

深度学习在分子生成中的应用进展

陈琳杰, 周瑞宁, 吕皓, 徐佳颖, 何正大, 陈亚东*

(中国药科大学理学院, 江苏 南京 211198)

[摘要] 分子设计中的药物设计是为了产生具有理想生物活性和物理化学性质的分子, 随着计算机科学与高性能计算的快速发展, 深度学习技术在药物设计领域的应用日益受到重视。生成式深度学习模型在自然语言、图像、音乐、视频等领域的表现卓越, 为分子生成提供了新的思路。越来越多的研究者开始尝试使用深度学习技术完成分子生成任务。综述总结了深度学习算法在分子生成中的研究进展, 重点介绍了常用的几种分子生成神经网络架构的原理、应用、分子表征形式及其技术细节。

[关键词] 人工智能; 深度学习; 药物设计; 分子生成

[中图分类号] TP39

[文献标志码] A

[文章编号] 1001-5094 (2023) 12-0950-11

DOI: 10.20053/j.issn1001-5094.2023.12.008

Application of Deep Learning in Molecule Generation

CHEN Linjie, ZHOU Ruining, LYU Hao, XU Jiaying, HE Zhengda, CHEN Yadong

(School of Science, China Pharmaceutical University, Nanjing 211198, China)

[Abstract] Drug design in molecular design aims to produce molecules with desirable biological activity and physicochemical properties. With the rapid development of computer science and high performance computing, the application of deep learning technologies in the field of drug design is gaining increasing recognition. Generative deep learning models have demonstrated remarkable performance in such fields as natural language, image, music, and video, providing new ideas for molecule generation. More and more researchers have started to use deep learning technologies to complete molecule generation tasks. This article summarizes the research progress of deep learning algorithms in molecule generation, focusing on the principles, applications, molecular representation forms, and technical details of several commonly used neural network architectures for molecule generation.

[Key words] artificial intelligence; deep learning; drug design; molecule generation

新药研发是具有成本高、风险大、周期长特点的复杂过程。从研发到最终上市需要投入数十亿美元和 10~15 年的时间, 然而新药研发的成功率很低, 仅约为 7.5%。通过使用计算机辅助药物设计手段, 可以生成并筛选候选化合物, 得到潜在的药物分子, 从而可以降低新药研发的成本、提高新药研发的成功率。计算机辅助药物设计不再特别依赖于药物化学家的经验, 但是非常依赖高质量的化合物库。目前有一些公开的化合物库如 ChEMBL^[1], PubChem^[2], ChemSpider^[3], 这些数据库的化合物数量一般在几百万个左右。然而, 潜在的类药物化合

物具有更为广阔的化学空间, 化合物数量估计在 $10^{23} \sim 10^{60}$ 之间。因此使用分子生成技术, 更有效地探索如此巨大的空间, 寻找潜在药物新分子非常有必要。

随着计算机科学与高性能计算的快速发展, 人工智能 (artificial intelligence, AI) 方法在图像处理、模式识别和自然语言处理等领域取得了成功。近年来, 机器学习尤其是深度学习, 被广泛应用于药物发现, 例如预测化合物的性质和活性以及它们与蛋白质靶点的相互作用。在过去几年里, 深度生成模型越来越受到关注, 其试图学习训练数据的概率分布, 提取有代表性的特征, 产生 1 个低维的连续表示, 最终通过从学习到的数据分布中采样来生成新的数据。生成模型已经应用在了图像、文本、语音和音乐等多个方向。生成模型的发展也为解决药物设计难题带来了新的思路, 被认为是最有前景的药物设

接受日期: 2023-03-20

项目资助: 国家自然科学基金 (No. 61806092)

*** 通信作者:** 陈亚东, 教授;

研究方向: 基于人工智能的药物分子设计及其应用研究;

Tel: 025-86185170; **E-mail:** ydchen@cpu.edu.cn

计方法之一。

当生成模型应用于生成分子时, 其本质是学习训练集中分子的分布, 从而获得与训练集中的分子相似但不同的分子集合; 也可通过结合进化算法或强化学习等算法, 生成具有特定生物活性或理化性质的分子。本文将按照模型进行分类, 对常见分子生成工作进行介绍, 并将介绍一些常用的分子生成数据集及分子生成的评价指标。

1 常见深度学习分子生成模型

深度学习分子生成模型的主要输入表示方式有 2 种: 文本序列和分子图。以文本作为输入的方式通常使用简化分子输入行记录系统 (simplified molecular input line entry system, SMILES) 进行表示。SMILES 具有唯一性和节省空间的特点, 可以使用非常少的数据表示出 1 个分子的信息。分子图表示的是分子的拓扑图, 与 SMILES 相比需要更多的空间, 因此包含了更多的信息。通常以 SMILES 为输入的生成模型包括: 循环神经网络 (Recurrent Neural Network, RNN); 以分子图作为输入的生成模型包括: 基于流的分子拓扑图生成模型; 对输入没有限制的生成模型包括: 自编码器模型、生成对抗网络 (Generative Adversarial Network, GAN)、Transformer 模型。

1.1 基于 RNN 的分子生成模型

1.1.1 RNN 理论知识 RNN 是自然语言处理领域兴

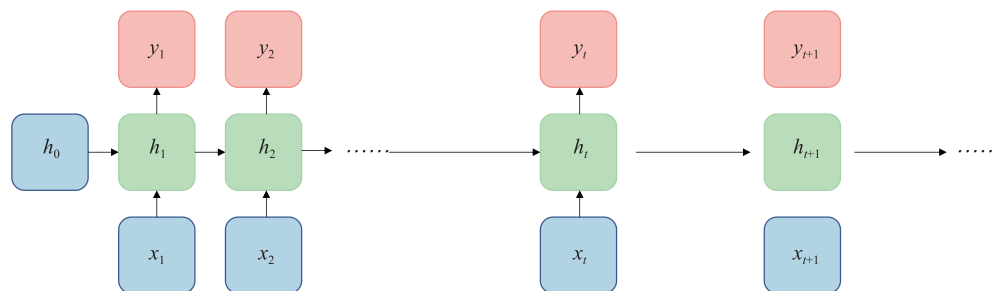


图 1 循环神经网络生成分子示意图

Figure 1 Illustration of molecule generation using recurrent neural networks

RNN 是特殊的神经网络, 它能够处理长序列数据, 并具有记忆性质, 能够记录之前时间步的信息并影响当前时间步的计算结果。RNN 是典型的递归网络, 其将上一时刻的隐层输出作为当前时刻的隐

