

深度学习辅助药物发现的研究进展

戴青青, 余俊霖, 李国菠*

(四川大学华西药学院药物化学系, 四川成都 610041)

[摘要] 深度学习技术近年来取得了重大突破, 被应用于医学、药学等多个领域。聚焦深度学习在创新药物发现中的发展和应用, 对深度学习被用于蛋白结构预测、药物靶标预测、药物-靶标相互作用预测、药物合成路线设计、从头药物分子设计以及药物吸收、分布、代谢、排泄和毒性 (ADMET) 预测等代表性案例进行详细综述, 同时总结了现有方法面临的问题和可能的解决思路, 以期对深度学习辅助药物发现相关方法的发展和應用提供借鉴与思考。

[关键词] 人工智能; 深度学习; 药物设计; 药物发现; 从头设计

[中图分类号] R914.2 **[文献标志码]** A **[文章编号]** 1001-5094 (2022) 01-0060-11

Recent Advances in Deep Learning Aided Drug Discovery

DAI Qingqing, YU Junlin, LI Guobo

(Department of Medicinal Chemistry, West China School of Pharmacy, Sichuan University, Chengdu 610041, China)

[Abstract] With its significant breakthroughs in recent years, deep learning technology has been used in medical, pharmaceutical and many other areas. This review focuses on the development and application of deep learning for innovative drug discovery, summarizes typical cases of deep learning for the prediction of protein structure, drug target and drug-target interaction, the design of drug synthesis route, *de novo* drug design, and the prediction of drug absorption, distribution, metabolism, excretion and toxicity (ADMET), and discusses the current problems and possible solutions, in the hope of providing some reference for the development and application of deep learning for drug discovery.

[Key words] artificial intelligence; deep learning; drug design; drug discovery; *de novo* design

人工智能 (artificial intelligence, AI) 概念始于 1956 年, 经过半个世纪的曲折探索, 于 2011 年进入蓬勃发展时期, 目前已成为一门新的技术科学, 推动人类进入智能时代。深度学习 (deep learning, DL), 又称为深度神经网络, 是 AI 领域中一个热门研究方向, 其通过对样本数据进行多层次的非线性信息处理和抽象, 挖掘内在规律, 用于解决特征学习、分类和模式识别等问题。当前主流的 DL 模型包括卷积神经网络 (convolutional neural network, CNN)、循环神经网络 (recurrent neural network, RNN) 和图神经网络 (graph neural network, GNN) 等, 以及这些模型的变体, 如残差卷积网络模型 (deep residual network, ResNet)、

变分自编码器 (variational autoencoder, VAE)、对抗自编码器 (adversarial autoencoder, AAE)、生成对抗网络模型 (generative adversarial network, GAN) 以及信息传递网络模型 (message passing neural network, MPNN) 等, 这些 DL 模型在图像识别、语音识别、机器翻译、人机对弈、无人驾驶等方面已取得了前所未有的成效, 深刻地改变着人们的生产生活方式^[1-2]。

同时, DL 技术在医学、药学、生命科学等领域也逐渐崭露头角。例如, 2018 年 Waller 团队通过 DL 网络对 1 240 万个单步反应进行化学转化规则提取, 再利用 3 种不同的神经网络与蒙特卡洛树搜索结合形成的新算法, 实现了化合物合成路线的高效设计^[3]。随后, Jensen 和 Jamison 团队又报道了一种集成合成路线设计和自动化合成的平台, 并完成了 15 个小分子药物的自动化合成, 进一步推动了该领域的发展^[4]。近期, Hassabis 团队报道了新蛋白结构预测工具 AlphaFold2, 通过将蛋白结构的物理和生物知识整合到 DL 方法中, 极大程度提高了蛋

接受日期: 2021-10-01

项目资助: 国家自然科学基金 (No.82122065); 四川省国际合作项目 (No.22GJHZ0253)

***通信作者:** 李国菠, 教授, 博士生导师;

研究方向: 药物设计与药物发现;

