogP, TPSA, HBD, HBA, RB and QED. Then use RDKit to calculate the 7 properties for each sampled SMILES.
- 5. Use natural language templates to integrate properties into natural language to form the final image-property answer pairs.
- 6. Divide the dataset into training and test set by the ratio of 95:5.

(3) img2smiles: The dataset and the construction steps are the same as the img2property task, the difference is that this task only apply templates for SMILES to construct image-SMILES answer pairs.
(4) property2img: This task is a reverse task of img2property. By swapping the question and answer of the img2property dataset, we construct property prompt-image pairs for property2img task.

(5) **img2img**: The dataset for this task is sourced from TDC (Huang et al., 2021). The original dataset contains SMILES-SMILES pairs. The previous molecule's LogP is lower while the latter molecule's LogP is higher. We also use Draw.MolToImage() function in RDKit (Landrum et al., 2006) to transfer them into image-image pairs. For data partitioning, we divide the training set and test set by the ratio of 9:1.

All the data used in this paper is publicly available.

# **B** Mathematical Notations

For ease of understanding, we list key mathematical notations in Table 10 in Appendix.

## C Data Examples

775

776

777

778

780

781

784

790

792

793

794

805

806

807

(1) **Molecule image captioning (img2caption).** As shown in Table 11, the input/output for img2caption task is text/image-text pair. The input is a question asking models to give captions for molecule images and the output is a caption concerning the source, functionality, structure feature and usage for each molecule image.

(2) Molecule image property prediction (img2property). As shown in Table 12, the input/output for img2property task is text/imagetext pair. The input is a question asking models to predict seven properties for given molecule images and the output is a natural language answer describing the seven properties.

(3) **Image-to-SMILES conversion (img2smiles).** As shown in Table 13, the input/output for img2smiles task is text/image-text/SMILES pair. The input is a question asking models to recognize the SMILES in the given molecule images and the output is a natural language answer describing the SMILES in the image.

(4) Controllable multi-objective molecule image design (property2img). As shown in Table 14, the input/output for property2img task is
text-text/image pair. The input is a question asking
models to generate images according to the given
values of the seven properties, and the output is the
generated molecule image with the given values of
the seven properties.

(5) Molecule image optimization (img2img). As
shown in Table 15, the input/output for img2img
task is text/image-text/image pair. The input is
a question describing the meaning of LogP and
then asking models to optimize the LogP property

for given molecule images, and the output is the optimized molecule image with better LogP.

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

# **D** Baseline Methods

• Qwen-VL-Chat (Bai et al., 2023) is an opensource multimodal conversational model developed by Alibaba Cloud, extending the Qwen-VL architecture to support complex visual-language interaction through instruction tuning. It integrates a frozen CLIP-ViT-G/14 vision encoder with the Qwen-7B language model via a trainable vision-language connector. Images are encoded into 1024-dimensional patch embeddings by the vision encoder, which are then linearly projected into the token embedding space of the language model (hidden size 4096). The language backbone consists of 32 transformer decoder layers, each employing multi-head masked self-attention with rotary position embeddings (RoPE), followed by a SwiGLU-activated feedforward network with an intermediate dimension of 11008. Owen-VL-Chat uses a 2048-token context window and supports dynamic multimodal prompts comprising text, images, and region boxes. It is instruction-tuned on a large-scale, GPT-generated multimodal dataset containing both single- and multi-turn visual conversations, enabling capabilities in visual question answering, dense captioning, document OCR, and multiimage reasoning. The model achieves high performance on benchmarks such as MME, SEED-Bench, and MMBench, demonstrating strong alignment between visual and linguistic modalities.

• InternVL-Chat-V1.5 (Chen et al., 2023) is an open-source vision-language instructionfollowing model developed by OpenGVLab, designed to support high-resolution, multilingual, and multi-turn visual conversations. The model integrates a powerful ViT-based vision encoder, InternViT-6B, with the InternLM2-Chat-20B language model via a trainable multi-layer perceptron (MLP) connector. InternViT encodes images into patch embeddings with dynamic resolution support, allowing the model to process up to 40 image tiles of size 448×448, effectively supporting 4K-level inputs. These embeddings are projected into the language model's token space to enable seamless multimodal interaction. The language backbone consists of 64 transformer decoder layers with rotary posi-

Table 10: Mathematical notations.