起的模型, 现在也已广泛应用于多个领域。Jordan^[4]和 Elman^[5]提出的简单循环网络模型被认为是当前 RNN 的基础版本。如图 1 所示, RNN 通过隐藏层上的循环连接, 可以在当前时间接收到前一时间隐藏单元的状态, 并且可以进一步更新当前时间的隐藏单元状态。RNN 隐藏单元 h_t 在时间 t 接收来自 2 个方面的数据, 分别是上一次的隐藏单元值 h_{t-1} 和当前输入数据 x_t , 2 个输出分别是通过下式计算得到的隐藏单元的值 h_t 和输出向量 y_t :

$$h_t = f(h_{t-1}, x_t)$$

$$y_t = O(h_t)$$

RNN 在通过反向传播更新网络中的参数时, 可能会发生“梯度爆炸”和“梯度消失”现象。这些问题已通过对 RNN 的结构进行调整得以解决。例如长短期记忆神经网络 (long short-term memory, LSTM) 和门控循环单元 (gate recurrent unit, GRU), 它们的内部结构更复杂, 有助于选择性地存储和更新信息。Schmidhuber 等^[6]提出了首个 LSTM 单元, 该单元具有用于输入、遗忘和输出的受控门。精心制作的“门”结构用于删除或增强隐藏信息。LSTM 单元使用更可控的信息流来确定哪些信息可以保留, 哪些可以丢弃。LSTM 实现了更精细的内部处理单元, 可以保持其内部状态以延长 RNN 中顺序输入的时间, 从而提高 RNN 的性能。GRU^[7]是 LSTM 架构的简化实现, 可以以较低的计算成本缓解梯度消失和爆炸的问题。

层输入, 因此必须按顺序计算。若未计算出上一时刻的隐层输出, 则无法计算下一时刻的隐层输出, 这也是 RNN 无法并行计算的原因。RNN 的深度和参数数量越大, 其处理能力就越强, 但同时也会增

加计算的复杂度, 故在大批量分子生成任务中效果可能一般。然而在小数据集的分子生成任务中, RNN 的优势就能得到体现。

1.1.2 RNN 的常见分子生成模型 RNN 在分子生成领域通常使用分子的文本序列作为输入, Bjerrum 等^[8]通过在 SMILES 序列上使用 RNN, 来找出合理的化学规则并生成可合成的分子。该工作产生的分子与它的训练集 ZINC 数据集相比, 分子的摩尔质量、LogP 值、氢键受体和供体的数量、可旋转键的数量和拓扑极性表面积方面具有非常相似分布。根据合成可及性 (synthetic accessibility, SA) 评分和 Wiley ChemPlanner 评估, Bjerrum 等人建立的模型所生成的化合物在大多数情况下是可合成的。Grisoni 等^[9]也在 SMILES 序列上使用了 RNN, 同时该作者提出了基于 SMILES 数据增强的新方法, 即交替学习的双向分子设计 (bidirectional molecule design by alternate learning, BIMODAL)。应用该策略使得生成的分子在新颖性、骨架多样性和化学-生物相关性方面都取得了更佳的效果。

RNN 还可以与迁移学习结合, 生成具有特定生物活性的分子。Shi 等^[10]基于 GRU 的神经网络结合迁移学习, 成功建立了 ADAM10 抑制剂的分子生成模型。该模型在使用 GRU 进行分子生成时

只需要化学配体的 SMILES 信息, 并可以有效地生成大量潜在的高生物活性新结构。

此外, RNN 也可以与强化学习相结合, Neil 等^[11]、Goel 等^[12]、Olivecrona 等^[13]所做的 3 项工作都是将强化学习与 RNN 相结合。使用强化学习可以生成对指定靶标具有高结合亲和力的分子并获得其他理想的理化特性。例如, Goel 等人设计了新颖的策略, 其中用于强化学习的奖励函数周期性改变, 使用不同的奖励策略, 最终可以生成针对目标靶点具有高亲和力的分子, 该方法可以针对需要的靶点生成潜在药物。

阿斯利康公司的工作将 RNN 与课程学习相结合, 该工作在 2 个任务上进行了实验, 分别是: 构建相对复杂的支架和满足分子对接约束^[14]。在这 2 项任务中阿斯利康公司提出的模型效果都好于其他强化学习策略。

1.2 基于自编码器的分子生成模型

1.2.1 自编码器原理 自编码器由 2 个网络组成 (见图 2), 编码器将高维数据映射到低维表示, 解码器在给定低维表示的情况下将原始输入重构为输出。自编码器被反复训练以最小化重建输出与原始输入之间的偏差, 自编码器的目标是找到更紧凑的样本表示。

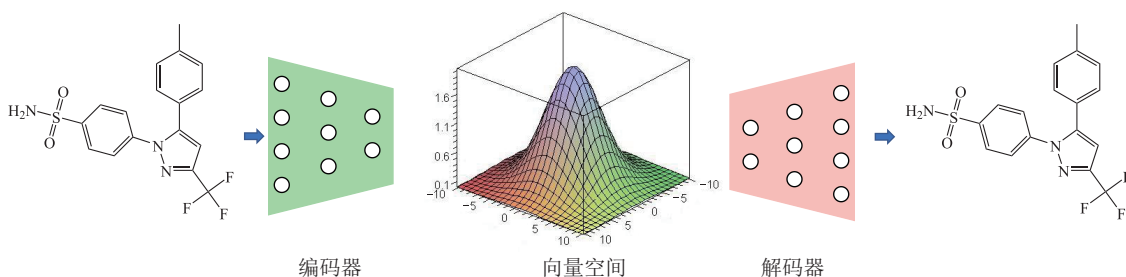


图 2 自编码器模型图

Figure 2 Illustration of molecule generation using autoencoder

变分自编码器 (variational auto-encoder, VAE) 通过一些额外的约束来修改经典自编码器, 从输入数据中学习潜在表示。与经典自编码器的目标不同, VAE 旨在学习数据集的概率分布, 从而生成与数据集相似但不同的样本。2013 年, Bengio 等^[15]提出了首个 VAE, VAE 中编码器和解码器的输出分别是数据在潜在空间和初始空间中的概率分布。在 VAE