Tel: 028-85503235; **E-mail:** liguobo@scu.edu.cn

白结构预测的准确性^[5];与此同时,Baker团队也报道了新蛋白结构预测工具RoseTTAFold^[6],其采用了注意力机制使整个DL能够同时学习到蛋白一级/二级/三级结构不同维度的信息,预测准确率与AlphaFold2不相上下。此外,近几年还发展了若干DL方法用于药物-靶标相互作用预测、药物靶标预测、药物从头设计、药物性质[主要包括吸收、分布、代谢、排泄、毒性(ADMET)]的预测,从而服务于创新药物研发的多个重要环节。这些工具或将改变创新药物研发进程,提升药物研发效率。鉴于此,本文聚焦DL在创新药物发现中的发展和应用,综述具有代表性的DL案例和研究思路,总结其应用特点、面临的问题及可能的解决策略,期望为DL在药物发现领域的发展提供借鉴和思考。

1 基于深度学习的蛋白结构预测

蛋白质三维结构是药物靶标功能研究与药物设计的重要基础,如何快速高效获得准确的蛋白质结构是需要解决的科学问题。早期阶段,研究人员基于统计的蛋白质进化信息,并采用传统的机器学习方法(如蒙特卡罗方法、支持向量机等)和全连接神经网络(fully-connected neural network, FNN)模型实现蛋白质三维结构的预测。例如,Bohr等^[7]和Fariselli等^[8]使用目标蛋白一级序列、同源蛋白序列以及关联突变等数据来训练FNN模型,实现对蛋白质主链结构的预测,但距离实现蛋白质三维结构精准预测仍有较大差距。

随着蛋白结构数据的不断增加和DL技术的迅猛发展,更复杂的深度网络模型和更丰富的蛋白质序列信息被应用于预测蛋白质的三维结构,突破了从蛋白质一级序列直接得到蛋白质三维结构的瓶颈,预测精度接近实验解析水平。基于DL的蛋白结构预测是研究人员一直在尝试和努力的方向,大致流程是通过序列比对得到进化相关的多序列比对(multiple sequence alignment, MSA)特征,联合蛋白序列编码作为输入,利用深度网络模型预测残基间的接触图或更具体的距离分布,以及蛋白骨架的二面角分布,然后将预测的空间结构信息作为约束条件,重构出蛋白三维结构(见图1)。

例如,Hassabis团队最新报道的蛋白结构预测工具AlphaFold2,在最近的蛋白质结构预测技术评估(即The 14th Edition of Critical Assessment of Structure Prediction, CASP14)比赛中取得最佳预测名次,全局距离测试(global distance test, GDT)中位数得分达92.4,达到实验解析水平。AlphaFold2是基于注意力机制的神经网络模型,由Evoformer网络模块和结构生成模块组成,通过给定的一级序列,结合学习蛋白结构的物理和生物知识,端对端直接生成蛋白的三维结构。Baek等^[6]也基于注意力机制开发了一种新的端到端蛋白结构预测工具RoseTTAFold。该工具是一种三轨网络模型,分别用逐级连接的网络来传递和处理来自蛋白一级、二级、三级结构的信息,轨道之间的多次连接让网络能够同时学习序列、残基间距离和原子坐标之间的关系。实验结果表明,RoseTTAFold不仅预测精度接近AlphaFold2,为未知结构蛋白生物学功能和机制提供一种解释,而且还能直接根据序列信息快速构建出准确的蛋白-蛋白复合物结构。在所需计算资源和计算时间方面,RoseTTAFold较AlphaFold2也显示出一定的优势,除去序列比对和模版搜索所用时间,其仅需1个图形处理器(graphic processing unit, GPU)就能在10 min之内生成蛋白3D主链结构。另外,Rahman等^[9]在ResNet模型的基础上进行了改进,提出了一种用来预测蛋白质残基间距离的DL模型,相对比以上方法使用更少的蛋白特征,包括2种共同进化特征和3种非进化特征,实现对蛋白质残基间真实距离的高精度预测,与最先进的同类方法相比,局部距离差测试平均分数提高了10%以上,为蛋白质结构预测提供了一种新的参考。