Notations	Descriptions
$x/\hat{x} \in \mathbb{R}^{H \times W \times 3}$	the input/reconstructed molecule image
H/W	the height/width of the input image
h/w	the height/width of the feature map
$n_z$	the channels of the feature map, same as the dimension of the codebook vector
E	the encoder of VQGAN, a convolutional neural network (CNN) that extracts
	features from original image
$\hat{z} \in \mathbb{R}^{h \times w \times n_z}$	the continuous feature map encoded by $E(x)$
$\hat{z}_{ij} \in \mathbb{R}^{n_z}$	the spatial code $\in \mathbb{R}^{n_z}$ at position $i, j$ in the feature map; $(i, j) \in \mathbb{R}^{n_z}$
	$\{0, 1, \dots, h\} \times \{0, 1, \dots, w\}$
q	vector quantization process in VQGAN, which transfers continuous feature into
	discrete feature
$Z = \{z_i\}_{i=1}^n \subset \mathbb{R}^{n_z}$	the codebook in VQGAN, a dictionary that represents the latent discrete space
$z_k \in \mathbb{R}^{n_z}$	entry in the codebook
$z_q \in \mathbb{R}^{h \times w \times n_z}$	the quantized feature map quantified by $\mathbf{q}(\hat{z})$
G	the decoder of VQGAN, a CNN that reconstructs image from latent discrete
	space
$sg[\cdot]$	the stop-gradient operation
$\mathcal{L}_{vqvae}$	origin VQVAE training loss
$\mathcal{L}_{rec}$	the reconstruction loss
$\mathcal{L}_{perceptual}$	the perceptual loss
$\mathcal{L}_{GAN}$	GAN loss
D	the discriminator to identify $x$ and $\hat{x}$ , a patch-based discriminator (Isola et al.,
0	
$\mathcal{L}_{LLM}$	next-token prediction cross-entropy loss with z-loss for training large language
	model the matchedility of a given a
$p_{\theta}(s_i s_1,\ldots,s_{i-1})$	the probability of $s_i$ given $s_1, \ldots, s_{i-1}$
$S_T$	text/SMILES token sequence tokenized by text tokenizer
ST	the concatenated sequence of image token sequence and text/SMILES token
D	sequence
L	the size of total sequence tokens
	the size of vocabulary
$z_{t-i} \in \mathbb{R}$	the logit at last layer
$\lambda_{\kappa,j} \subset \mathbb{R}$	hyper-parameters or adaptively calculated parameters to adjust the weight of
~, ~1, ~2	different loss functions
	1

tional embeddings (RoPE), multi-head masked 874 875 876 877 882 883 884

self-attention, and SwiGLU-activated feedforward layers. The model uses a maximum context length of 4096 tokens and is instruction-tuned on a high-quality bilingual dataset containing document images, natural images, and complex multimodal dialogues. InternVL-Chat-V1.5 achieves strong performance on benchmarks including MME, MMBench, and AI2D, demonstrating robust capabilities in visual question answering, document OCR, visual reasoning, and bilingual understanding. The total parameter count is approximately 25.5 billion, with both the vision encoder and language model jointly fine-tuned during the instruction-following phase.

• LLaVA-v1.5-7B (Large Language and Vision

Assistant) (Liu et al., 2023) is an open-source vision-language instruction-tuned model that integrates a pre-trained CLIP-ViT-L/14 vision encoder with the LLaMA-7B language model through a projection network. LLaVA-7B processes visual inputs by encoding images into 1024-dimensional patch embeddings via the vision encoder, which are then projected into the language model's token space through a trainable linear layer (hidden size 4096). The language model comprises 32 transformer decoder blocks, each with masked self-attention (key/query size 4096, 32 attention heads), followed by a feed-forward network with SwiGLU activation and an intermediate dimension of 11008. The self-attention layers use RoPE (rotary positional

input	output
Give an overview of this	The molecule is a steroid ester that is methyl (17E)-pregna-4,17-dien-21-
molecule's properties.	oate substituted by oxo groups at positions 3 and 11. It is a 3-oxo-Delta(4) steroid, an 11-oxo steroid, a steroid ester and a methyl ester. It derives
A the second sec	from a hydride of a pregnane.
Explain the key traits of this molecule.	The molecule is an epoxy(hydroxy)icosatrienoate that is the conjugate base of 11 hydroxy-(14R,15S)-epoxy-(5Z,8Z,12E)-icosatrienoic acid.

Provide a detailed description of this molecular structure.

a soft

List the notable attributes of this molecule.



906

907

908

909

910

911

912

913

The molecule is an epoxy(hydroxy)icosatrienoate that is the conjugate base of 11 hydroxy-(14R,15S)-epoxy-(5Z,8Z,12E)-icosatrienoic acid, obtained by deprotonation of the carboxy group; major species at pH 7.3. It is a conjugate base of an 11 hydroxy-(14R,15S)-epoxy-(5Z,8Z,12E)-icosatrienoic acid.

The molecule is a tetracyclic triterpenoid that is 4,4,8-trimethylandrosta-1,14-diene substituted by an oxo group at position 3, an acetoxy group at position 7 and a furan-3-yl group at position 17. Isolated from Azadirachta indica, it exhibits antiplasmodial and antineoplastic activities. It has a role as an antineoplastic agent, an antiplasmodial drug and a plant metabolite. It is an acetate ester, a cyclic terpene ketone, a member of furans, a limonoid and a tetracyclic triterpenoid.

The molecule is a member of the class of N-nitrosoureas that is urea in which one of the nitrogens is substituted by methyl and nitroso groups. It has a role as a carcinogenic agent, a mutagen, a teratogenic agent and an alkylating agent.

### Table 11: Example for img2caption task.

embeddings) and support a context window of 2048 tokens. LLaVA-7B leverages instructiontuning on 558K GPT-4 generated multimodal instruction-following samples, aligning visual and textual representations for tasks such as visual QA and image captioning. The total number of trainable parameters is approximately 7 billion, with the vision encoder frozen during fine-tuning.

