

引用格式: 王纪峰, 汪莹. 生成式深度学习在目标导向分子设计中的应用进展[J]. 中国材料进展, 2025, 44(5): 424-435.
WANG J F, WANG Y. Application of Generative Deep Learning in Object-Oriented Molecular Design[J]. Materials China, 2025, 44(5): 424-435.

特约专栏

生成式深度学习在目标导向分子设计中的应用进展

王纪峰, 汪莹

(复旦大学高分子科学系 聚合物分子工程国家重点实验室, 上海 200438)

摘要: 分子设计作为化学与材料科学中的一项核心任务, 面临着在庞大的化学空间中高效筛选并开发具备特定功能的分子的问题, 传统方法在效率和探索性方面存在明显局限。近年来, 生成式深度学习的兴起为分子设计提供了自动化与智能化的新路径。综述了生成式深度学习在分子设计中的应用进展, 首先对不同分子表示方法(如 SMILES、分子图和三维结构表示)进行比较, 分析了各自的优缺点。随后, 综合评估了3种主流生成式模型: 生成对抗网络(GAN)、变分自动编码器(VAE)和去噪扩散概率模型(DDPM), 并探讨了生成式模型在目标导向分子设计中的应用, 重点分析不同模型在分子生成质量与性质优化方面的差异。最后, 基于现有技术的研究进展, 提出了未来生成式模型在分子设计领域的研究方向。

关键词: 分子生成; 生成式深度学习; 生成对抗网络; 变分自动编码器; 去噪扩散概率模型; 模型性能评估框架; 分子表示
中图分类号: TQ317; TP18 **文献标识码:** A **文章编号:** 1674-3962(2025)05-0424-12

Application of Generative Deep Learning in Object-Oriented Molecular Design

WANG Jifeng, WANG Ying

(State Key Laboratory of Molecular Engineering of Polymers, Department of Macromolecular Science, Fudan University, Shanghai 200438, China)

Abstract: Designing molecules with specific functions within an immense chemical space is a fundamental challenge in chemistry and materials science, as traditional methods often lack efficiency and exploratory capacity. The advent of generative deep learning has introduced automated and intelligent approaches that promise to transform molecular design. In this review, we summarize the advancements in applying generative deep learning to molecular design. We first compare molecular representation methods, including SMILES notation, molecular graphs, and three-dimensional structural representations, highlight their respective advantages and limitations. We then critically evaluate three leading generative models: generative adversarial network (GAN), variational autoencoder (VAE), and denoising diffusion probabilistic model (DDPM), and discuss applications of generative models in object-oriented molecular design, with a particular focus on the differences among various models in terms of molecular generation quality and property optimization. Finally, we propose future research directions for leveraging generative models in molecular design, aiming to inspire further advancements in this rapidly evolving field.

Key words: molecule generation; generative deep learning; generative adversarial network (GAN); variational autoencoder (VAE); denoising diffusion probabilistic model (DDPM); model performance evaluation framework; molecular representation

收稿日期: 2024-11-14 修回日期: 2025-05-01
基金项目: 国家自然科学基金重大研究计划培育项目(92372126)
第一作者: 王纪峰, 男, 2000年生, 博士研究生
通讯作者: 汪莹, 女, 1989年生, 青年研究员, 博士生导师,
Email: wying@fudan.edu.cn
DOI: 10.7502/j.issn.1674-3962.202411011

1 前言

在现代化学与材料科学领域, 如何设计具有特定功能的分子是一个核心问题。传统的分子设计大多依赖于科学家的经验积累和直觉判断, 但这种方式不仅耗时耗力, 而且面对巨大的化学空间, 难以系统性地探索潜在的分子结构^[1]。这使得通过试错法找到具备特定化学或

生物学性质的分子成为一项极具挑战的任务。因此迫切需要一种自动化、智能化的分子设计方法，以提升新分子发现的效率^[2, 3]。

生成式深度学习作为近年来兴起的人工智能方法，在分子设计中的应用逐步显现出独特的优势。生成模型可以在无监督或弱监督的条件下，从已有的分子数据中学习复杂的分布规律，从而生成具有类似特性的全新分子^[4]。与传统的枚举法或基于反应的分子合成方法相比，生成式模型不仅能够高效地扩展分子空间，还可以在不依赖特定反应路径的情况下生成具有较高合成可行性的分子^[5]。此外，深度学习模型的可拓展性和灵活性也使得它们能够融合化学、材料和生物学领域的多种特征，为不同应用场景下的目标导向分子的生成提供了可能^[6, 7]。

目前，基于深度学习的生成模型在分子设计中的应用主要包括功能材料分子设计和药物先导化合物研发等方面。具体而言，生成对抗网络 (generative adversarial network, GAN)^[8]、变分自动编码器 (variational autoencoder, VAE)^[9]、去噪扩散概率模型 (denoising diffusion probabilistic model, DDPM)^[10] 等生成式模型在分子生成中取得了显著进展。例如，GAN 通过对抗训练生成结构合理的分子，而 VAE 可以构建具有连续特征的潜在空间以便于分子优化，DDPM 则通过逐步学习分子分布中的噪声消除过程，生成符合分子物理和化学特性的结构。这些方法的融合进一步推动了分子设计从经验驱动的模式向数据驱动的模式转变，使得新分子的发现更加系统化、自动化。

此外，随着生成式模型的发展，特定目标导向的分子生成已成为研究的热点。通过加入强化学习和条件生成等优化方法，生成式模型能够生成具有特定分子活性或化学性质的分子^[11]。这种特征控制能力让生成模型在功能材料优化等实际应用中展现出巨大潜力。然而，在生成式深度学习用于分子生成的研究领域中，仍存在诸多亟待解决的核心问题，主要集中在以下几个方面。

(1) 生成分子的合成可行性：生成式模型虽然能够产生大量结构多样的分子，但这些分子在实验中往往难以合成，限制了其实际应用价值^[12, 13]。

(2) 目标特性控制的准确性和有效性：目前的生成式模型在控制生成分子的特性方面仍存在较大不确定性，难以精确调整分子以满足特定要求^[14]。

(3) 化学空间的高效探索：化学空间的庞大和复杂性使得生成式模型难以全面、高效地覆盖其中的潜在分子，可能导致重要的新颖分子被遗漏^[15]。

(4) 模型泛化能力的提升：生成式模型通常依赖于特定的数据集进行训练，导致模型对新化合物或不同数据分布的泛化能力不足。

近年来，已有大量综述聚焦于生成式模型在药物分子、材料分子等不同领域的应用进展，涵盖了分子生成策略、性能评估及相关算法进展。然而，目前尚缺乏针对“分子处理流程—分子表示—数据库—深度学习生成模型”这一完整链路的系统性分析，尤其缺少面向初学者的、可迁移于不同细分领域的通用框架总结。本文系统梳理了生成式深度学习模型生成分子的全流程，从分子表示方法、主流分子数据库，到基于深度学习的分子生成技术，并围绕目标导向的分子设计的评估体系和最新进展进行了结构化归纳。本文不仅剖析了不同生成式模型的优势与不足，还力求为初步涉足该领域的研究者提供一条清晰、易操作的入门路径。

2 分子处理方法

2.1 分子表示

在深度学习的分子生成任务中，分子的表示方式至关重要，不同的表示方式会直接影响模型的输入处理方式、生成效果及生成分子的特性。以下介绍了几种在深度学习中常用的分子表示方法，如图 1 所示，并评估了它们的优势和缺点。

SMILES 表示 (simplified molecular input line entry system)^[16]：SMILES 是一种将分子结构转化为一维字符串的表示方法，记录了原子及其连接方式。SMILES 表示形式紧凑、便于存储，并且广泛应用于化学式的表达。此外，SMILES 字符串可以通过独热编码 (one-hot encode) 处理，将字符转化为模型可处理的稀疏向量形式，使其适用于多种深度学习算法^[17]。然而，SMILES 表示也存在一定的冗余性，同一个分子可以对应多个不同的 SMILES 表示形式；且在生成过程中，模型可能生成不符合化学规则的无效 SMILES 字符串，导致生成分子的有效性较低^[18]。

分子图 (molecular graph)^[19]：在分子图表示中，原子被视为节点，化学键作为边，可分别构建节点矩阵和邻接矩阵。分子图能够完整表达分子的结构信息，广泛应用于生成式深度学习模型中。然而，由于图的稀疏性和不规则性，图神经网络的计算复杂度较高，且对大型分子处理的效率较低^[20]。