此外,Yang等^[10]首次提出利用GAN模型预测蛋白质残基-残基接触图,并在基准测试集上表现出不错的预测效果。该模型被命名为GANcon,GANcon通过对抗性学习策略训练生成模型和判别模型,最终能够生成接近真实数据分布的接触图。其中,生成模型采用编码器-解码器框架从多种蛋白质序列特征中捕捉潜在的残基间接触信息,从而生成仿真的残基接触图;判别模型则选用基于残基块的CNN,以生成的或真实的接触图——蛋白质序列

特征样本作为输入, 识别生成的接触图与真实接触图之间的差异, 驱动生成模型生成更准确的接触图。他们还引入了一种新的对称焦点损失函数, 用来解

决接触图内数据不平衡问题。但 GANcon 在训练过程中的不稳定性以及输入特征的选择等方面仍有改进空间。

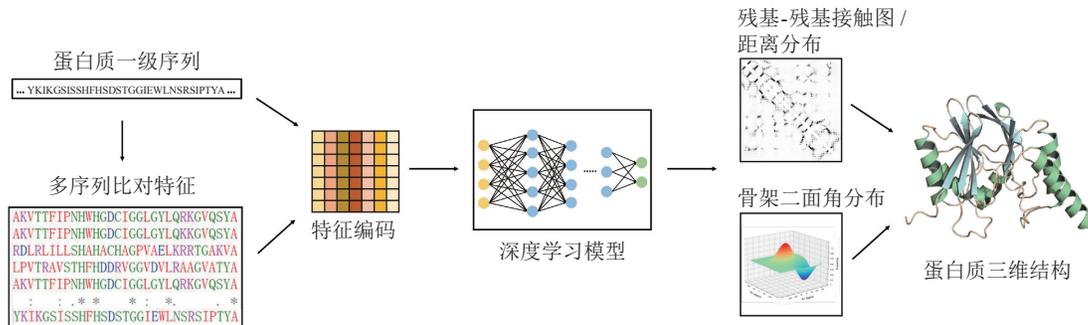


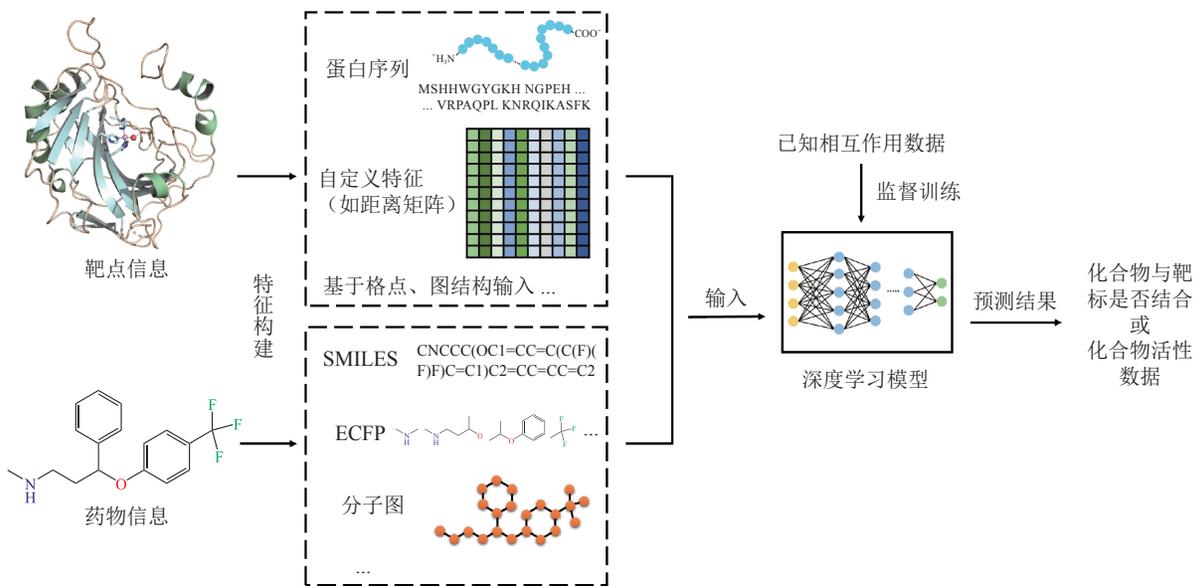
图 1 基于深度学习的蛋白质三维结构预测流程

Figure 1 The general process of deep learning-based 3D protein structure prediction