• GPT-40 (OpenAI, 2024) (Generative Pre-trained 914 Transformer 4 Omni) is a state-of-the-art multi-915 modal foundation model developed by OpenAI, 916 designed to natively process and reason across 917 text, images, and audio modalities. Unlike previ-918 ous GPT-4 variants that rely on separate vision 919 encoders, GPT-40 employs a unified transformer 920 architecture that jointly encodes multimodal in-921 922 puts, enabling low-latency and high-fidelity interactions. The model supports up to 128k tokens of context and exhibits strong performance 924 across a wide range of tasks, including natural language understanding, image captioning, docu-926

ment analysis, and spoken language comprehension. GPT-40 achieves significant improvements in visual reasoning (*e.g.*, charts, diagrams, OCR), math problem solving, and multilingual capability, surpassing the capabilities of GPT-4-turbo while operating with lower inference latency. It is instruction-tuned on a diverse and extensive corpus of text and multimodal data, and optimized for both conversational fluency and factual grounding. We also utilized GPT-40 to assist in the writing of this paper. 927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

• ChemLLM-7B-Chat (Zhang et al., 2024a) is an open-source domain-specific large language model framework tailored for chemical sciences, designed to address the limitations of generalpurpose LLMs in structured scientific domains. The model is instruction-tuned using ChemData, a templated dataset that transforms structured chemical knowledge (*e.g.*, molecular properties, SMILES strings, compound databases) into natural language instructions across a variety of chem-

input	output
Predict the following properties of the molecule:	The MW is 1022.34, LogP is 4.8, TPSA is 304.6,
MW, LogP, TPSA, HBD, HBA, RB and QED.	HBD is 9, HBA is 11, RB is 24 and QED is 0.05

Predict the following properties of the molecule: The MW is 305.42, LogP is 5.42, TPSA is 25.5, HBD MW, LogP, TPSA, HBD, HBA, RB and QED. is 0, HBA is 2, RB is 3 and QED is 0.63

MW, LogP, TPSA, HBD, HBA, RB and QED. is 0, HBA is 2, RB is 1 and QED is 0.64

Predict the following properties of the molecule: The MW is 181.3, LogP is 2.56, TPSA is 3.24, HBD

MW, LogP, TPSA, HBD, HBA, RB and QED. HBD is 1, HBA is 4, RB is 7 and QED is 0.49

Predict the following properties of the molecule: The MW is 479.99, LogP is 4.88, TPSA is 61.88,

Table 12: Example for img2property task.

ical tasks. ChemLLM adopts a decoder-only transformer architecture based on InternLM2-Base-7B (Cai et al., 2024), and is fine-tuned to handle both single- and multi-turn dialogues. The model supports tasks such as molecular property prediction, compound generation, synthetic route planning, and reaction condition recommendation.

• ChemVLM-8B is an open-source multimodal domain-specific large language model specifically designed for chemistry-related tasks, aiming to bridge the gap between vision and language understanding in the chemical domain. The model adopts a ViT-MLP-LLM architecture, integrating a vision encoder based on Vision Transformer (ViT), a multi-layer perceptron (MLP) as a projection layer, and a 20billion-parameter decoder-only language model 965 ChemLLM-20B (Zhang et al., 2024a) as the backbone. ChemVLM processes visual inputs such as molecular structures, chemical reaction schemes,

and spectra by encoding images into patch em-969 beddings through ViT, which are then linearly 970 projected to the token space of the language 971 model. It is instruction-tuned on a constructed 972 dataset of 1.2M multimodal samples covering 973 tasks like molecule captioning, reaction classifi-974 cation, and chemical structure understanding. 975

976

977

978

979

980

981

982

983

984

985

986

987

#### **More Visual Result** Ε

To better visualize results for different run on five tasks, we draw metric bar for each tasks. The metric bar for img2caption, img2property, img2smiles, property2img, img2img task is shown in Figure 7, Figure 8, Figure 10, Figure 11, Figure 9, respectively.

For the two image generation task, we provide more examples to better display the ability of ChemMLLM. More examples for property2img task are shown in Figure 12. More examples for img2imgtask are shown in Figure 13.

951

953

955

957

960

961

962

963 964

input	output
What is the SMILES of the molecule in this image?	The SMILES of the molecule in the image is C1=CC
galag-	=C(C=C1)COC(=O)/C=C/C2=CC=C(O2)[N+](=O) [O-].
What is the SMILES of the molecule in this image?	The SMILES of the molecule in the image is CC1=
	CC2=C(C=C1C)OC(=O)C=C2CN3CCC[C@@H] 3C4=CC=CN4C.
What is the SMILES of the molecule in this image?	The SMILES of the molecule in the image is CC(C)( C)OC(=O)C1CC(C(N1C(=O)CNC(=O)NC2=CC=
office	CC(=C2)C(=O)O)C3=CC=CC=C3)S(=O)(=O)C4= CC=CC=C4.
What is the SMILES of the molecule in this image?	The SMILES of the molecule in the image is COC1 =CC(=C(C=C1)C=N/N=C(\\N)/NO)OC.

Table 13: Example for img2smiles task.