三维分子结构 (3D structure)^[21]：通过将分子的空间坐标与连接关系相结合，表示原子在三维空间中的排布。在笛卡尔坐标系中为连接关系中的每个原子分配三维空间坐标 (x, y, z)，从而直接描述分子的几何构型。除此之外，可基于内部坐标，例如键长、键角和二面角，描述每个原子相对于其邻近原子的空间位置，而无需使用全局坐标系。这种描述方式有效捕捉了分子的三维结构特征，使其能够更准确地表征分子间的空间关系和远程

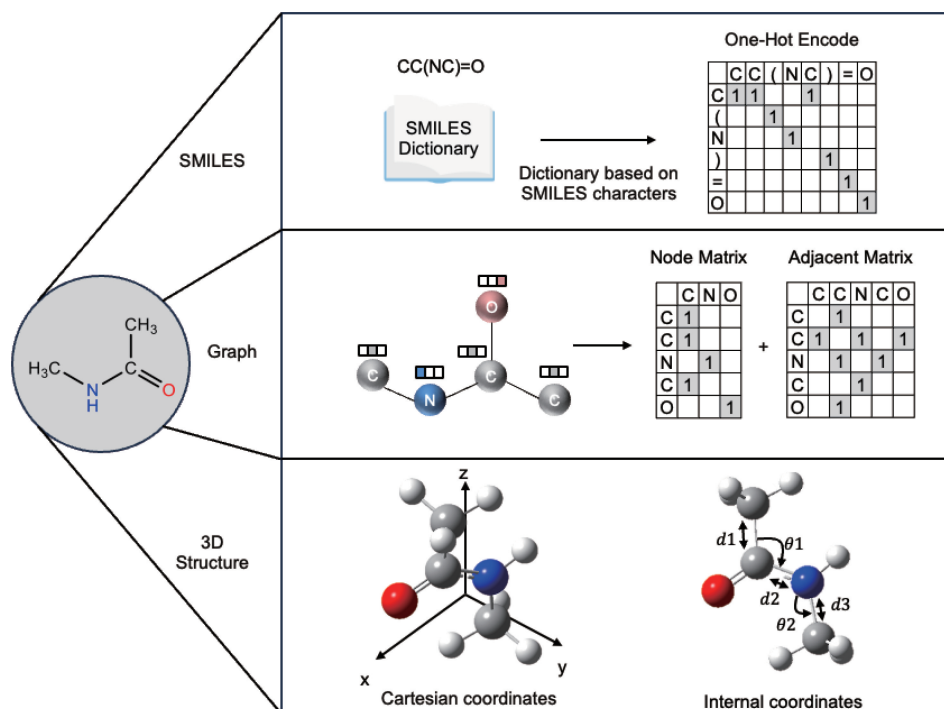


图 1 以 N-methylacetamide 为例, 分别展示了基于 SMILES、分子图和三维分子结构的分子表示: 在 SMILES 分子表示中, 首先需要根据 SMILES 中所有存在的不同种类的字符设定 SMILES 字典, 然后以当前分子的 SMILES 作为独热编码的横轴、SMILES 字典作为纵轴构建矩阵; 在分子图表示中, 仅考虑分子中的重元素, 然后分别构建分子的节点矩阵和邻接矩阵; 在三维分子结构表示中, 基于三维坐标(笛卡尔坐标系)表示节点或者将三维信息作为节点属性(内坐标表示)

Fig. 1 Illustration of molecular representations for N-methylacetamide, including SMILES, molecular graph, and 3D structure: for the SMILES representation, a SMILES dictionary is constructed based on all unique characters present, with a one-hot encoding matrix where the SMILES string is set as the horizontal axis and the dictionary as the vertical axis, filling in 1 at matching positions and 0 elsewhere; in the molecular graph representation, only heavy atoms are considered, constructing both a node matrix and an adjacency matrix; in the 3D structure representation, spatial information is based on 3D coordinates in a Cartesian coordinate system or incorporated either as node attributes (internal coordinates)

相互作用^[22]。缺点为三维分子结构的计算复杂度较高, 可能增加模型的计算负担。其次, 由于分子具有多种构象, 在生成任务中采用三维描述需要考虑构象的多样性和动态性, 这给分子生成的稳定性和效率带来一定挑战^[23]。为了解决该问题, Zheng 等^[24]开发了分布式图变换生成器(distributional graphormer), 旨在预测分子系统的平衡分布, 可提供比传统方法快几个数量级的分子构象生成与采样。

需要注意, 分子的表示方式选择需要根据任务要求和模型特点进行权衡。SMILES 使用便捷, 但分子生成有效性较低; 分子图适合捕捉分子结构关系, 但计算成本较高并且表示过程较为复杂; 三维分子结构描述能够准确捕捉分子的空间结构和相互作用特性, 但其数据量大、计算成本较高^[25, 26]。因此, 选择合适的分子表示方式将有助于改善生成式模型的表现, 从而在分子设计中完成更为高效和精准的生成任务^[27]。

2.2 分子数据库

在训练用于分子生成的基于深度学习的生成式模型

时, 选择合适的分子数据库至关重要。表 1 给出了从小分子到复杂分子的常用分子数据库。ChEMBL^[28] 数据库不仅包含大量结构数据, 还提供了丰富的元信息, 如靶点类别、测定类型及实验条件等。该数据库特别适用于基于性质约束的分子生成任务。ZINC^[29] 数据库以其丰富的化学结构和详细的可购性标签广泛应用于分子生成任务, 常用作模型语法预训练语料库。PubChem^[30] 数据库数据来源广泛, 涵盖天然产物、材料化学等多个领域, 适合用于跨领域迁移学习和多任务预训练。然而, 由于其数据异质性强, 需在模型训练前对数据进行标准化处理和去偏采样, 以降低数据分布漂移带来的影响。GDB-17^[31] 数据库收录了大量仅含结构信息的分子, 尽管缺乏实验性质数据, 但可为生成式模型的预训练提供丰富而广泛的分子构象基础, 有助于提升模型的结构感知能力与泛化性能。该数据库常与小规模有标注数据联合, 用于低标签场景下的对比学习或自监督对齐。QM9^[32] 数据库则聚焦于分子的量子化学性质, 是发展物性引导分子生成和机器学习-量子化学混合模型的重要基准。QM9 也

广泛用于评估模型在电子结构性质预测任务中的泛化与外推能力。

值得注意的是，数据库的规模与多样性是影响生成式分子设计模型泛化能力的关键因素。大规模数据库为模型提供了丰富的训练样本，能够有效覆盖更广阔的化学空间，从而提升模型在未见分子结构上的预测和生成能力^[33]。同时，数据库的多样性，即所包含分子的结构类型、化学属性及应用领域的广泛性，直接决定了模型对不同化学环境的适应能力。多样化的数据能够促使模型学习到更具代表性的结构特征和性质关联，减少模型对特定类别分子的过拟合风险。相反，规模有限或类型单一的数据库容易导致模型泛化能力不足，难以在实际应用中生成具有创新性或满足特定功能需求的新分子。因此，构建大规模且多样化的分子数据库，不仅是推动分子生成模型性能提升的基础，也是实现分子设计创新与应用落地的前提条件^[34]。

表1 训练分子生成式模型常用的数据库

Table 1 Commonly used databases for training molecular generation models

Database	Data number	Description
ChEMBL ^[28]	2M	Diverse bioactive compounds, especially drug-like small molecules.
ZINC ^[29]	750M	Commercially available drug-like molecules.
PubChem ^[30]	111M	Millions of small molecules with chemical and biological annotations.
GDB-17 ^[31]	116B	Structure-only molecules without experimental data.
QM9 ^[32]	134K	Small molecules (≤ 9 heavy atoms) with quantum chemical annotations.