中, 连续表示输入 z 被解释为潜在变量, $P(z)$ 是遵循高斯分布的先验分布。概率解码器由具有参数 θ 的似然函数 $P\theta(x|z)$ 定义, 编码器使用由 ϕ 参数化的模型 $q\phi(z|x)$ 近似后验分布。VAE 的目标是通过使用概率公式最大化训练集中每个 x 的概率。该概率公式如下:

$$P(x) = \int_z P(x)P(x|z)dz$$

自编码器在训练时, 需要花费更多的时间和计算资源。由于自编码器为无监督学习方法, 其学到的特征可能不具有物理意义, 即不能很好地被解释, 因此, 自编码器在分子生成任务中的可解释性较差。不过, 如果数据集更大, 自编码器的性能往往更好, 原因在于其可以更好地捕捉到数据的潜在结构和特征。

1.2.2 基于自编码器分子生成的工作 自编码器既能够以分子拓扑图作为输入, 也能够以分子序列信息作为输入。Berenger 等^[16]提出了生成具有优化特性的模型, 可以将评分函数与分子生成器相结合, 以设计具有所需特性的新分子。同时该工作具有较高的生成速度, 在使用 DeepSMILES 时该方法达到了峰值性能(每秒 340 000 余个分子)。Alperstein 等^[17]提出了使用 1 组堆叠的 RNN 对单个分子的多个 SMILES 字符串进行编码, 在 SMILES 表示之间汇集每个原子的隐藏表示, 并使用注意力池来构建最终的固定-长度潜在表示。这样使得该模型在较低计算复杂度的条件下, 取得了与以分子图作为输入媲美的性能。Polykovskiy 等^[18]提出了对抗式自编码器(adversarial autoencoder, AAE), 在这项工作中, 他们应用 AAE 模型生成 Janus 激酶 3(Janus kinase 3, JAK3) 抑制剂, 并发现了 1 种有开发前景的高活性化合物, 该化合物显示出良好的体外活性和选择性。Zavoronkov 等^[19]提出用于从头小分子设计的深层生成模型——生成张量强化学习模型(Generative Tensorial Reinforcement Learning, GENTRL)。该模型在自编码器的基础上结合了强化学习、变分推理和张量分解算法。GENTRL 成功地用于发现盘状结构域受体 1(discoidin domain receptor 1, DDR1) 的有效抑制剂, DDR1 是与纤维化和其他疾病相关的激酶靶点。从新生成的化合物中鉴定出 6 种候选先导化合物, 其中 1 种化合物在小鼠中显示出良好的疗效和药代动力学特性。在这项工作中, 使用 GENTRL 在 21 天内发现了 DDR1 的有效抑制剂(见图 3), 并在 46 天内完成了设计、合成和实验测试, 证明了该方法用于快速有效分子设计的潜力。

Jin 等^[20]实现了以分子图作为自编码器的输入分子生成工作。该工作提出的连接树 VAE 分 2 个

阶段生成分子图, 首先在化学子结构上生成树结构支架, 然后将它们组合成具有图消息传递网络的分子。这种方法能够逐步扩展分子, 同时在每一步都保持化学有效性。Liu 等^[21]也提出了 VAE 模型, 其中编码器和解码器都是图结构的。该工作的解码器假设图形扩展步骤按顺序排列, 且该模型通过使用潜迁移学习的方法, 可以生成具有特定性质的分子。Chenthamaraksha 等^[22]提出了名为受控分子生成(Controlled Generation of Molecule, CogMol)的生成模型, 通过在自编码器中引入受控采样模式, 设计一系列靶向严重急性呼吸综合征冠状病毒 2(severe acute respiratory syndrome coronavirus 2, SRAS-CoV-2)的分子。该模型生成的分子同时受到靶向性和选择性、药物相似性、合成可行性和毒性的限制。结果表明, 生成的分子能够很好地结合到目标结构的相关药物袋中, 并显示出低的预测代谢物毒性和高的合成可行性。



图 3 GENTRL 模型生成的 DDR1 激酶抑制^[19]

Figure 3 DDR1 kinase inhibition generated by the GENTRL model^[19]

1.3 基于 GAN 的模型

1.3.1 GAN 模型原理 GAN 的概念由 Makhzani 等^[23]于 2015 年首次提出。如图 4 所示 GAN 模型包括生成器 G 和判别器 D。通常, 生成器是二元分类器, 其学习将随机噪声映射到需要接近数据分布的特定分布, 判别器则确定输入是真实数据还是生成器生成的样本。一旦模型训练完成, 就可以从生成器中获得新的样本。具体来说, 在对抗过程中, 同时训练 2 个神经网络模型, 生成器 G 和判别器 D, 使 D

可以在输入数据中找到隐藏模式, 从而准确地区分真实数据和 G 生成的数据, 而 G 将通过优化数据采样的矩阵乘法的权重进行迭代, 以学习欺骗训练有素的 D。总的来说, GAN 模型的本质是 D 和 G 相互竞争。下面展示了 GAN 的目标函数:

$$\min_G \max_D E_{x \sim p_{data}(x)} [\log D(x)] + E_{z \sim p(z)} [\log (1 - D(G(z)))]$$

在上式中, $p_{data}(x)$ 是真实数据分布, $p(z)$ 是先验概率分布。训练判别器 D 以最大化正确区分训练

样本和来自生成器 G 的样本的概率, 同时训练生成器 G 以最小化判别器 D 能够区分真假样本的概率。条件生成对抗网络 (ConditionalGAN) 是 GAN 的变体, 它通过在生成器和判别器中添加额外信息 c 来调节。条件向量 c 和输入噪声 z 被输入到生成器中, 在鉴别器中, 与训练样本连接的条件向量 c 为输入。目标函数表示为:

$$\min_G \max_D E_{x \sim p_{data}(x)} [\log D(x|c)] + E_{z \sim p(z)} [\log (1 - D(G(z|c)))]$$

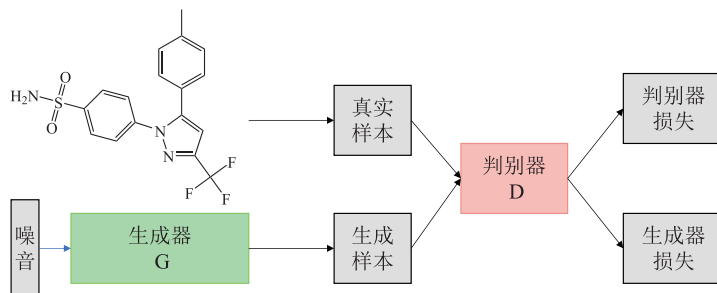


图 4 生成对抗网络生成分子示意图

Figure 4 Illustration of molecule generation using generative adversarial networks