MH, CH

input	output			
Generate an image of 256x256 according to the fol-	The molecular diagram is shown below:			
lowing prompt:\n Image of a molecule whose MW				
is 354.46, LogP is 3.91, TPSA is 81.57, HBD is 2,				
HBA is 5, RB is 6 and QED is 0.66	CD-CD			
Generate an image of 256x256 according to the fol-	See the molecular depiction:			
lowing prompt:\n Image of a molecule whose MW				
18 461.36, LogP 18 2.34, TPSA 18 84.91, HBD 18 3,	corre			
HDA IS 5, KB IS 6 and QED IS 0.62	·· ,			
Generate an image of 256x256 according to the fol-	The molecular diagram is shown below:			
lowing prompt. Image of a molecule whose MW is	The molecular diagram is shown below.			
353.4. LogP is 1.37. TPSA is 90.85. HBD is 0. HBA				
is 9, RB is 2 and QED is 0.45				
	×2 /			
Generate an image of 256x256 according to the fol-	See the molecular depiction:			
lowing prompt:\n Image of a molecule whose MW is	L.			
594.8, LogP is 6.15, TPSA is 76.15, HBD is 0, HBA	A m			
is 6, RB is 13 and QED is 0.19	A tro			

Table 14: Example for property2img task.

## input

LogP (Partition Coefficient) measures a molecule's solubility in fats versus water by quantifying its distribution between octanol (fat-like) and water phases. Calculated as the logarithm of the concentration ratio (LogP = log[octanol]/[water]), it predicts drug absorption and permeability—higher values (>0) indicate greater fat solubility, while lower values (<0) suggest water solubility. Ideal drug candidates typically have LogP between 0-3 for optimal bioavailability. Here is an image of a molecule, please generate an image of a new similar molecule whose LogP is better.

LogP (Partition Coefficient) measures a molecule's solubility in fats versus wa-

ter by quantifying its distribution between octanol (fat-like) and water phases....

LogP (Partition Coefficient) measures a molecule's solubility in fats versus

water by quantifying its distribution between octanol (fat-like) and water

phases....Here is an image of a molecule, please generate an image of a new

LogP (Partition Coefficient) measures a molecule's solubility in fats versus

water by quantifying its distribution between octanol (fat-like) and water

phases....Here is an image of a molecule, please generate an image of a new

Here is a new similar molecule with better LogP.

output

Here is a new similar molecule with better LogP.



Here is a new similar molecule with better LogP.



Here is a new similar molecule with better LogP.

yord

\_\_\_\_\_

# **F** Evaluation Metrics

(1) img2caption: We use BLEU-2/4, ROUGE-1/2/L, and METEOR to evaluate the quality of generated captions against reference texts; (2) img2property: We extract seven molecular properties from LLM-generated outputs and evaluate the accuracy using Mean Squared Error (MSE), Mean Absolute Error (MAE), and Pearson correlation; (3) img2smiles: We extract SMILES strings from the LLM outputs and adopt Tanimoto similarity and

similar molecule whose LogP is better.

similar molecule whose LogP is better.

Tanimoto hit 1.0 (tanimoto@1.0) that measures the 998 percentage of exact matches (similarity = 1.0); (4) 999 property2img: Generated molecular images are converted to SMILES via MolScribe (Qian et al., 1001 2023), from which properties are computed and 1002 evaluated using MSE, MAE, and Pearson correla-1003 tion. Each model is run five times, and the best 1004 result is reported; (5) img2img: We use Increased 1005 LogP, as well as molecular diversity and novelty 1006 to measure the optimized molecule. Other settings 1007

Here is an image of a molecule, please generate an image of a new similar molecule whose LogP is better.

 $\langle 1 | 1 \rangle$ 

989

Table 15: Example for img2img task.



Figure 7: Metric bar for different runs of img2caption task.



Figure 8: Metric bar for different runs of img2property task. pc\_mw means the Pearson correlation between the predicted molecule weight and groundtruth and pc\_tpsa means the Pearson correlation between the predicted topological polar surface area and groundtruth.

are the same as property2img.

1008

1009

1010

1012

1013

1014

1015

1016

1017

1018

1019

1021

The detailed explanations of each metric are as follows:

• **BLEU-N** (Bilingual Evaluation Understudy) is an automatic evaluation metric for machinegenerated text that assesses how closely a candidate sentence matches one or more reference sentences. It uses the modified precision of n-grams up to length N. It is defined as

BLEU-N = BP · exp 
$$\left(\frac{1}{N}\sum_{n=1}^{N}\log p_n\right)$$
, (5)

where  $p_n$  denotes the modified n-gram precision for n-grams of size n, and BP is the brevity penalty, which penalizes short candidate sentences to prevent artificially high



Figure 9: Metric bar for different runs of img2imgtask. LogP Improve means Increased LogP, which is the increase in LogP of the optimized molecule relative to the original molecule.



Figure 10: Metric bar for different runs of img2smiles task.

scores. The BLEU-N score ranges from 0 to 1, where a higher score indicates better overlap with the reference text in terms of n-gram content. A higher BLEU-N value generally reflects better fluency and adequacy in the generated text. In our experimental evaluation, we use the word\_tokenize() function from the NLTK library (Bird et al., 2009) to do tokenization and employ the sentence\_bleu() metric with uniform weights for all n-gram precision calculations (*i.e.*, equal weights for 1- to 4-gram contributions).