2.3 分子处理工具

分子数据的高效处理需要专门的工具将分子数据转换为适合计算机分析的格式和结构。这些工具不仅能完成分子文件格式转换、生成分子图与三维构象等核心功能，还可在自动化处理和数据可视化方面为研究者提供支持，大幅减少了手动操作的时间与工作量。

RDKit^[35]是一款功能强大的开源化学信息学工具包，能够对分子数据进行读取、表示、修改、输出和可视化处理。该工具不仅支持多种分子文件格式，还可以进行分子属性分数计算、分子片段化、手性中心识别以及高效的子结构搜索等操作，因此在化学结构处理领域具有广泛应用。

Open Babel^[36]则是一款支持超过100种分子文件格式的转换工具，可以实现不同文件格式间的双向转换。它还提供了分子指纹计算和分子相似性比较等功能。

DeepChem^[37]是一个面向化学和生物学研究的深度学

习框架，集成了多个适用于分子和药物发现的深度学习模型。作为专为化学研究设计的深度学习库，DeepChem包含多种预训练模型和数据处理管道，为分子特征提取、活性预测、毒性预测等任务提供了丰富的算法支持。

3 生成式深度学习模型

GAN(图2a)是一种通过对抗训练生成数据的模型，由生成器和判别器两个神经网络组成^[8]。生成器负责生成假分子样本，而判别器则用于区分生成样本与真实样本。通过生成器和判别器的对抗博弈，生成器逐渐学习生成更具真实性的分子结构。GAN在分子设计中的应用主要聚焦于生成同时满足结构合理性与化学特征约束的分子结构。其优势在于其生成结果的多样性和丰富性，尤其在分子结构生成和目标性质优化任务中具有显著效果。然而，GAN模型在训练过程中存在较为突出的模式崩溃(mode collapse)问题，即生成器倾向于输出有限的几种分子结构，导致生成样本的多样性降低。这种现象的产生主要源于对抗训练过程中生成器和判别器博弈不平衡，生成器捕捉到某些能够“骗过”判别器的模式后，可能会忽略其他潜在结构。为缓解模式崩溃，目前提出了多种改进方案。例如，Wasserstein GAN(WGAN)^[38]通过引入Earth-Mover距离，提高了训练过程的稳定性；Least Squares GAN(LSGAN)^[39]通过最小均方损失缓解了梯度消失问题。这些改进损失函数在实际分子生成任务中有助于提升模型的收敛速度和生成多样性。在训练策略上，GAN通常采用交替更新生成器和判别器参数的方式，常见做法为每训练一次生成器就训练多次判别器(如1:5)，以保持判别器相对优势。需要注意的是，GAN模型对超参数较为敏感，建议通过多次实验对比确定最优配置。

VAE(图2b)是一种基于概率建模的生成模型，通过神经网络实现对复杂分布的建模^[9]。VAE主要由编码器和解码器两部分构成，编码器将输入分子映射为潜在空间中的概率分布(通常为多维高斯分布)，随后从该分布中采样潜在向量，经解码器重构原始分子结构。VAE的潜在空间具有良好的连续性，这意味着对潜在空间的线性操作可以平滑地对应到分子的生成与优化，因而在分子设计和性质优化等任务中表现出独特优势。为了保证潜在空间分布的结构合理性，VAE引入了潜在空间正则化机制。常用的方法是通过KL散度(Kullback-Leibler divergence)对每个样本的近似后验分布加以约束，使其尽可能接近预设的先验分布(通常为标准正态分布)，从而提升模型的生成能力和潜在空间的可解释性。KL散度项有效防止模型陷入异常分布，促使生成分子的多样性和

连续性。损失函数方面, VAE 采用重构损失与正则化损失的加权和。重构损失通常采用均方误差或交叉熵, 用于衡量重构分子与原始输入之间的相似度; 正则化损失则是上述 KL 散度, 衡量编码器输出分布与先验分布(如标准高斯分布)之间的差异。两者的加权系数需要结合实际任务进行调整, 以平衡生成样本的多样性和准确性。值得注意的是, VAE 潜在空间维度过小可能导致模型表达能力不足, 过大则易过拟合, KL 散度权重的合理调整对于平衡生成分子的多样性与保真度尤为重要。

DDPM(图 2c)是一种基于逐步去噪过程的生成式模型^[10], 通过逐步引入噪声将分子结构扰动到简单的高斯分布, 随后通过反向去噪逐步恢复出符合原始分布的分子样本。尽管 DDPM 的生成效率较低, 但在分子生成中表现出极高的生成质量, 尤其适用于捕捉复杂的分子结构特征。其优势在于生成过程的稳定性和生成样本的高质量; 此外, 通过调整去噪过程中的参数, 可以精确控制生成分子的特性。在实际应用中, DDPM 通常采用固定数量的去噪步数(例如 1000 步), 每一步都对分子结构进行微小扰动。步数的选择会直接影响模型生成效率和样本质量: 步数越多, 反向生成过程越平滑, 生成样本的质量越高, 但生成速度也会降低。损失函数方面, DDPM 一般采用均方误差作为优化目标。具体而言, 模型在每个扩散步骤中尝试预测所添加的噪声, 并通过最小化预测噪声与真实噪声之间的均方差来更新网络参数。训练过程中, 对噪声的加入步数进行均匀采样, 确保模型充分学习各个去噪阶段的特征。此外, 还可根据具体任务加入数据增强等手段。超参数的设置对模型性能影响较大, 尤其是去噪步数(如 $T=1000$)以及噪声调度方式(例如线性调度与余弦调度)。合理设置去噪步数以及噪声调度方式有助于平衡生成样本的多样性与生成效率。

4 性能评估

4.1 模型性能评估框架

近年来, 多个标准化框架和基准数据集相继被提出, 以便在不同生成式模型之间进行性能比较。Brown 等^[40]提出了一个用于从头分子设计的基准框架 GuacaMol。该框架包含两类主要基准: 分布学习和目标导向生成。分布学习基准通过有效性、独特性和新颖性来评估模型生成新分子的能力, 而目标导向生成则着重于评估模型在优化特定分子特征或性质上的表现。此外, GuacaMol 框架收录了标准化数据集及基线模型的表现, 有助于客观比较新模型的性能。

Polykovskiy 等^[41]提出了另一套基准测试框架 MOSES, 主要用于评估生成式模型的分布学习能力。该框架通过有效性、独特性、新颖性、内部多样性以及通过一

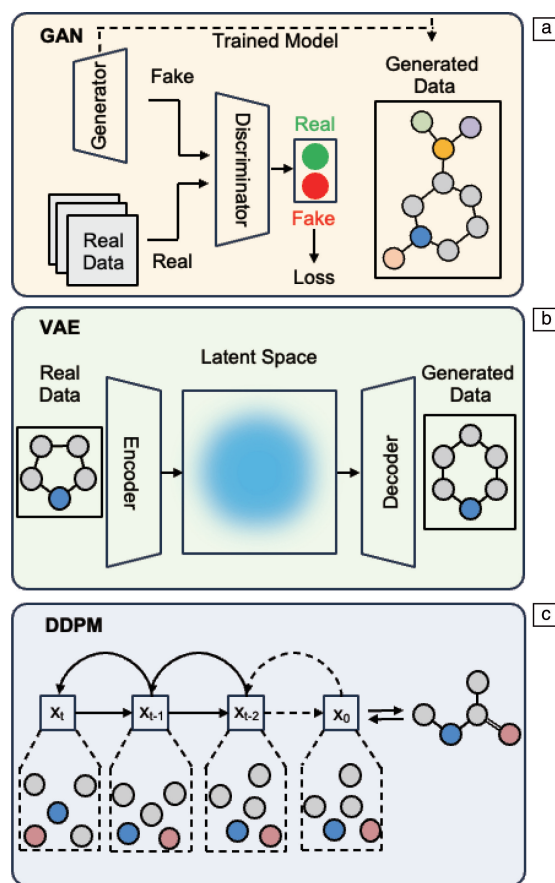


图 2 3 种生成式深度学习模型: (a) GAN: 由生成器和判别器组成, 分子生成时使用训练后的生成器; (b) VAE: 通过将分子编码到潜在空间中, 分子生成时在潜在空间中进行采样, 并通过解码器得到分子表示; (c) DDPM: 由正向过程(弯曲箭头)和反向过程(直箭头)构成, 通过反向过程逐步从标准正态噪声 x_1 中去除噪声, 生成分子

Fig. 2 Three generative deep learning models: (a) GAN model: it consists of a generator and a discriminator, with the trained generator used for molecular generation; (b) VAE model: it encodes molecules into a latent space, where molecular generation occurs by sampling in the latent space, followed by decoding to obtain molecular representations; (c) DDPM model: it involves a forward (curved arrows) and reverse (straight arrows) process, where molecular generation is achieved in the reverse process by gradually removing noise from standard normal noise x_1 to construct the molecule

系列分子质量和结构过滤器的比例来衡量模型生成新分子的能力。MOSES 还提供了一组指标, 以评估模型对训练数据集中分子特征的学习程度。这些指标包括 Fréchet 化学距离(FCD)、不同物理化学性质分布差异的距离度量、最近邻 Tanimoto 相似度、Bemis-Murcko 支架以及 BRICS 片段的余弦相似度。前两个指标用于评估更抽象的化学和生物相似性^[42], 后两个指标则有助于在子结构水平上进行分子比较。此外, 该框架还包含各种基线模型及标准化数据集, 以支持性能测试。