训练 GAN 需要达到纳什均衡, 即生成器和判别器的损失函数的收敛达到最优值。然而, 达到纳什均衡是较为困难的过程, 原因是 GAN 是对抗学习模型, 其中生成器和判别器同时进行训练, 并试图对方训练得更好。有时候, 可以用梯度下降法训练 GAN 以达到纳什均衡, 但有时因为训练环境的不稳定性, 该法并不能达到纳什均衡。由于 GAN 模型具有一定的不可控性, 因此需要对其进行严格的调整以使其训练效果更佳。此外, 由于 GAN 模型是易于添加约束条件的生成模型, 其可以生成具有特定性质的分子, 因此, GAN 模型适合作为条件生成模型, 用于生成具有特定性质的分子。但是, 由于 GAN 模型的不可控性, 需对其进行严格的调整, 以确保生成的分子具有特定的性质。

1.3.2 基于 GAN 模型的工作 GAN 模型可以以分子拓扑图或 SMILES 作为输入。Putin 等^[24]提出了以 SMILES 作为输入的基于深度神经网络对抗性阈值神经计算机 (Adversarial Threshold Neural Computer, ATNC) 的模型, 该模型是结合 GAN 架构和强化学习实现分子生成工作。ATNC 使用可微分神经计算机作为生成器, 并有 1 个新的特定模

块, 称为对抗阈值 (adversarial threshold, AT)。AT 充当代理 (生成器) 和环境 (判别器 + 目标奖励函数) 之间的过滤器。对关键分子描述符和化学统计特征的分析表明, ATNC 产生的分子具有较高成药性。Guimaraes 等^[25]提出了以 SMILES 字符作为输入的对抗神经网络 SeqGAN, 该工作结合强化学习, 提出了利用对抗性训练和基于专家奖励的生成对抗网络框架, 将数据生成器建模为强化学习设置中的随机策略; 该工作还扩展了训练过程, 以包括判别器奖励之外的特定领域目标。以上 2 种奖励的混合可以通过可调参数来控制。同时该模型提出了用 Wasserstein 距离作为判别器的损失函数, 使得其训练过程更加稳定。

Cao 等^[26]提出了用于小分子图的隐式、无似然生成模型 MolGAN (Molecular GAN), 其规避了先前基于似然的方法对昂贵的图匹配程序或节点排序启发式的需求。该模型采用 GAN 直接对图结构数据进行操作, 同时将对抗生成网络与强化学习目标相结合, 以生成具有特定所需化学性质的分子。

1.4 基于 Transformer 的分子生成模型

1.4.1 Transformer 的理论知识 Transformer 是在自然

语言处理领域提出的新模型^[27]。Transformer 的原始版本由编码器和解码器组成, 如图 5 所示, 在进行分子生成任务时通常只有解码器部分参与。该模型的关键是注意力机制, 其可以考虑序列中的长期依赖关系。详细来说, 有 key (K), query (Q) 和 value (V) 3 个向量, 对应的 attention 表示如下:

$$\text{attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

其中 d_k 是键和查询向量的维度, 用于缩放这些向量的点积。

考虑到注意力机制中不包含序列中标记的顺序, 因此将额外的位置信息注入到输入中。具体来说, 正弦和余弦函数以下列公式的形式使用:

$$PE_{(\text{pos}, 2i)} = \sin \left(\frac{\text{pos}}{10\,000^{2i/d_{\text{model}}}} \right)$$

$$PE_{(\text{pos}, 2i+1)} = \cos \left(\frac{\text{pos}}{10\,000^{2i/d_{\text{model}}}} \right)$$

其中 pos 代表位置, i 代表维度, d_{model} 是嵌入

的大小。Transformer 比传统序列模型如 RNN/LSTM 具备优势在于其强大的并行计算能力。

对于 RNN 来说, 任意时刻 t 的输入是时刻 t 的输入 $x(t)$ 和上一时刻的隐藏层输出 $h(t-1)$, 经过运算后得到当前时刻隐藏层的输出 $h(t)$, RNN 的历史信息需要通过这个时间步一步一步向后传递, 而这就意味着 RNN 序列后面的信息只能等到前面的计算结束后, 将历史信息通过隐藏层传递给后面才能开始计算, 形成链式的序列依赖关系, 无法实现并行。

相对而言, Transformer 是具有很高计算效率的神经网络结构, 它在 self-attention 层可以同步计算所有单词之间的注意力关系, 无论序列长度如何, 都可以实现并行计算。另外, Transformer 由于采用了 Multi-head attention 结构和计算机制, 具有比 RNN 更强大的特征抽取能力。然而, Transformer 由于对计算设备的要求更高, 模型的训练时间也更长, 并且由于其强大的特征抽取能力, 生成分子的新颖性往往较 RNN 类模型低。因此, Transformer 更适合大数据集分子生成任务。