1022

1023

1024

1025

1026

1027

1028

1030

1031

1032

1033

 ROUGE-N (Recall-Oriented Understudy for Gisting Evaluation) is a recall-based metric
 that measures the overlap of n-grams between
 a candidate text and one or more reference



Figure 11: Metric bar for different runs of property2img task. pc\_tpsa means the Pearson correlation between the predicted topological polar surface area and ground truth, and pc\_rb means the Pearson correlation between the predicted rotatable bond and ground truth.

texts. It is defined as

$$\operatorname{ROUGE-N} = \frac{\sum_{S \in \{\operatorname{Ref}\}} \sum_{\operatorname{gram}_n \in S}}{\sum_{S \in \{\operatorname{Ref}\}} \sum_{\operatorname{gram}_n \in S}} \quad (6)$$
$$\frac{\operatorname{Count}_{\operatorname{match}}(\operatorname{gram}_n)}{\operatorname{Count}(\operatorname{gram}_n)},$$

where *n* denotes the length of the n-grams (*e.g.*, ROUGE-1 for unigrams, ROUGE-2 for bigrams), and Count<sub>match</sub>(gram<sub>n</sub>) is the number of n-grams in the reference that also appear in the candidate text. ROUGE-N values range from 0 to 1 and a higher ROUGE-N value indicates better performance. In our experimental evaluation, we employ the rouge\_scorer() metric from the rouge\_score (Lin, 2004) library.

• **ROUGE-L** evaluates the quality of generated text by measuring the longest common subsequence (LCS) between the candidate and reference texts. It is defined as

$$\text{ROUGE-L} = \frac{(1+\beta^2) \cdot \text{Precision} \cdot \text{Recall}}{\text{Recall} + \beta^2 \cdot \text{Precision}},$$
(7)

where Precision =  $\frac{LCS(X,Y)}{|X|}$  and Recall =  $\frac{LCS(X,Y)}{|Y|}$ , with X and Y denoting the candidate and reference sequences, respectively.  $\beta$  is typically set to favor recall ( $\beta = 1.2$  in common settings). ROUGE-L scores range from 0 to 1, with higher values indicating better preservation of the reference's sequence and structure.

METEOR (Metric for Evaluation of Translation with Explicit ORdering) is a metric designed to evaluate the quality of machine-generated text by aligning it to one or more reference texts. It is defined as

$$METEOR = F_{mean} \cdot (1 - Penalty), \quad (8)$$

where  $F_{\text{mean}} = \frac{10 \cdot P \cdot R}{R+9P}$ , with *P* and *R* denoting unigram precision and recall, respectively. The penalty is a function of the number of chunks in the alignment, designed to penalize disordered matches. METEOR scores range from 0 to 1, with higher scores indicating better alignment with the reference text. In our experimental evaluation, we perform tokenization using the word\_tokenize() function and employ the METEOR metric (meteor\_score()), both implemented in the NLTK library (Bird et al., 2009).

• Mean Squared Error (MSE) measures the average of the squares of the difference between the forecasted value and the actual value. It is defined as

MSE = 
$$\frac{1}{N} \sum_{i=1}^{N} (y_i - \hat{y}_i)^2$$
, (9)

where N is the size of the test set;  $y_i$  and  $\hat{y}_i$ denote the ground truth and predicted score of the *i*-th data sample in the test set, respectively. MSE value ranges from 0 to positive infinity. A lower MSE value indicates better performance.

• Mean Absolute Error (MAE) measures the absolute value of the difference between the predicted value and the actual value. It is defined as

$$MAE = \frac{1}{N} \sum_{i=1}^{N} |y_i - \hat{y}_i|, \qquad (10)$$

where N is the size of the test set;  $y_i$  and  $\hat{y}_i$ denote the ground truth and predicted score of the *i*-th data sample in the test set, respectively. MAE value ranges from 0 to positive infinity. It emphasizes the ranking order of the prediction instead of the absolute value. A lower MAE value indicates better performance.

• **Pearson Correlation** (PC) is defined as the covariance of the prediction and the ground 1105



Figure 12: More examples for property2img task.

truth divided by the product of their standard deviations. For two random variables x and y, Pearson Correlation is formally defined as

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115 1116

1117

1118 1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

$$PC = \frac{\mathbb{E}[(x - \mu_x)(y - \mu_y)]}{\sigma_x \sigma_y}, \quad (11)$$

In the regression task, suppose there are Ndata points in the test set,  $y_i$  is the ground truth of the *i*-th data sample,  $\hat{y}_i$  is the prediction for *i*-th data, Pearson Correlation becomes

$$PC = \frac{\sum_{i=1}^{N} \left( (y_i - \mu_y)(\widehat{y}_i - \mu_{\widehat{y}}) \right)}{\sigma_y \sigma_{\widehat{y}}}, \quad (12)$$

where  $\mu_y = \frac{1}{N} \sum_{j=1}^N y_j$  and  $\mu_{\widehat{y}} = \frac{1}{N} \sum_{j=1}^N \widehat{y}_j$  are mean of ground truth and prediction, respectively.  $\sigma_y = \sum_{i=1}^{N} (y_i - \frac{1}{N} \sum_{j=1}^{N} y_j)^2$  and  $\sigma_{\widehat{y}} = \sum_{i=1}^{N} (\widehat{y}_i - \frac{1}{N} \sum_{j=1}^{N} \widehat{y}_j)^2$  are the standard deviations of ground truth and prediction, respectively. The value ranges from -1 to 1. A higher Pearson correlation value indicates better performance.