4.2 合成可及性评估

在分子生成领域，合成可及性评估是一个关键且具有挑战性的问题。由于生成式模型提出的分子通常未明确考虑其合成可行性，许多候选分子可能缺乏已知的合成路线。基于启发式的评分方法，SAscore (Synthetic Accessibility score)^[43] 和 SCScore (Synthetic Complexity score)^[44] 是两种广泛使用的评估指标。

SAscore 通过分析分子的结构复杂度和子结构在已知合成分子中的出现频率，对分子的合成难度进行评分。评分越高，表示分子越容易合成。它综合考虑了分子中的罕见结构和整体复杂性，为研究人员提供了快速评估工具。SCScore 利用机器学习模型，根据分子的结构特征预测其合成复杂度。该评分反映了从简单前体化合物合成目标分子的可能性，评分越低，表示合成所需的步骤可能越少，合成路径越简单。这些方法可作为评估分子合成可行性的初步手段，有助于在早期阶段排除难以合成的分子，提高材料研发与设计的效率。

5 目标导向分子生成

自动生成新分子的方法通常针对特定的性质和特征，例如溶解度或生物活性等。因此，能高效地生成既具有创新性又满足性能要求的分子的生成方法备受研究者的关注。表2总结了目标导向分子生成的先进方法，下文将对此详细描述。

5.1 GAN 模型

GAN 通过生成器与判别器之间的对抗性训练，能够捕捉复杂的分子结构和属性分布，从而生成具有化学意

表2 分子生成模型及每种模型所采用的分子表示与模型架构

Table 2 Molecular generation models and the molecular representations and model architectures employed by each model

Model name	Molecular representation	Model architecture
MolGAN ^[45]	Graph	GAN
LatentGAN ^[46]	SMILES	GAN
MatGAN ^[47]	Graph	GAN
Mol-CycleGAN ^[48]	Graph	GAN
Gene-GAN ^[49]	SMILES	GAN
DeepICL ^[50]	3D Structure	VAE
SDVAE ^[51]	SMILES	VAE
MGCVAE ^[52]	Graph	VAE
POLYGON ^[53]	SMILES	VAE
Drug-CVAE ^[54]	SMILES	VAE
MDM ^[55]	3D Structure	DDPM
EDM ^[56]	3D Structure	DDPM
GCDM ^[57]	3D Structure	DDPM
HGLDM ^[58]	Graph	DDPM
GEOLDM ^[59]	3D Structure	DDPM

义的新分子。其优势在于生成器可以根据判别器的反馈不断优化生成策略，使得生成的分子更有效性和新颖性，同时能够实现一定的分子性质优化^[60]。然而，GAN 在分子生成中的缺点也较为明显。由于对抗训练过程中的不稳定性，GAN 生成的分子易出现模式崩溃现象，即生成的分子多样性可能不足^[61]。为缓解这一问题，用图3所示方法在生成-判别对之外串联了基于图卷积网络 (graph

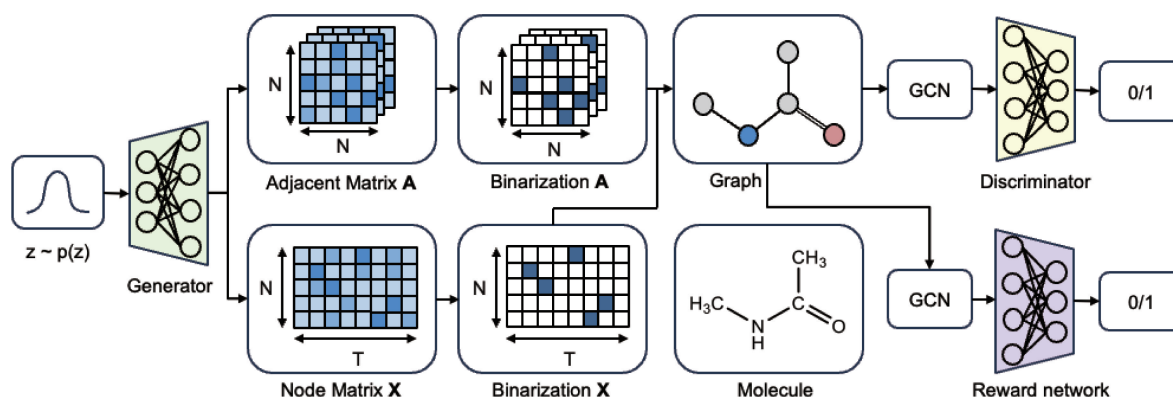


图3 采用分子图编码的条件 GAN 模型训练过程示意图：在该条件 GAN 中，生成器被训练以生成逼真的分子，而判别器则旨在辨别给定的分子是真实的还是生成的；随后，分别训练两个图卷积神经网络 (GCN)，用于评估生成分子的真实性以及其分子性质是否符合预期要求；经过充分训练后，生成器能够产生高度逼真且性质符合目标标准的分子

Fig. 3 Schematic representation of the training process of a conditional GAN model using molecular graph encoding: in this conditional GAN, the generator is trained to produce realistic synthetic data, while the discriminator aims to determine whether the given data is real or synthetic; subsequently, two graph convolutional networks (GCNs) are trained separately to evaluate the authenticity of the generated molecules and whether their molecular properties meet the desired criteria; after sufficient training, the generator is capable of producing highly realistic molecules with properties that align with the target specifications

convolutional network, GCN) 性质预测器, 并将其预测得分作为奖励信号嵌入强化学习回路。这种方式不仅抑制了模式崩溃, 还能引导生成过程朝着满足目标性质区间的分子空间收敛, 从而同时获得结构多样性与预期的目标性能。

de Cao 等^[45]提出的 MolGAN 模型专注于小分子的图结构生成, 成功避开了传统方法中的图匹配和节点顺序问题。MolGAN 利用 GAN 生成小分子的图结构, 并结合强化学习优化目标性质。其生成的化合物有效率接近 100%, 在确保分子有效性的同时, 提供了合理的多样性, 为小分子生成提供了高效的解决方案。在分子表示方面, Prykhodko 等^[46]提出了基于 SMILES 字符串潜在向量的 LatentGAN 模型, 融合了自编码器和 GAN 的优势。该模型在生成药物分子和靶向分子方面表现出色, 生成分子的新颖性超过 95%, 且其“药物相似性”得分与训练集相当。与传统的循环神经网络方法相比, LatentGAN 在化学空间覆盖上具有显著优势, 并在多个靶点上验证了其卓越的生成能力。针对无机材料的化学组成生成, Dan 等^[47]提出了 MatGAN 模型, 采用分子图表示方法生成无机材料的化学组成。MatGAN 通过在 ICSD 等数据库上训练, 生成的假设无机材料中有 84.5% 满足电中性和平衡电负性要求, 新颖性达 92.53%。此外, MatGAN 在采样效率方面比穷举法提高了 77 倍, 展示了其在化学设计空间采样

中的独特优势。在分子优化领域, Maziarka 等^[48]提出了 Mol-CycleGAN 模型, 基于 CycleGAN 架构, 使用分子图表示实现分子性质的优化。该模型通过学习特定分子性质的转换规则, 在优化物理化学性质(如 logP 值)方面表现出色, 同时保持了分子结构的相似性, 生成分子的有效性高达 100%。此外, Méndez-Lucio 等^[49]开发了一种基于基因表达谱的条件 GAN 模型, 利用 SMILES 表示生成能够诱导特定基因表达特征的分子。该模型通过条件 GAN 生成与特定靶点相关的潜在活性分子, 生成的分子在生物相似性和合成可行性方面表现优异, 与已知活性分子的相似性得分达到 0.64。该模型无需先验的活性分子数据, 即可生成满足特定生物活性的分子。

5.2 VAE 模型

VAE 的优势在于其通过潜在空间的连续分布对分子结构进行编码, 允许在该潜在空间中插值生成新分子, 从而能够控制生成分子的特定理化性质。这种方法可以帮助发现与已知分子结构相似但具有独特属性的新分子, 支持在广阔的化学空间内进行探索^[62]。然而, VAE 也存在一些局限, 例如生成分子时容易出现无效结构, 且在生成特定属性分子方面的精确控制能力仍有待提升^[63]。此外, 由于其生成结果依赖于数据分布, 模型在小样本数据集或新领域应用中的表现可能受到限制, 图 4 展示了条件生成 VAE 模型的基本架构。