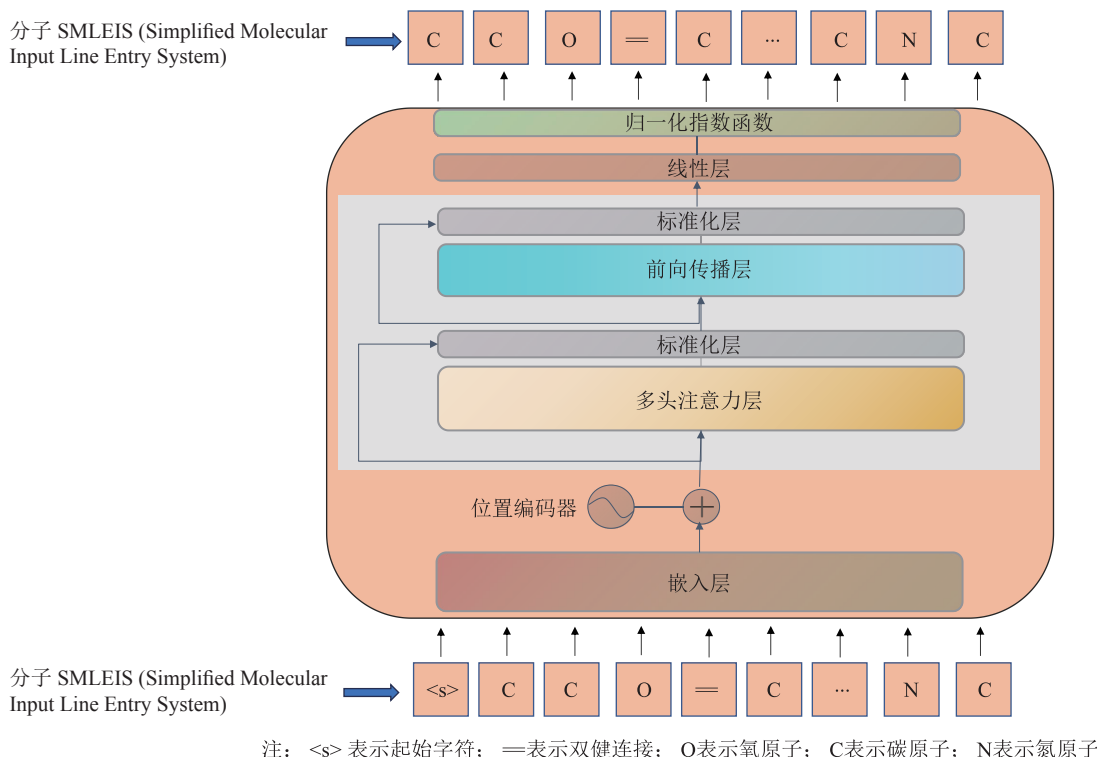


图 5 Transformer 模型生成分子示意图

Figure 5 Illustration of molecule generation using Transformer

1.4.2 基于 Transformer 的分子生成模型 Transformer 生成模型既可以使用分子序列作为输入, 也可以使用分子拓扑图作为输入。Bagal 等^[28]提出的分子生成预训练 (Molecular Generative Pre-training Transformer, MolGPT) 模型将 SMILES 输入给 Transformer 用于一般的分子设计。受已被证明可成功生成有意义文本的生成预训练 (Generative Pre-training Transformer, GPT) 模型的启发, MolGPT 中使用掩蔽自注意力来训练 Transformer 解码器以处理下一字符任务, 从而生成药物分子。MolGPT 在生成有效、独特和新颖的分子方面与其他先前提出的用于分子生成的现代机器学习框架的性能相当。此外, 该模型可以有条件地训练以控制生成分子的多种特性。

Wang 等^[29]提出了多约束分子生成模型 (Multi-Constraint Molecular Generation, MCMG), 该方法可以通过知识蒸馏 (knowledge distillation, KD) 结合条件 Transformer 和强化学习算法来满足多个约束。条件 Transformer 通过有效地学习并将结构性关系合并到有偏差的生成过程中来训练分子生成模型; 然后使用 KD 模型来降低模型的复杂性, 以便可以通过强化学习有效地对其进行微调并增强生成分子的结构多样性。Zhu 等^[30]提出了端到端的 SMILES 转换器 (SMILES Transformer, ST)-KD, 用于 KD 促进的分子表示学习。为了进行从图 Transformer 到 ST-KD 的知识转移, 研究者重新设计了注意力层并引入了预转换步骤来标记 SMILES 字符串, 同时注入基于结构的位置嵌入。无需昂贵的预训练, ST-KD 在最新的标准分子数据集 PCQM4M-LSC 和 QM9 上显示出具有竞争力的结果, 与现有图模型相比, 推理速度提高了 3~14 倍。

Mitton 等^[31]提出使用拓扑图作为 Transformer 的输入, 该模型充分利用图卷积和图池化层直接对图进行操作。图 Transformer 模型实现了新颖的节点编码层, 取代了通常在 Transformer 中使用的位置编码, 以创建没有位置信息的 Transformer, 对图进行操作, 将相邻节点的属性编码到边缘生成过程中。作者所提出的模型系在对具有边缘特征的图进行操作的图生成工作之上所创建, 该模型通过图中的节点数量提供了改进的可扩展性。此外, 该模型能够

通过潜在变量和图属性之间的映射来学习表示图属性的解耦和可解释的潜在空间。

1.5 基于流的生成模型

1.5.1 基于流的生成模型理论知识 基于流的生成模型是由一系列的可逆变换器组成^[32]。如图 6 所示, 基于流的生成模型可以使得模型能够更加精确地学习到数据分布, 它的损失函数是负对数似然函数。令 $f: \epsilon \rightarrow Z$ 是 1 个可逆变换, 其中 $\epsilon \sim p_\epsilon(\epsilon)$ 是基分布, 那么可以通过变量的变化来计算真实世界数据的密度函数, 即公式:

$$p_z(z) = p_\epsilon(f_\theta^{-1}(z)) \left| \det \frac{\partial f_\theta^{-1}(z)}{\partial z} \right|$$

基于流的生成模型有 2 个关键过程: 其一为计算数据似然性, 即给定 1 个数据点 z , 可以通过反转变换 $f, \epsilon = f_\theta^{-1}(z)$ 来计算精确密度 $p_z(z)$; 其二为采样, z 可以从分布 $p_z(z)$ 中通过首个样本 $\epsilon \sim p_\epsilon(\epsilon)$ 进行采样, 然后进行前馈变换 $z = f_\theta(\epsilon)$ 为了有效地执行上述操作, 需要 f_θ 是可逆的, 具有易于计算的雅可比行列式。Papamakarios 等^[33]提出的自回归流 (Autoregressive Flow, AF) 模型是满足这些标准的变体, 它包含 1 个三角雅可比矩阵, 并且行列式可以线性计算。形式上, 给定 $z \in R^D$ (D 是观察数据的维度), 自回归条件概率可以参数化为高斯分布:

$$p(z_d | z_{1:d-1}) = N(z_d | \mu_d, (\alpha_d)^2),$$

$$\text{where } \mu_d = g_\mu(z_{1:d-1}; \theta),$$

$$\alpha_d = g_\alpha(z_{1:d-1}; \theta)$$

其中 g_μ 和 g_α 分别是 $z_{1:d-1}$ 的无约束和正标量函数, 用于计算均值和偏差。在实践中, 这些功能可以实现为神经网络。AF 的仿射变换可以写成:

$$f_\theta(\epsilon_d) = z_d = \mu_d + \alpha_d \cdot \epsilon_d; f_\theta^{-1}(z_d) = \epsilon_d = \frac{z_d - \mu_d}{\alpha_d}$$