2 基于深度学习的药物-靶标相互作用预测

药物-靶标相互作用 (drug target interaction, DTI) 是药物发现的重要基础, 准确有效的 DTI 预测能极大地助力药物研发, 加速先导或苗头化合物发现。近几年, 基于 DL 预测 DTI 的方法陆续被报道,

其一般工作流程如图 2 所示, 研究人员针对药物和靶标的结构以及理化性质构建各具特色的描述符, 并采用不同的 DL 网络模型, 学习 DTI 规律, 最终预测出 DTI 的可能性或者相互作用强度。



SMILES: simplified molecular input line entry system (简化分子线性输入规范); ECFP: extended connectivity fingerprints (扩展连接指纹)

图 2 基于深度学习的药物靶标相互作用预测一般流程

Figure 2 The general process of deep learning-based drug-target interaction prediction

早期研究人员倾向于使用简单直接的输入数据和结构单一的网络框架。例如采用药物结构信息和靶标的序列信息, 通过基础版本的 RNN、CNN 等模型学习相互作用特征^[11-12], 但预测结果并不理想。

研究人员分析发现只是纯粹地使用药物-靶标相关信息套用 DL 模型不能从根本上解决问题, 需在 DL 和药物发现的双重理论指导下, 根据药物、靶标的各种性质合理构建输入描述符, 同时搭建适应药物-

靶标体系的神经网络框架,才能有效提高模型的预测能力和结果可靠性。在此基础上,发展出了一系列基于格点、基于图结构以及新算法的DL网络,并合理引入注意力机制等算法增强模型的可解释性。

基于格点的特征构建方法蕴含更加丰富的空间信息,比较适应于DTI预测体系。由此方法构建的特征可以视作一幅三维图片,可配合使用三维CNN模型进行训练、学习,但存在参数量大、计算成本高等问题。Li等^[13]借鉴ShuffleNet、Xception等轻量级三维CNN模型^[14]并构建了DeepAtom模型,用于预测药物-靶标亲和力。除了具备三维CNN模型的各种优势,DeepAtom模型同时通过深度可分离卷积解决了三维CNN模型参数过多的问题,并利用多个小的卷积核代替单个大卷积核,达到减少参数的同时增加网络复杂度的目的。该模型在PDBbind(2016版)核心测试集预测的皮尔森相关系数达0.831,表现出较强的预测能力。

Zheng等^[15]对DTI预测有着不同理解,他们将DTI预测抽象成虚拟问答(visual question answering, VQA)问题,采用药物SMILES和靶标残基距离矩阵作为输入,并基于CNN与RNN模型构建了DrugVQA模型,同时引入了注意力机制以增加模型的可理解性。经过训练及超参数优化,DrugVQA模型最终在数据库DUD-E上表现出不凡的预测能力,受试者工作特征曲线下面积(area under the receiver operating characteristic curve, ROC-AUC)达到0.972。

GNN模型在此领域也备受关注,Cho等^[16]采用了一种特殊的GNN模型,提出了InteractionNet框架,用于预测药物-靶标之间的结合常数。InteractionNet模型是一种非常规的GNN模型,在对药物-靶标体系建模时除了考虑共价键外,还考虑了非共价作用,最后基于PDBbind数据集采用20折交叉方法进行验证,其均方根误差(root mean square error, RMSE)为1.321,优于PoteintialNet模型(RMSE为1.343)。