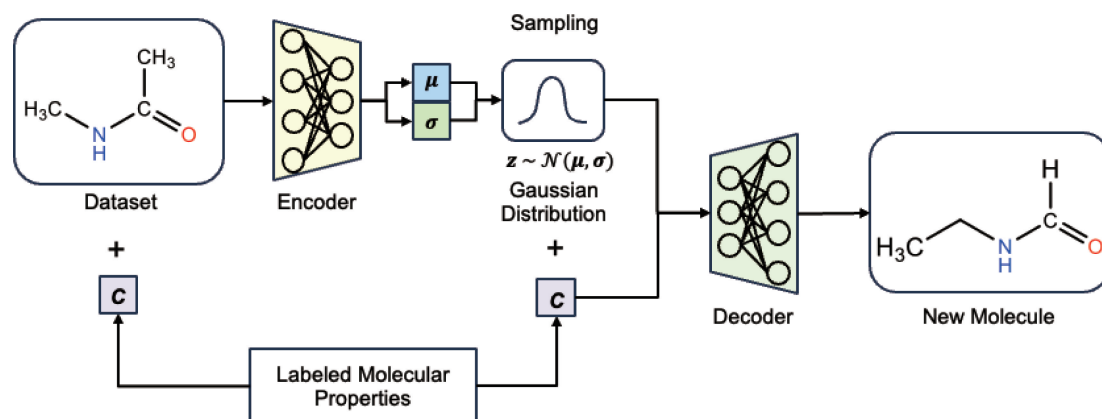


图 4 条件变分自编码器 (VAE) 模型的示意图: 将训练集中的分子通过编码器映射到潜在空间, 将条件信息同时嵌入到原始数据和潜在表示中; 在生成阶段, 从潜在空间中随机采样, 然后由解码器解码, 生成满足特定条件的分子

Fig. 4 Schematic illustration of a conditional variational autoencoder (VAE) model; molecules from the training set are mapped into the latent space by the encoder, with conditional information embedded in both the original data and the latent representation; during the generation stage, random samples are drawn from the latent space and decoded by the decoder to generate molecules that satisfy specific conditions

Zhung 等^[50]提出了一种基于三维条件生成的 CVAE 模型 DeepICL, 专为基于蛋白质-配体相互作用的药物设计所提出。该模型采用 3D 分子结构作为表示方法, 将蛋白结合位点信息纳入生成过程, 以指导配体生成, 使其更具靶向性和高结合亲和力。Liu 等^[51]提出了面向离子

液体设计的 SDVAE 模型, 主要用于高效 CO₂ 捕集。该模型使用 SMILES 字符串表示离子液体, 通过 VAE 学习分子结构和性质的分布, 生成的离子液体在 CO₂ 溶解度上较数据集中最佳分子提高了 35.3%, 展现出在稀疏数据集上良好的生成效率和性能, 同时在合成难度上具有一

定的可行性。Lee 等^[52]提出了一种基于分子图的多目标优化生成式模型 MGCVAE，旨在实现分子物理性质的多目标设计。该模型利用条件 VAE，通过将目标属性嵌入到分子图表示的潜在空间中，从而生成满足多重理化性质的分子。在优化特定物理性质（如 $\log P$ 和折射率）的任务中，MGCVAE 实现了 25.89% 的目标分子生成率，显著高于无条件生成式模型。Munson 等^[53]设计了一个结合 VAE 和强化学习的多靶点生成式模型 POLYGON，专为多靶向药物设计提出。该模型通过 SMILES 字符串表示分子，并在潜在空间中引入目标属性优化，通过强化学习实现对靶点的特异性生成。该模型在 MEK1 和 mTOR 双靶点上表现出 82.5% 的靶点预测准确率，并在生成的分子中筛选出多个在实验中具备活性的候选药物。Romanelli 等^[54]提出的 Drug-CVAE 模型结合 SMILES 表示，专注于多靶点药物设计。该模型通过条件编码的方

式生成符合 CDK2、PPAR γ 、DPP-IV 等靶点需求的分子，评估表明生成的分子在有效性、新颖性及药物相似性等方面表现出色。

5.3 DDPM 模型

DDPM 在分子生成领域代表了生成建模的前沿进展。DDPM 通过逐步添加和去除噪声的方式，将分子生成过程分解为多步去噪，从而实现生成样本的高质量和高多样性，有效避免了生成式模型中常见的模式崩溃问题^[64]。这一特性使 DDPM 在生成高度真实且多样化的分子结构方面具有独特优势，尤其适用于复杂分子的精细生成^[65]。然而，DDPM 的主要缺点在于采样时间较长，生成过程需要经过多步去噪操作，相比于 GAN 和 VAE 的单步生成，其计算成本显著增加。此外，DDPM 模型的训练过程复杂，对计算资源的要求较高^[66, 67]，图 5 展示了条件生成 DDPM 模型的基本架构。

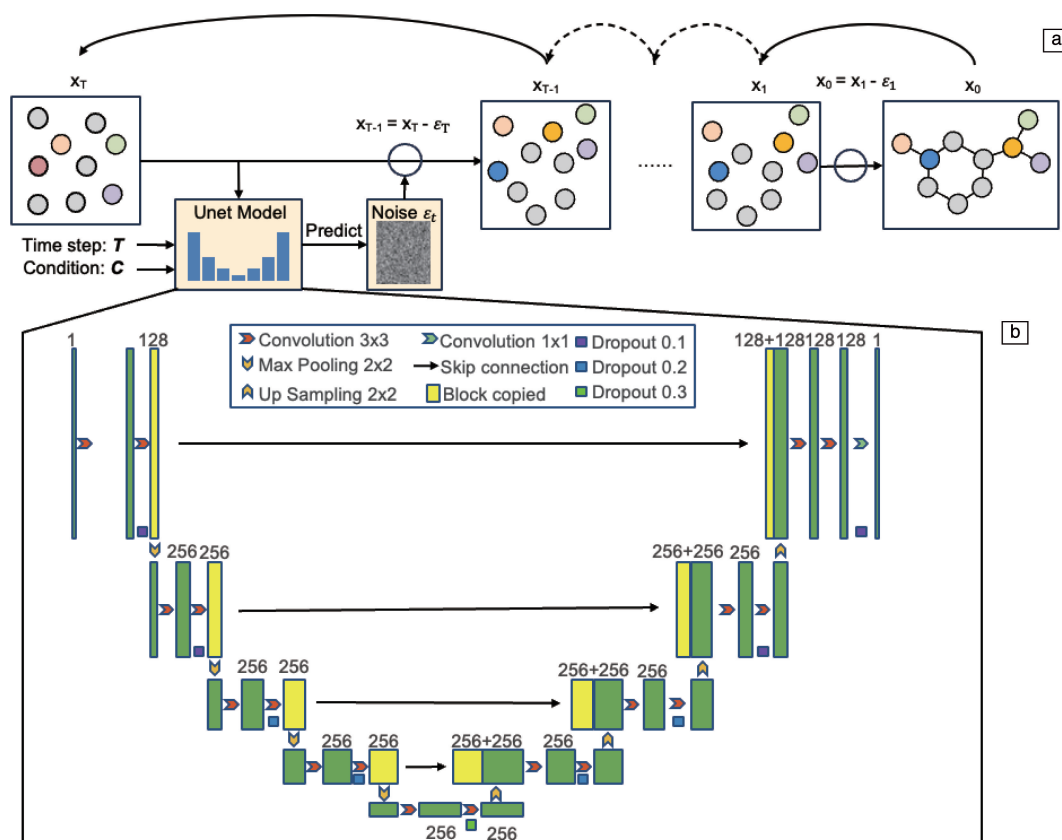


图 5 条件生成 DDPM 模型的示意图：在前向扩散过程中（弯曲箭头），扩散模型将逐渐向数据引入噪声，直到数据遵循各向同性高斯分布；在相反的过程中（直箭头），从噪声分布中采样获得的数据生成样本 (a)；DDPM 的噪声预测网络 U-Net：利用跳跃连接将编码器层的特征与相应的解码器层相连接；U-Net 接受带噪声的数据、时间步编码以及条件向量作为输入，输出对噪声的估计值 (b)

Fig. 5 Schematic illustration of conditional DDPM model; in the forward diffusion process (curved arrows), the diffusion model progressively adds noise to the data until it conforms to an isotropic Gaussian distribution, in the reverse process (straight arrows), data samples are generated by sampling from the noise distribution (a); noise prediction network U-Net used in DDPM; U-Net accepts noisy data, time step encodings, and condition vectors as inputs and outputs the estimated noise, skip connections are employed to link features from the encoder layers to the corresponding decoder layer (b)