AF 中的雅可比矩阵是三角形的, 因为 $\partial z_i / \partial \epsilon_j$ 仅当 $j \leq i$ 时非零。因此, 可以有效地计算行列式。具体来说, 为了进行密度估计, 可以并行应用所有单独的标量仿射变换来计算基础密度, 每个基础密度都取决于先前的变量 $z_{1:d-1}$; 在对 z 进行采样时, 可以先采样 $\epsilon \in R^D$ 并通过仿射变换计算 z_1 , 然后可以

根据之前观察到的 $z_{1:d-1}$ 依次计算每个后续 z_d 。

流模型在分子生成任务中表现出色, 是近年来新兴生成模型中的佼佼者。它与自编码器和 GAN

相比, 具有一定的性能优势, 但也对计算资源有更高的要求。尽管如此, 流模型的优点使其成为分子生成等任务中最有潜力的模型之一。

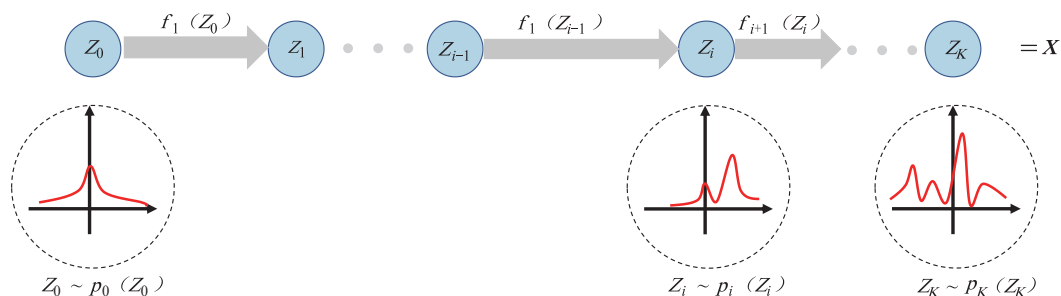


图 6 流模型生成分子示意图

Figure 6 Illustration of molecule generation using flow-based models

1.5.2 常见基于流的分子生成模型 基于流的分子生成模型通常以分子拓扑图作为输入, Shi 等^[34]提出的基于流的自回归模型图生成模型 GraphAF, 该模型同时具有自回归和流模型的优点; 同时, 其还可以使用强化学习对目标导向特性优化模型进行微调, 产生具有特定生物活性的分子。Luo 等^[35]提出了基于归一化流方法的新型离散潜变量模型 GraphDF, 该模型可用于分子图生成, 解决了离散图结构的建模不准确的问题。GraphDF 使用可逆的模移位变换将离散的潜在变量映射到图节点和边。其结果表明, 使用离散潜变量能降低计算成本并能消除去量化的负面影响。Kuznetsov 等^[36]提出了用于生成分子图的层次归一化流模型。该模型通过递归地将每个节点一分为二, 从单节点图生成新的分子结构。所有操作均可逆, 可以作为即插即用模块使用。该模型还可以对化学性质进行全局和约束优化, 生成具有特殊性质的分子。

表 1 列举了以上各种生成模型的经典工作。

2 分子生成常用数据集

分子生成任务需要通过学习现有数据集的分布从而生成新的分子, 因此使用的数据集质量与最终生成分子质量息息相关, 本章将介绍 ZINC, QM9 和 ChEMBL 3 个在分子生成任务中常用的分子数据集。

ZINC 是用于虚拟筛选的商用化合物的免费数据库^[37]。该数据集包含 2.3 亿个分子, 每个分子都具有 3D 结构。这些分子包含了生物学相关的质子

化状态, 并包含化学信息如相对分子量、计算的 LogP 值和可旋转键的数量等。ZINC 数据集集中的每个分子都包含供应商和采购信息, 并且可以使用很多对接程序进行对接。该数据库可免费下载 (<http://zinc.docking.org>), 有多种常见文件格式, 包括 SMILES, mol2, 3D SDF 和 DOCK flexibase 格式。

QM9 数据集报告了由碳、氢、氧、氮、氟原子组成的 134 000 个稳定有机小分子, 以及它们的计算几何、能量、电子和热力学特性^[38]。这个数据集提供了相关的、一致的和全面的有机小分子化学空间的量子化学特性。这个数据库可以作为现有方法的基准数据, 或用以开发新的方法, 如量子力学/机器学习, 以及构效关系的系统识别。

ChEMBL 数据集是经人工整理的具有类药物特性的生物活性分子的化学数据库。其由位于英国欣克斯顿 Wellcome Trust Genome Campus 的欧洲分子生物学实验室 (European Molecular Biology Laboratory, EMBL) 的欧洲生物信息学研究所 (European Bioinformatics Institute, EBI) 维护。ChEMBL 数据集包含针对药物靶点的化合物生物活性数据 (如 K_i , K_d , IC_{50} 和 EC_{50}), 包括 240 万次测量数据, 涵盖 622 824 种化合物 (含 24 000 种天然产物)。以上数据来源于 12 种药物化学期刊的 34 000 多篇文献。

3 分子生成常用评价指标

优秀的分子生成工作需要同时具备 2 个特点:

其一, 生成的分子是有效唯一的新分子; 其二, 生成的分子具有合理的化学属性, 因此通常采用有效性、唯一性、新颖性等指标对生成分子进行评价。

有效性是指生成的分子是否满足基本的规则, 是否是理论上存在的分子。通常通过 Rdkit 库能否识别作为分子是否有效的衡量方式。

表 1 基于深度学习的分子生成经典工作

Table 1 Tasks for molecular generation based on deep learning

作者及参考文献	模型	表征	数据集	发表时间/年份
Bjerrum E J ^[8]	RNN	SMILES	ZINC	2017
Grisoni F ^[9]	RNN	SMILES	CHEMBL	2020
Shi T ^[10]	RNN	SMILES	ChEMBL	2020
Neil D ^[11]	RNN	SMILES	CHEMBL	2018
Goel M ^[12]	RNN	SMILES	CHEMBL	2021
Olivecrona M ^[13]	RNN	SMILES	CHEMBL	2017
Berenger F ^[16]	自编码器	SMILES	ChEMBL	2021
Alperstein Z ^[17]	自编码器	SMILES	ZINC	2019
Jin W G ^[20]	自编码器	Graph	ZINC	2019
Liu Q ^[21]	自编码器	Graph	ZINC	2019
Putin E ^[24]	GAN	SMILES	ZINC	2018
Guimaraes G L ^[25]	GAN	SMILES	ZINC	2018
Cao N D ^[26]	GAN	Graph	QM9	2018
Bagal V ^[28]	Transformer	SMILES	MOSES	2022
Wang J ^[29]	Transformer	SMILES	MOSES	2021
Zhu W H ^[30]	Transformer	SMILES	QM9	2021
Mitton J ^[31]	Transformer	Graph	QM9	2021
Shi C ^[34]	流	Graph	ZINC	2020
Luo Y ^[35]	流	Graph	ZINC	2021
Kuznetsov M ^[36]	流	Graph	MOSES	2021