> • Tanimoto similarity is to measure the similarity between two molecules. Tanimoto similarity is also known as the Jaccard coefficient, *i.e.*, the ratio of their intersection to their union over two chemical fingerprint vectors.

$$\sin(X,Y) = \frac{|\mathbf{b}_X \cap \mathbf{b}_Y|}{|\mathbf{b}_X \cup \mathbf{b}_Y|}, \qquad (13)$$

where  $\mathbf{b}_X$  is the binary fingerprint vector for 1130 the molecule X. Tanimoto distance between two molecules is defined as one minus Tani-1132 moto similarity.

Tanimoto-distance $(X,$	$Y) = 1 - \sin(X, Y),$
	(14)

Also, given a set of chemical compounds, 1135 we are typically interested in their diversity, 1136 which is defined based on Tanimoto distance. 1137 Specifically, diversity is defined as the aver-1138 age pairwise Tanimoto distance between the 1139 molecular fingerprints, 1140

diversity(
$$\mathcal{Z}$$
) =1 -  $\frac{1}{|\mathcal{Z}|(|\mathcal{Z}|-1)}$ .  

$$\sum_{X,Y\in\mathcal{Z},X\neq Y} sim(X,Y),$$
(15)

1141

where  $\mathcal{Z}$  is the set of generated molecules to 1142 evaluate. 1143

• Tanimoto@1 is to measures the proportion 1144 of generated molecules that exactly match the 1145 ground-truth molecule in terms of Tanimoto 1146 similarity. Specifically, it computes the ratio 1147 of generated molecules whose Tanimoto sim-1148 ilarity with the corresponding ground-truth 1149

1131

1134



Figure 13: More examples for img2img task.

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

 $\begin{aligned} \text{Tanimoto} @1 &= \frac{1}{N} \sum_{i=1}^{N} \\ \mathbb{I}\left[\text{Tanimoto}(f_i^{\text{gen}}, f_i^{\text{true}}) = 1\right], \end{aligned} \tag{16}$ 

molecule is equal to 1. It is defined as

where N is the number of generated molecules,  $f_i^{\text{gen}}$  and  $f_i^{\text{true}}$  are the Morgan fingerprints of the *i*-th generated and groundtruth molecules, respectively, and  $\mathbb{I}[\cdot]$  is the indicator function. The score ranges from 0 to 1, where a higher Tanimoto@1 indicates better exact matching performance between generated and reference molecules.

• **Increased LogP** is to evaluate molecular optimization performance. It measures the average increase in the LogP of molecules after optimization. For each molecule, the improvement is computed as the difference between the LogP value of the optimized molecule and that of the original molecule. The final score is the mean improvement across all molecule pairs. It is defined as

Increased 
$$\text{LogP} = \frac{1}{N} \sum_{i=1}^{N} (17)$$
  
 $\left(\text{LogP}(m_i^{\text{opt}}) - \text{LogP}(m_i^{\text{orig}})\right),$ 

where N is the number of molecule pairs,<br/> $m_i^{\text{orig}}$  and  $m_i^{\text{opt}}$  denote the *i*-th original and<br/>optimized molecules, respectively. A higher<br/>LogP Improvement value indicates a greater<br/>enhancement of the LogP property through<br/>the optimization process.1170<br/>1171

1176

1177

1178

1179

1180

• **Diversity** is a metric used to quantify the structural variety within a set of molecules. It is defined as the average pairwise Tanimoto distance between the Morgan fingerprints of the molecules:

Diversity 
$$= \frac{2}{N(N-1)} \sum_{i=1}^{N} \sum_{j=i+1}^{N} (18)$$
$$(1 - \operatorname{Tanimoto}(f_i, f_j)),$$

where N is the number of molecules in the 1182 set, and  $f_i$  and  $f_j$  are the Morgan fingerprints 1183 of the *i*-th and *j*-th molecules, respectively. 1184 The Tanimoto similarity measures the over-1185 lap between two binary fingerprints, and the 1186 distance is computed as 1 - Tanimoto. The 1187 diversity values range from 0 to 1, with higher 1188 values indicating greater chemical diversity. 1189

• Novelty evaluates the proportion of generated 1190 molecules that are not present in the training 1191



Figure 14: Training curve from start to the best check point, we apply GAN loss after training in E,G, and Z for a period of steps for stability. As shown in the validation loss curve, GAN loss is introduced at 45000 steps and cause the oscillation of validation loss and finally converge to stable result.

set. It reflects the ability of a generative model to produce novel chemical structures, rather than simply memorizing and replicating the training data. It is defined as

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1204

1205

1207

1208

1209

1210

1211

1212

Novelty = 
$$\frac{|\mathcal{G} \setminus \mathcal{T}|}{|\mathcal{G}|}$$
, (19)

where  $\mathcal{G}$  denotes the set of generated molecules, and  $\mathcal{T}$  denotes the set of molecules in the training set. The numerator counts the number of molecules in  $\mathcal{G}$  that are not in  $\mathcal{T}$ . The score ranges from 0 to 1, with higher values indicating greater novelty.