Huang 等^[55]提出的 MDM 模型主要用于 3D 分子生成,特别针对较大分子的生成和多样性不足的问题。该模型通过引入原子间的力关系和分布控制变量,改善了分子结构的多样性和生成效果。在 QM9 和 GEOM-Drugs 数据集上的实验显示,MDM 在生成分子的独特性(31.4%)和新颖性(4.1%)方面超过了前沿的 EDM 模型^[56]。Morehead 等^[57]提出的 GCDM 模型针对 3D 分子生成过程中的几何稳定性进行了改进,特别适用于较大分子的生成。GCDM 在 GEOM-Drugs 数据集上生成的分子有效性达到了 92.5%,并在 QM9 数据集上的独特性和有效性分别达到 95%和 99.8%,表现出较高的生成稳定性和结构合理性。Bian 等^[58]提出的 HGLDM 是一种分层潜在扩散模型,可用于条件分子生成任务。该模型通过层次嵌入框架捕捉分子局部和整体的结构特征,支持多层次的条件生成。在特定分子属性的生成任务中,HGLDM 的生成效率较仅有图级别嵌入的 GLDM 模型提升了 24%,并在目标属性的生成率上提高了 15%。Xu 等^[59]提出的 GEOLDM 模型采用了基于潜在空间的扩散框架,以增强生成过程中的控制能力。GEOLDM 在 QM9 数据集上达到了 97%的有效性,并在大分子生成中实现了更高的结构新颖性。相比其他直接在原子空间中操作的模型,GEOLDM 通过更为简洁的潜在空间表示实现了有效性和多样性的提升。

5.4 生成模型实际应用案例

(1) 基于 VAE 的 CO₂ 高溶解度离子液体设计

Chen 等^[68]以离子液体 CO₂ 溶解度为设计目标,通过引入 VAE-ANN-PSO 框架(如图 6 所示)对离子液体分子进行逆向设计。他们首先利用超过 200 万条分子结构数据,将分子结构编码为高维潜在向量,训练了专用于阳离子与阴离子生成的 VAE 模型。随后,结合神经网络模型(ANN),实现了基于潜在空间的 CO₂ 溶解度高精度预测(测试集 MAE = 0.022),有效提升了模型的泛化能力。在此基础上,通过粒子群优化算法(PSO)于潜在空间中筛选生成新分子,共获得 5120 种候选离子液体分子。结果表明,新生成的离子液体分子不仅在结构上具有高新颖性(97.8%),且其 CO₂ 溶解度分布显著高于已知数据库中的。以 [BMIM][AC] 为对照,最佳生成离子液体分子的预测溶解度提升了近 60%,并经独立模型(ILTransR)再次验证。值得注意的是,所设计的离子液体分子在合成可行性评分(SAScore)方面与已知离子液体相当,显示出良好的潜在实用性。

(2) 基于 DDPM 设计具有高离子电导率的聚合物

Yang 等^[69]利用包含 6024 种不同无定形聚合物电解质的 HTP-MD 数据库,训练了 DDPM 模型实现新型聚合

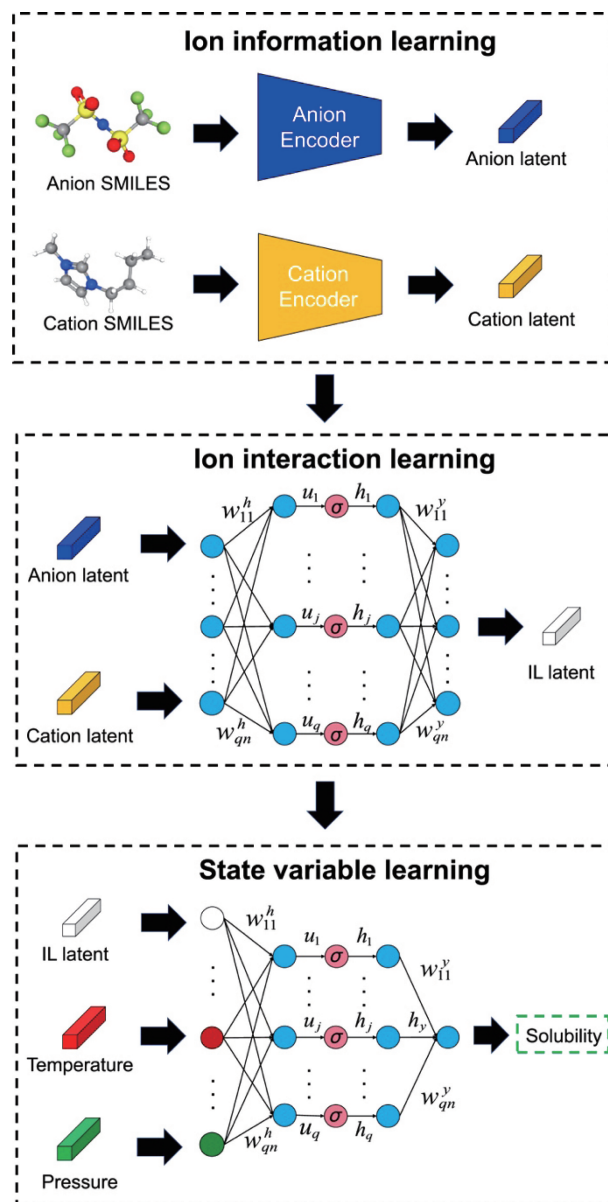


图 6 VAE-ANN-PSO 预测模型示意图^[68]

Fig. 6 Schematic diagram of the process in VAE-ANN-PSO prediction model^[68]

物分子的生成,并采用分子动力学(MD)模拟对生成分子的离子电导率进行高通量计算验证。在有条件生成任务中,通过引入高电导率标签,模型能够定向生成潜在高电导率的聚合物分子。最终,从 10 万个条件生成候选中筛选出 46 种最优分子,并通过 MD 模拟进一步验证其性能,如图 7a 所示。结果显示,17 种候选分子的离子电导率优于训练集中所有已知聚合物,其中最佳分子的电导率达到 1.13×10^{-3} S/cm,较已知最佳提高 1 倍以上,如图 7b 所示。此外,生成分子的可合成性与实际聚合物相当,进一步证实了生成模型在实际材料设计中的实用潜力。在模型生成和 MD 验证过程中亦存在部分失败案例。