RNN: Recurrent Neural Networks (循环神经网络); GAN: Generative Adversarial Networks (生成对抗网络); SMILES: simplified molecular input line entry system (简化分子输入行记录系统)

唯一性是指在分子生成任务中, 研究者们希望生成的分子尽可能不一致, 这样可以覆盖更大的化学空间。因此采用唯一性指标判断生成唯一分子占总分子数的比例。

新颖性指生成与训练集不同分子占生成总分子数比例, 在分子生成任务数据集中会包含已知药物分子, 因此研究者希望生成的分子学习到了药物分子的部分结构, 同时又与已知药物结构不同。

定量评估类药性 (quantitative estimate of drug-likeness, QED) 是将药物相似性量化为介于 0 和 1 之间的数值的性质。QED 由以下 8 个指标计算得到: 相对分子质量、亲脂性、氢键供体的数量、氢键受体的数量、极表面积、可旋转键数目、芳环数目、警报结构数目。

合成可行性分数 (synthetic accessibility score, SA SCORE) 计算方式为化合物的片段贡献减去复

杂度。片段贡献值根据 PubChem 数据库中上百万分子计算共性进行计算, 复杂度则考虑大环、非标准环、立体异构和相对分子质量大小等方面。

4 总结与展望

使用深度学习搭建生成模型进行分子生成研究, 在减少药物研发成本和时间方面有着显著优势。同时, 生成类人工智能模型仍在不断发展, 如最近在自然语言领域大火的 ChatGPT 和图像领域大火的扩散模型, 都为分子生成研究提供了新的方向。然而, 生成模型在药物设计领域的应用仍然处于初步阶段, 未被完全普及。虽然与传统的遗传算法相比, 基于深度学习的生成模型也可以获得可比较的结果, 但目前研究仍主要集中在生成模型本身, 而鲜有关于实际应用场景的研究。尽管已有很多关于使用生成模型生成新化合物的报告, 但对生成的化合物进行

实验评估的数据相对较少。因此, 未来的研究需要关注模型方法的发展与特定靶点的实际应用相结合。

深度学习作为从头药物设计的重要支柱, 有助于加速药物发现周期, 并为药物化学家带来新的创意。

[参考文献]

- [1] Anna G, Bellis L J, Patricia B A, *et al.* ChEMBL: a large-scale bioactivity database for drug discovery[J/OL]. *Nucleic Acids Res*, 2012, 40(Database issue): D1100–D1107[2023-03-20]. <https://pubmed.ncbi.nlm.nih.gov/21948594/>. DOI: 10.1093/nar/gkr777.
- [2] Bolton E E, Wang Y, Thiessen P A, *et al.* Chapter 12–PubChem: integrated platform of small molecules and biological activities[J/OL]. *Annu Rep Comput Chem*, 2008, 4: 217–241[2023-03-20]. [https://doi.org/10.1016/S1574-1400\(08\)00012-1](https://doi.org/10.1016/S1574-1400(08)00012-1).
- [3] Pence H E, Williams A. ChemSpider: an online chemical information resource[J]. *J Chem Educ*, 2010, 87(11): 1123–1124.
- [4] Jordan M I. Attractor dynamics and parallelism in connectionist sequential machine[M/OL]//Diederich J. *Artificial Neural Networks: Concept Learning*. IEEE Press, 1990: 112–127[2023-03-20]. <https://dl.acm.org/doi/abs/10.5555/104134.104148>.
- [5] Elman J L. Finding structure in time[J]. *Cognitive Sci*, 1990, 14(2): 179–211.
- [6] Schmidhuber J. Gradient flow in recurrent nets: the difficulty of learning long-term dependencies[M/OL]//Kolen J F, Kremer S C. *A Field Guide to Dynamical Recurrent Networks*. Wiley-IEEE Press, 2001: 237–243[2023-03-20]. <https://ieeexplore.ieee.org/document/5264952>.
- [7] Chung J, Gulcehre C, Cho K H, *et al.* Empirical evaluation of gated recurrent neural networks on sequence modeling[CP/OL]. (2014-12-11)[2023-03-20]. <https://doi.org/10.48550/arXiv.1412.3555>.
- [8] Bjerrum E J, Threlfall R. Molecular generation with recurrent neural networks (RNNs)[CP/OL]. (2017-05-11)[2023-03-20]. <https://doi.org/10.48550/arXiv.1705.04612>.
- [9] Grisoni F, Moret M, Lingwood R, *et al.* Bidirectional molecule generation with recurrent neural networks[J]. *J Chem Inf Model*, 2020, 60(3): 1175–1183.
- [10] Shi T, Huang S, Chen L, *et al.* A molecular generative model of ADAM10 inhibitors by using GRU-based deep neural network and transfer learning[J/OL]. *Chemom Intell Lab Syst*, 2020, 205: 104122[2023-03-20]. <https://doi.org/10.1016/j.chemolab.2020.104122>.
- [11] Neil D, Segler M, Guasch L, *et al.* Exploring deep recurrent models with reinforcement learning for molecule design[CP/OL]//6th International Conference on Learning Representations, April 30–May 3, 2018, Vancouver Convention Center, Vancouver, BC, Canada[2023-03-20]. <https://openreview.net/forum?id=Bk0xi1IDz>.
- [12] Goel M, Raghunathan S, Laghuvarapu S, *et al.* MoleGuLAR: molecule generation using reinforcement learning with alternating rewards[J]. *J Chem Inf Model*, 2021, 61(12): 5815–5826.
- [13] Olivecrona M, Blaschke T, Engkvist O, *et al.* Molecular *de-novo* design through deep reinforcement learning[CP/OL]. (2017-08-29)[2023-03-20]. <https://doi.org/10.48550/arXiv.1704.07555>.
- [14] Guo J, Fialková V, Arango J D, *et al.* Improving *de novo* molecular design with curriculum learning[J]. *Nat Mach Intell*, 2022, 4(6): 555–563.
- [15] Bengio Y. Learning deep architectures for AI[M/OL]. *Boston-Delft: Now Foundations and Trends*, 2009: 136[2023-03-20]. <https://ieeexplore.ieee.org/document/8187120>. DOI: 10.1561/22000000006.
- [16] Berenger F, Tsuda K. Molecular generation by fast assembly of (Deep) SMILES fragments[J]. *J Cheminform*, 2021, 13(1): 1–10.
- [17] Alperstein Z, Cherkasov A, Rolfe J T. All SMILES variational autoencoder[CP/OL]. (2019-01-03)[2023-03-20]. <https://doi.org/10.48550/arXiv.1905.13343>.
- [18] Polykovskiy D, Zhebrak A, Vetrov D, *et al.* Entangled conditional adversarial autoencoder for *de novo* drug discovery[J]. *Mol Pharmaceut*, 2018, 15(10): 4398–4405.
- [19] Zhavoronkov A, Ivanenkov Y A, Aliper A, *et al.* Deep learning enables rapid identification of potent DDR1 kinase inhibitors[J]. *Nat Biotechnol*, 2019, 37(9): 1038–1040.
- [20] Jin W G, Barzilay R, Jaakkola T. Junction tree variational autoencoder for molecular graph generation[CP/OL]. (2018-03-29)[2023-03-20]. <https://doi.org/10.48550/arXiv.1802.04364>.
- [21] Liu Q, Allamanis M, Brockschmidt M, *et al.* Constrained graph variational autoencoders for molecule design[CP/OL]. (2019-05-07)[2023-03-20]. <https://doi.org/10.48550/arXiv.1805.09076>.
- [22] Chenthamarakshan V, Das P, Hoffman S C, *et al.* CogMol: target-specific and selective drug design for COVID-19 using deep