• Valid % is to evaluate the structural and syntactic validity of model outputs in instructionfollowing tasks involving molecule generation. It measures the proportion of outputs that are both (1) successfully parsed according to a predefined structured format (*i.e.*, instructionfollowing), and (2) contain syntactically valid SMILES strings, if any are present. It is defined as

Valid Rate 
$$=\frac{1}{N}\sum_{i=1}^{N}$$
 (20)

 $\mathbb{I}\left[\mathsf{structured}(o_i) \land \mathsf{valid}(o_i)\right],$ 

where N is the total number of model outputs,  $o_i$  is the *i*-th output, structured( $\cdot$ ) checks whether the output follows the expected structured format, and valid( $\cdot$ ) verifies the syntactic validity of any SMILES strings present in the output. The score ranges from 0 to 1, with a higher Valid% indicates better adherence to the required output format and chemical validity.

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1224

1225

1226

1227

1228

1230

Also, we conduct statistical testing to check if the improvement is statistically significant.

# **G** Molecular Properties

- **MW**: Molecular Weight (MW) is the sum of atomic masses of all atoms in a molecule (units: g/mol or Da). It influences physicochemical properties such as solubility, diffusion rate, and bioavailability. MW should be below 500 Da for optimal oral bioavailability.
- LogP: Octanol-water Partition Coefficient 1231 (LogP) assesses the solubility and synthetic 1232 accessibility of a chemical compound. The 1233 LogP score of a molecule ranges from  $-\infty$  to 1234  $+\infty$ . 1235
- TPSA: Topological Polar Surface Area 1236

1237(TPSA) quantifies the surface area contributed1238by polar atoms, typically oxygen and nitrogen,1239including their attached hydrogens. The the-1240oretical TPSA ranges from 0 to several hun-1241dreds or even thousands of Å<sup>2</sup> for highly polar1242or large biomolecules.

1243

1244

1245

1246

1247

1248

1249

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1280

1281

1282

1284

- HBD: Hydrogen Bond Donor (HBD) counts the number of polar functional groups (*e.g.*, -OH, -NH) in a molecule that can donate hydrogen atoms to form hydrogen bonds.
  - **HBA**: Hydrogen Bond Acceptor (HBA) counts the number of atoms (*e.g.*, O, N, S, F) in a molecule capable of accepting hydrogen bonds via lone electron pairs. Typical small-molecule drugs containing 2–10 HBA sites.
  - **RB**: Rotatable Bond (RB) counts the number of non-ring single bonds (*e.g.*, C-C, C-O, C-N) in a molecule that allow free rotation at room temperature. Optimal drug-like compounds typically contain ≤10 rotatable bonds (RB).
  - **QED**: Quantitative Estimate of Drug-likeness (QED) is an integrative score to evaluate compounds' favorability to become a drug. The QED value ranges from 0 to 1. A higher value is more desirable.

# H Implementation Details

For the first training stage, Mol-VQGAN is trained on 8 NVIDIA A800×80G GPUs for two epochs. The batch size is set to 16 and the base learning rate is set to 4.5e-06. We use Adam (Kingma and Ba, 2014) as optimizer and  $\mathcal{L}_{rec} + \lambda_1 \mathcal{L}_{perceptual}$ as monitor to validate the model and save the best checkpoint during 2-epoch training. For the second training stage, ChemMLLM is trained on 8 NVIDIA A800×80G GPUs for three epochs. The batch size is set to be 16, the base learning rate is set to be 2e-5 and z-loss weight is set to be 1e-5. We use AdamW (Loshchilov and Hutter, 2017) as optimizer and employ mixed-precision training, utilizing brain floating point 16 precision (bf16) (Kalamkar et al., 2019) for forward propagation and 32-bit floating point precision (fp32) for backward propagation to balance the training efficiency and stability. To handle distributed training, we apply PyTorch Fully Shared Data Parallel (FSDP) (Zhao et al., 2023) strategy. We train three variants of ChemMLLM for different tasks.

Mol-VQGAN Training For Mol-VQGAN train-1285 ing code, we use the official implementation code 1286 of original VQGAN (Esser et al., 2021). Specifi-1287 cally, we utilize the synthesized image datasets to 1288 let the original Chameleon VQGAN learn how to 1289 understand and generate molecule images. We sam-1290 ple 1,000,000 molecule images from PubChem syn-1291 thesized by RDKit (Landrum et al., 2006) and com-1292 bine them with all images synthesized in 5 down-1293 stream tasks as training dataset for Mol-VQGAN. 1294 All parameters of Encoder, Decoder and codebook 1295 of VQGAN are trained. We first train VQGAN 1296 on this 1 million-level molecule image dataset for 1297 two epochs and save the best check point according 1298 to  $\mathcal{L}_{rec} + \lambda_1 \mathcal{L}_{perceptual}$ . After the initial training, 1299 we find that it can not reconstruct image well on 1300 the dataset with less data size like img2caption 1301 dataset, so we do continuous training based on the 1302 best checkpoint using small-size datasets. Specifi-1303 cally, we continue training the 1-million best check 1304 point on img2caption images for five epochs and 1305 save the best checkpoint. Finally, we get the well-1306 trained Mol-VQGAN to tokenize molecule im-1307 ages. The original VQGAN will result in unclear 1308 and distorted images when encoding and decod-1309 ing molecule images. After training, Mol-VQGAN 1310 can encode and decode molecule image almost the 1311 same with original image. Several examples are 1312 shown in Figure 15. The validation curve from start 1313 to the best checkpoint is shown in Figure 14. 1314

**ChemMLLM Supervised FineTuning Training** For ChemMLLM training code, we utilize LuminamGPT framework (Liu et al., 2024). We train three variants ChemMLLM ChemMLLM-pro2img and ChemMLLM-img2img.