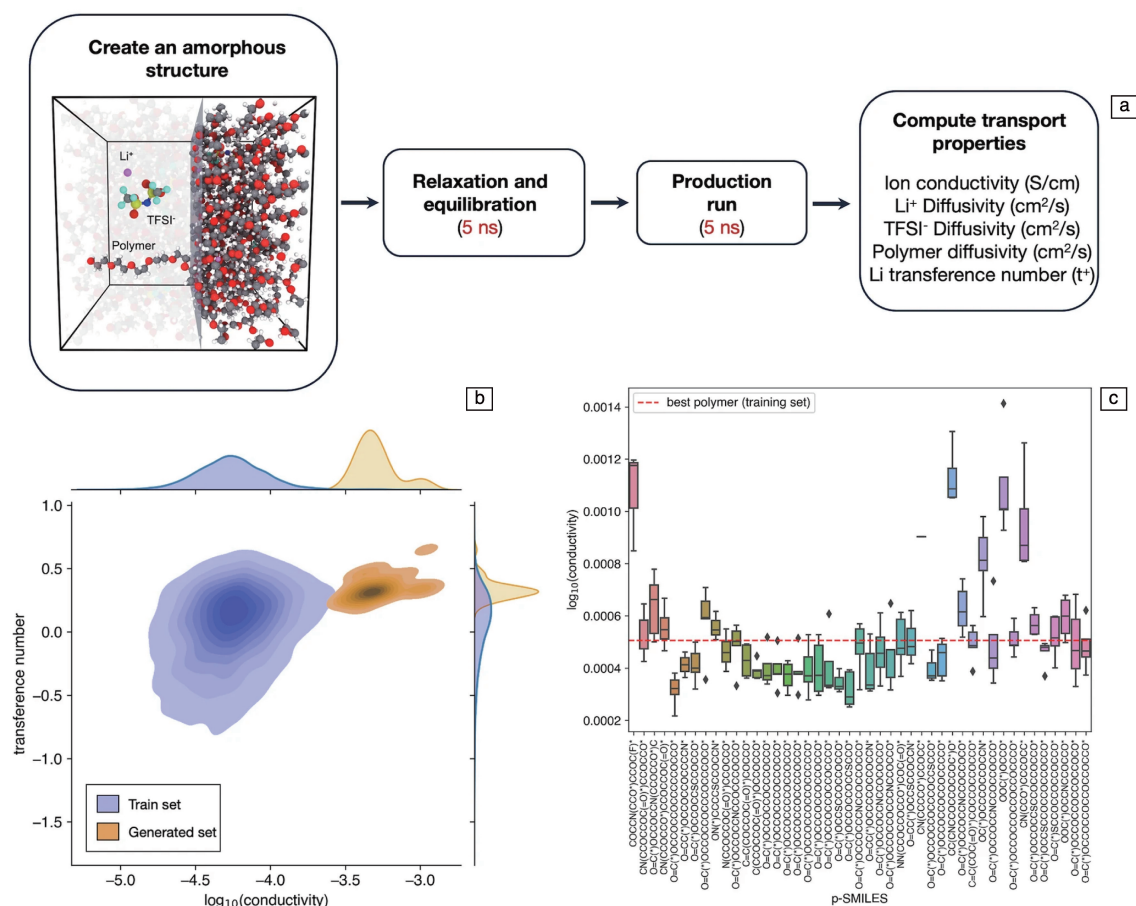


图 7 分子动力学模拟流程示意图：该流程首先构建无定形聚合物-盐体系，随后依次进行体系弛豫与平衡，接着进入正式的生产模拟阶段；在生产阶段采集的数据基础上，采用集群 Nernst-Einstein 方法计算离子传输性质 (a)。生成聚合物候选的电导率与迁移数分布：对最终筛选出的 46 种候选聚合物，通过分子动力学模拟获得其离子电导率与阳离子迁移数，并与训练集数据进行对比。整体候选集共包含 10 万个分子，文中展示的是其中性能最优的 46 种 (b)。各候选分子电导率的箱线图：绘制了 46 个候选聚合物的电导率箱线图，结果显示，有 17 个候选分子的电导率超过了训练集中的最佳聚合物 (5.1×10^{-4} S/cm) (c) [69]

Fig. 7 Schematic illustrating the MD simulation workflow: the process begins with the creation of an amorphous polymer-salt system, subsequent steps include system relaxation and equilibration, followed by the production run; ion transport properties are then determined using the cluster Nernst-Einstein method based on data collected during the production phase (a). Distribution of conductivity and transference number calculated using MD simulations for the top 46 candidates in the generated set (10^5 samples in total) in comparison to the training data (b). Box-and-whisker plot showing the conductivity values for each candidate; there are 17 out of 46 candidates showing better conductivity than the optimal polymer (conductivity = 5.1×10^{-4} S/cm) in the train set (c) [69]

例如，部分生成的分子因初始构型不稳定或体系剧烈波动而导致仿真失败。此类失败与局限提示我们，实际生成分子的可用性仍受限于分子 3D 结构而不仅仅是 2D 结构，未来可在模型中嵌入结合 3D 坐标迭代优化生成流程。

6 不同生成模型的计算资源需求

生成式深度学习模型在分子生成领域展现出显著的能力，但不同模型训练与应用过程中对计算资源的需求存在明显差异，对实际应用带来不同程度的限制。下文对 GAN、VAE 与 DDPM 模型在计算成本和实际应用方面进行简单比较。

首先，VAE 因其模型结构相对简单、收敛速度较快，通常训练时间较短，对计算资源的需求也较低，适合在有限算力条件下的大规模分子生成和快速原型验证。

GAN 在实际训练中由于判别器和生成器需交替优化，且易出现梯度不稳定、模式崩溃等问题，往往需要更多的训练轮次和调参时间。但整体而言，GAN 在单次训练迭代中的计算量仍低于 DDPM，且可通过网络结构优化和正则化手段在一定程度上降低资源消耗。

相较之下，DDPM 由于其生成和训练过程包含大量逐步正向和反向的噪声扰动，通常需数百至上千步的逐步更新。这一特性显著增加了模型的训练时间和显存消

耗,使其在资源有限或需实时响应的实际场景下应用受到较大限制。例如,在常用公开分子数据集上的实验表明,DDPM 的训练时长可为 VAE 的数倍,显存占用亦明显提升。

总体来看,VAE 在效率与可扩展性方面具有优势,适用于快速迭代和大数据量场景;GAN 在样本多样性和分布拟合能力上表现突出,但训练过程需平衡稳定性与资源开销;DDPM 则以生成质量和复杂结构建模能力见长,但高昂的计算成本对实际部署提出了更高要求。因此,在具体应用中,模型的选择需结合实际算力、任务目标与应用场景进行权衡。

7 结 语

生成式深度学习在分子设计领域的应用已取得显著进展。在分子表示方面,初期的研究多基于 SMILES 表示法,而后逐渐向更加直观的分子图表示拓展,近年来则进一步探索分子的三维结构生成。随着三维分子生成方法的出现,模型能够更精确地捕捉分子的立体特征,弥补了二维表示的不足。而不同的生成式模型在多样性、稳定性和特定属性控制方面各具优势。GAN 能够高效生成结构合理的分子;VAE 凭借其潜在空间的连续性在分子优化任务中表现突出;而 DDPM 则通过逐步去噪过程实现了高质量的分子生成,尤其适用于复杂分子结构的生成与优化。然而,这些模型在生成分子的合成可行性和属性控制的精确度等方面仍然存在挑战。未来的研究方向可能包括以下几点:首先,进一步改进生成式模型的合成可行性,确保生成的分子在实验上具有可操作性,这是推动生成式模型在实际应用中落地的关键。其次,探索多模型融合的生成框架,以利用不同模型的特长,更全面地覆盖化学空间。最后,构建更加全面、实用的评估体系,为生成式模型的评价和优化提供参考。

参考文献 References

- [1] ELTON D C, BOUKOUVALAS Z, FUGE M D, *et al.* *Molecular Systems Design & Engineering*[J], 2019, 4(4): 828-849.
- [2] PANG C, QIAO J, ZENG X, *et al.* *Journal of Chemical Information and Modeling*[J], 2023, 64(7): 2174-2194.
- [3] LI K, WANG J, SONG Y, *et al.* *Nature Communications*[J], 2023, 14(1): 2789.
- [4] POLLICE R, DOS PASSOS GOMES G, ALDEGHI M, *et al.* *Accounts of Chemical Research*[J], 2021, 54(4): 849-860.
- [5] CHOUDHARY K, DECOST B, CHEN C, *et al.* *npj Computational Materials*[J], 2022, 8(1): 59.
- [6] WALTERS W P, BARZILAY R. *Accounts of Chemical Research*[J], 2021, 54(2): 263-270.
- [7] DU Y, JAMASB A R, GUO J, *et al.* *Nature Machine Intelligence*[J], 2024, 6(6): 589-604.
- [8] GOODFELLOW I, POUGET-ABADIE J, MIRZA M, *et al.* *Advances in Neural Information Processing Systems 27*[C]. New York: Curran Associates Inc., 2014.
- [9] KINGMA D P. *Auto-Encoding Variational Bayes*[J/OL]. (2022-12-10)[2024-11-14]. <https://arxiv.org/abs/1312.6114>
- [10] HO J, JAIN A, ABBEEL P. *Advances in Neural Information Processing Systems 33*[C]. New York: Curran Associates Inc., 2020: 6840-6851.
- [11] BAGAL V, AGGARWAL R, VINOD P, *et al.* *Journal of Chemical Information and Modeling*[J], 2021, 62(9): 2064-2076.
- [12] MEYERS J, FABIAN B, BROWN N. *Drug Discovery Today* [J], 2021, 26(11): 2707-2715.
- [13] BILODEAU C, JIN W, JAAKKOLA T, *et al.* *Wiley Interdisciplinary Reviews: Computational Molecular Science*[J], 2022, 12(5): e1608.
- [14] ZENG X, WANG F, LUO Y, *et al.* *Cell Reports Medicine*[J], 2022, 3(12): 100794.
- [15] ZHANG J, MERCADO R, ENKQVIST O, *et al.* *Journal of Chemical Information and Modeling*[J], 2021, 61(6): 2572-2581.
- [16] WEININGER D. *Journal of Chemical Information and Computer Sciences*[J], 1988, 28(1): 31-36.
- [17] MSWAHLI M E, JEONG Y S. *Heliyon*[J], 2024, 10(20): e39038.
- [18] NGUYEN L A, HE H, PHAM-HUY C. *International Journal of Biomedical Science: IJBS*[J], 2006, 2(2): 85.
- [19] MORGAN H L. *Journal of Chemical Documentation*[J], 1965, 5(2): 107-113.
- [20] LI M M, HUANG K, ZITNIK M. *Nature Biomedical Engineering*[J], 2022, 6(12): 1353-1369.
- [21] ENGEL T, GASTEIGER J. *Cheminformatics: Basic Concepts and Methods*[M]. Weinheim: John Wiley & Sons, 2018: 200-215.
- [22] LI Z, JIANG M, WANG S, *et al.* *Drug Discovery Today*[J], 2022, 27(12): 103373.
- [23] DAVID L, THAKKAR A, MERCADO R, *et al.* *Journal of Cheminformatics*[J], 2020, 12(1): 56.
- [24] ZHENG S, HE J, LIU C, *et al.* *Nature Machine Intelligence* [J], 2024, 6(5): 558-567.
- [25] WIGH D S, GOODMAN J M, LAPKIN A A. *Wiley Interdisciplinary Reviews: Computational Molecular Science*[J], 2022, 12(5): e1603.
- [26] BAI Q, XU T, HUANG J, *et al.* *Drug Discovery Today*[J], 2024, 29(7): 104024.
- [27] RAGHUNATHAN S, PRIYAKUMAR U D. *International Journal of Quantum Chemistry*[J], 2022, 122(7): e26870.
- [28] GAULTON A, BELLIS L J, BENTO A P, *et al.* *Nucleic Acids Research*[J], 2012, 40(1): 1100-1107.
- [29] IRWIN J J, STERLING T, MYSINGER M M, *et al.* *Journal of Chemical Information and Modeling*[J], 2012, 52(7): 1757-1768.
- [30] KIM S, CHEN J, CHENG T, *et al.* *Nucleic Acids Research* [J], 2019, 47(1): 1102-1109.
- [31] RUDDIGKEIT L, van DEURSEN R, BLUM L C, *et al.* *Journal of Chemical Information and Modeling*[J], 2012, 52(11): 2864-2875.