- generative models[CP/OL]. (2020-06-24)[2023-03-20]. <https://arxiv.org/abs/2004.01215>.
- [23] Makhzani A, Shlens J, Jaitly N, *et al.* Adversarial autoencoders[CP/OL]. (2015-05-25)[2023-03-20]. <https://doi.org/10.48550/arXiv.1511.05644>.
- [24] Putin E, Asadulaev A, Vanhaelen Q, *et al.* Adversarial threshold neural computer for molecular *de novo* design[J]. *Mol Pharmaceut*, 2018, 15(10): 4386–4397.
- [25] Guimaraes G L, Sanchez-Lengeling B, Outeiral C, *et al.* Objective-reinforced generative adversarial networks (ORGAN) for sequence generation models[CP/OL]. (2018-02-07)[2023-03-20]. <https://doi.org/10.48550/arXiv.1705.10843>.
- [26] Cao N D, Kipf T. MolGAN: an implicit generative model for small molecular graphs[CP/OL]. (2022-09-27)[2023-03-20]. <https://doi.org/10.48550/arXiv.1805.11973>.
- [27] Vaswani A, Shazeer N, Parmar N, *et al.* Attention is all you need[CP/OL]. (2022-09-27)[2023-03-20]. <https://doi.org/10.48550/arXiv.1706.03762>.
- [28] Bagal V, Aggarwal R, Vinod P K, *et al.* MolGPT: molecular generation using a transformer-decoder model[J]. *J Chem Inf Model*, 2022, 62(9): 2064–2076.
- [29] Wang J, Hsieh C Y, Wang M, *et al.* Multi-constraint molecular generation based on conditional transformer, knowledge distillation and reinforcement learning[J/OL]. *Nat Mach Intell*, 2021, 3: 914–922[2023-03-20]. <https://doi.org/10.1038/s42256-021-00403-1>.
- [30] Zhu W H, Li Z Y, Cai L S, *et al.* Stepping back to SMILES transformers for fast molecular representation inference[CP/OL]. (2021-12-26)[2023-03-20]. <https://doi.org/10.48550/arXiv.2112.13305>.
- [31] Mitton J, Senn H M, Wynne K, *et al.* A graph VAE and graph transformer approach to generating molecular graphs[CP/OL]. (2021-04-09)[2023-03-20]. <https://doi.org/10.48550/arXiv.2104.04345>.
- [32] Kobyzev I, Prince S J D, Brubaker M A. Normalizing flows: an introduction and review of current methods[J]. *IEEE Trans Pattern Anal Mach Intell*, 2020, 43(11): 3964–3979.
- [33] Papamakarios G, Nalisnick E T, Rezende D J, *et al.* Normalizing flows for probabilistic modeling and inference[J]. *J Mach Learn Res*, 2021, 22(57): 1–64.
- [34] Shi C, Xu M, Zhu Z, *et al.* GraphAF: a flow-based autoregressive model for molecular graph generation[CP/OL]. (2020-02-27)[2023-03-20]. <https://doi.org/10.48550/arXiv.2001.09382>.
- [35] Luo Y, Yan K, Ji S. GraphDF: a discrete flow model for molecular graph generation[CP/OL]. (2021-06-02)[2023-03-20]. <https://doi.org/10.48550/arXiv.2102.01189>.
- [36] Kuznetsov M, Polykovskiy D. MolGrow: a graph normalizing flow for hierarchical molecular generation[CP/OL]. (2021-02-03)[2023-03-20]. <https://doi.org/10.48550/arXiv.2106.05856>.
- [37] Dietl T, Ohno H, Matsukura F, *et al.* Zener model description of ferromagnetism in Zinc-Blende magnetic semiconductors[J]. *Science*, 2000, 287(5455):1019–1022.
- [38] Ramakrishnan R, Dral P O, Rupp M, *et al.* Quantum chemistry structures and properties of 134 kilo molecules[J/OL]. *Sci Data*, 2014, 1: 140022[2023-03-20]. <https://pubmed.ncbi.nlm.nih.gov/25977779/>. DOI: 10.1038/sdata.2014.22.



【专家介绍】陈亚东: 中国药科大学教授, 博士生导师。江苏省“青蓝工程”优秀青年骨干教师, 江苏省“青蓝工程”中青年学术带头人, 美国密歇根大学医学院综合癌症中心访问学者。主持和参与了多项国家自然科学基金、国家重大科技专项“重大新药创制”等科研项目; 课题组与国内多个药企合作进行新药研发。申请国内专利 14 项, PCT 专利 3 项; 发表 SCI 论文 100 多篇。2015 年研究团队发现的 1.1 类抗肿瘤新药以 1.5 亿元人民币里程碑金转让给上海复星医药, 目前在美国、澳大利亚和中国大陆地区进行 I 期临床试验, 2019 年底获美国 FDA 授予孤儿药资格认定。