1315

1316

1317

1318

1319

1320

1321

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

For ChemMLLM, we train it on img2caption, img2property, img2smiles and property2img Apart from the four datasets, we datasets. also train it on SMILES molecule image generation (smiles2img) and text based drug design dataset (caption2img), which is the inverse task of img2smiles and img2caption. smiles2img requires model to receive a SMILES and output the corresponding image and caption2img ask model to output molecule image according to molecule caption. We train ChemMLLM on the six datasets for three epochs; For property2img task, we train and evaluate ChemMLLM-pro2img. We find that directly training the model on raw data can not achieve good performance, so we do data augmentation by rotating images by  $90^{\circ}$ ,  $180^{\circ}$ , and  $270^{\circ}$  so as to generate augmented data that is four times the size of the



Figure 15: Examples of origin VQGAN and Mol-VQGAN. The origin images is shown in the first row. As shown in the second row, the origin VQGAN can not encode and decode molecule image clearly and accurately, resulting in distorted atoms and bonds which are hard to distinguish. As shown in the third row, our well-trained Mol-VQGAN can encode and decode molecule images with clear atoms and bonds, which is almost the same as the origin molecule images.

original dataset. We train ChemMLLM-pro2img on the augmented dataset for two epochs; For img2img task, we train and evaluate ChemMLLMimg2img. we train it on img2img dataset for two epochs. The training hyper-parameters primarily follow Lumina-mGPT (Liu et al., 2024).

The details of data information for training different ChemMLLM variants are shown in Table 16. The details of training settings for training different ChemMLLM variants are shown in Table 17.

Model	task	# Training	# Test
	img2caption	69,799	3,300
ChemMLLM	img2property	95,000	5,000
	img2smiles	95,000	5,000
	property2img		5,000
	smiles2img	95,000	5,000
	caption2img	72,143	3301
ChemMLLM-pro2img	property2img	380,000	20,000
ChemMLLM-img2img	img2img	157,673	17,520

Table 16: Detailed dataset information for training different ChemMLLM variants.

## I Ablation Study

We conduct ablation study on property2img task with a focus on evaluating the impact of Mol-VQGAN training and data augmentation. The ef-

Settings	ChemMLLM	ChemMLLM-pro2img	ChemMLLM-img2img	
z-loss weight	1e-05	1e-05	1e-05	
warmup epochs	0.01	0.01	0.01	
learning rate	2e-05	2e-05	2e-05	
weight decay	0.1	0.1	0.1	
drop rate	0.05	0.05	0.05	
total bacth size	$16 \times 8 \times 1$	$16 \times 8 \times 1$	$8 \times 4 \times 1$	
GPUs for training	8×A800 (80G)	8× A800 (80G)	4× A800 (80G)	
GPUs hours(h)	65	30.6	25.4	

Table 17: Detailed training settings for training different ChemMLLM variants.

fectiveness of the generated molecular images is assessed through the Pearson correlation coefficients between the predicted and ground truth values of several molecular properties, including MW, LogP, TPSA, Hbd, Hba, Rb, and QED. we use a subset of 200 samples of he property2img task and do 5-shot evaluation.

1352

1353

1354

1355

1356

1357

1359

1360

1362

1363

1364

1365

1367

The result is shown in Table 18. When both fine-tuned VQGAN and data augmentation are employed, the model achieves the highest correlation scores across all evaluated properties. Notably, the correlations for MW (0.71), TPSA (0.71), Hba (0.66), and Rb (0.62) indicate that the generated images capture molecular structure and features that align well with the original textual descriptions. This demonstrates the efficacy of our approach in enhancing the semantic fidelity and chemical rele-

1348

1349

1350

data augmentation	VQGAN	MW Pearson $(\uparrow)$	$LogP \ Pearson \ (\uparrow)$	TPSA Pearson $(\uparrow)$	Hbd Pearson $(\uparrow)$	Hba Pearson $(\uparrow)$	Rb Pearson $(\uparrow)$	QED Pearson $(\uparrow)$
✓	~	0.71	0.42	0.71	0.45	0.66	0.62	0.34
х	$\checkmark$	0.46	0.04	0.40	0.08	0.53	0.17	0.26
×	×	0.55	0.06	-0.05	0.06	-0.06	0.31	0.35

Table 18: Ablation study.

vance of the generated visual representations. Re-1368 moving data augmentation while retaining the fine-1369 tuned VQGAN leads to a significant drop in per-1370 formance, especially for properties such as LogP 1371 (reduced to 0.04) and TPSA (0.40), highlighting 1372 the importance of data augmentation in improv-1373 ing the model's generalization and robustness dur-1374 ing training. The performance further deteriorates 1375 when both fine-tuning and data augmentation are removed. In this setting, the model yields the low-1377 est correlations, with some properties (e.g., TPSA 1378 at -0.05 and Hba at -0.06) exhibiting negative corre-1379 lation, suggesting that the general VQGAN trained 1380 on natural images fails to preserve critical molecu-1381 lar features necessary for reliable property predic-1382 tion. 1383

In summary, these results clearly demonstrate that both fine-tuning Mol-VQGAN and applying data augmentation play complementary and crucial roles in enhancing the quality and accuracy of chemical multimodal tasks.

1384

1385

1386