- [32] RAMAKRISHNAN R, DRAL P O, RUPP M, *et al.* *Scientific Data* [J], 2014, 1(1): 1–7.
- [33] KANAKALA G C, DEVATA S, CHATTERJEE P, *et al.* *Current Opinion in Biotechnology*[J], 2024, 89: 103175.
- [34] WANG Y. *npj Computational Materials*[J], 2025, 11(1): 89.
- [35] WOŁOS A, KOSZELEWSKI D, ROSZAK R, *et al.* *Nature*[J], 2022, 604(7907): 668–676.
- [36] O'BOYLE N M, BANCK M, JAMES C A, *et al.* *Journal of Cheminformatics*[J], 2011, 3(1): 33.
- [37] RAMSUNDAR B, EASTMAN P, WALTERS P, *et al.* *Deep Learning for the Life Sciences: Applying Deep Learning to Genomics, Microscopy, Drug Discovery, and More*[M]. California: O'Reilly Media, Inc., 2019.
- [38] ADLER J, LUNZ S. *Advances in Neural Information Processing Systems 31*[C]. New York: Curran Associates Inc., 2018.
- [39] MAO X, LI Q, XIE H, *et al.* *Proceedings of the IEEE International Conference on Computer Vision 2017*[C]. Venice, Italy: Computer Vision Foundation, 2017: 2794–2802.
- [40] BROWN N, FISCATO M, SEGLER M H, *et al.* *Journal of Chemical Information and Modeling*[J], 2019, 59(3): 1096–1108.
- [41] POLYKOVSKIY D, ZHEBRAK A, SANCHEZ-LENGELING B, *et al.* *Frontiers in Pharmacology*[J], 2020, 11: 565644.
- [42] BAJUSZ D, RÁCZ A, HÉBERGER K. *Journal of Cheminformatics* [J], 2015, 7(1): 1–13.
- [43] ERTL P, SCHUFFENHAUER A. *Journal of Cheminformatics* [J], 2009, 1(1): 8.
- [44] COLEY C W, ROGERS L, GREEN W H, *et al.* *Journal of Chemical Information and Modeling*[J], 2018, 58(2): 252–261.
- [45] DE CAO N, KIPF T. MolGAN: An Implicit Generative Model for Small Molecular Graphs[J/OL]. (2022–9–27)[2024–11–14]. <https://arxiv.org/abs/1805.11973>
- [46] PRYKHODKO O, JOHANSSON S V, KOTSIAS P C, *et al.* *Journal of Cheminformatics*[J], 2019, 11(1): 74.
- [47] DAN Y, ZHAO Y, LI X, *et al.* *npj Computational Materials* [J], 2020, 6(1): 84.
- [48] MAZIARKA Ł, POCHA A, KACZMARCZYK J, *et al.* *Journal of Cheminformatics*[J], 2020, 12(1): 2.
- [49] MÉNDEZ-LUCIO O, BAILLIF B, CLEVERT D A, *et al.* *Nature Communications*[J], 2020, 11(1): 10.
- [50] ZHUNG W, KIM H, KIM W Y. *Nature Communications*[J], 2024, 15(1): 2688.
- [51] LIU X, CHU J, HUANG S, *et al.* *ACS Sustainable Chemistry & Engineering*[J], 2023, 11(24): 8978–8987.
- [52] LEE M, MIN K. *Journal of Chemical Information and Modeling*[J], 2022, 62(12): 2943–2950.
- [53] MUNSON B P, CHEN M, BOGOSIAN A, *et al.* *Nature Communications*[J], 2024, 15(1): 3636.
- [54] ROMANELLI V, ANNUNZIATA D, CERCHIA C, *et al.* *ACS Omega* [J], 2024, 9(43): 43963–43976.
- [55] HUANG L, ZHANG H, XU T, *et al.* *Proceedings of the AAAI Conference on Artificial Intelligence*[J], 2023, 37(4): 5105–5112.
- [56] HOOGEBOOM E, SATORRAS V G, VIGNAC C, *et al.* *Proceedings of International Conference on Machine Learning 2022*[C]. Baltimore, USA: IMLS, 2022: 8867–8887.
- [57] MOREHEAD A, CHENG J. *Communications Chemistry*[J], 2024, 7(1): 150.
- [58] BIAN T, NIU Y, CHANG H, *et al.* *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management* [C]. New York, United States: Association for Computing Machinery, 2024: 130–140.
- [59] XU M, POWERS A S, DROR R O, *et al.* *Proceedings of International Conference on Machine Learning 2023*[C]. New York, USA: IMLS, 2023: 38592–38610.
- [60] VALDEBENITO MATURANA C N, SANDOVAL OROZCO A L, GARCÍA VILLALBA L J. *Applied Sciences* [J], 2023, 13(19): 10637.
- [61] LIANG K J, LI C, WANG G, *et al.* *Generative Adversarial Network Training is a Continual Learning Problem*[J/OL]. (2018–11–27)[2024–11–14]. <https://arxiv.org/abs/1811.11083>.
- [62] NADERI H, SOLEIMANI B H, MATWIN S. *Proceedings of 2020 International Joint Conference on Neural Networks(IJCNN)* [C]. Glasgow, Scotland UK: IEEE, 2020: 1–8.
- [63] MIAO Y, YU L, BLUNSOM P. *Proceedings of International Conference on Machine Learning 2016* [C]. Seoul, South Korea: IMLS, 2016: 1727–1736.
- [64] GUO Z, LIU J, WANG Y, *et al.* *Nature Reviews Bioengineering*[J], 2024, 2(2): 136–154.
- [65] BASTEK J H, SUN W, KOCHMANN D M. *Physics-Informed Diffusion Models*[J/OL]. (2024–3–21)[2024–11–14]. <https://arxiv.org/abs/2403.14404>
- [66] CAO H, TAN C, GAO Z, *et al.* *IEEE Transactions on Knowledge and Data Engineering*[J], 2024, 36(7): 2814–2830.
- [67] CROITORU F A, HONDURU V, IONESCU R T, *et al.* *IEEE Transactions on Pattern Analysis and Machine Intelligence*[J], 2023, 45(9): 10850–10869.
- [68] CHEN X, CHEN G, XIE K, *et al.* *Green Chemical Engineering*[J], 2025, 6(3): 335–343.
- [69] YANG Z, YE W, LEI X, *et al.* *npj Computational Materials* [J], 2024, 10(1): 296.