Zeng等^[17]认为通过拼接药物和靶标的特征向量来表征二者的相互作用,并不能准确描述二者复杂作用体系,需要某种特殊的算法或网络来解决。

据此,他们提出了一种多注意力模块MATT_DTI,首先通过相对自注意模块提取药物的化合物原子间联系,用CNN模块分别学习药物和靶标的隐含信息,最后通过多头注意力模块和全连接层提取相互作用信息并给出预测结果。该方法在KIBA和Davis数据集上表现良好,均比同类模型有更好的预测效果,如用KIBA数据集进行测试,MATT_DTI模型平均标准误差(mean squared error, MSE)在0.15左右,低于其他基准模型的MSE指标。Sajadi等^[18]以药物指纹矩阵和药物-靶标矩阵为输入,构建了一个无监督去噪自编码器(denoising autoencoder, DAE)模型,并将其命名为AutoDTI⁺⁺。该方法在G蛋白偶联受体(G protein-coupled receptor, GPCR)数据集上预测随机药物靶点对时,ROC-AUC值达0.85,与类似算法的模型测试结果相比有明显提升。

3 基于深度学习的药物靶标预测

药物靶标预测可以帮助研究人员确定已知药物或活性分子的潜在靶标,从而有助于实现老药新用、药物重定位、毒性预测等。上述DTI预测方法也可以用于药物靶标预测。除此之外,基于异质网络等DL方法也被用于药物靶标预测,其特点在于利用药物-疾病信息、靶标-靶标信息、药物-靶标信息等多维度信息(见图3)作为网络输入特征,将其进一步转化为一组DL模型可处理的特征矩阵,实现对药物靶标的预测。

自编码器(autoencoder, AE)及其变体,如DAE等在基于异质网络的靶标预测方法中较为主流,研究人员通过收集药物、靶标相关的各种信息,构建异质网络,利用各种AE变体进行学习,最终分析和预测药物的潜在靶标。Zeng等^[19]收集了药物-疾病、药物-不良反应、药物-靶标、药物-药物相关信息,以此构建异质网络,从中提取药物与靶标之间的关系,使用随机游走算法计算得到概率共生矩阵(probabilistic co-occurrence matrix, PCO),再计算正点互信息矩阵(positive pointwise mutual information, PPMI)来表征异质网络整体结构,用于训练DL网络模型,由此发展了deepDR模型。该

模型在基准模型上, deepDR 预测效果更佳, ROC-AUC 达 0.908。后来, 他们又进一步做出了改进^[20], 设计了一个新的模型 (deepDTnet), 该模型在输入和框架方面都进行了优化, 丰富了异质网络所蕴含的信息, 加入了更多靶标相关信息, 如靶标-靶标相似性、靶标-疾病信息, 同时保留 PCO 矩阵和 PPMI 矩阵的表征方式, 采用多层 DAE 学习异质网络的隐含信息。与 deepDR 相比, deepDTnet 具有更强的预

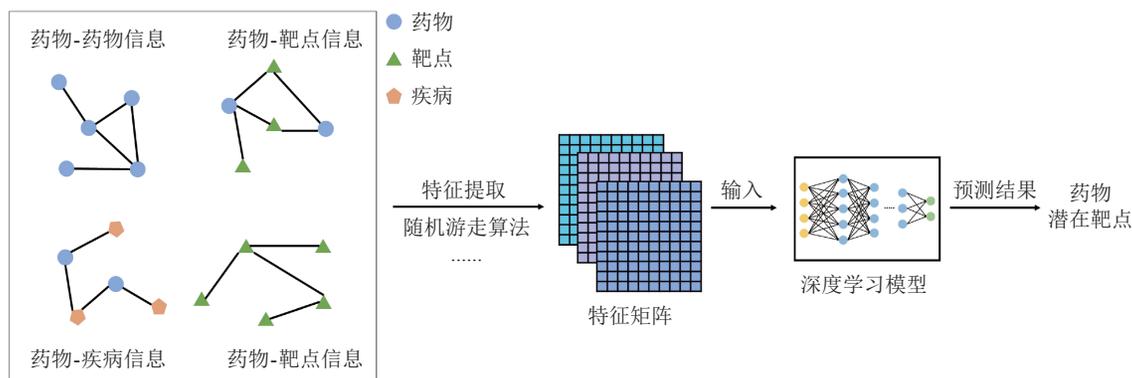


图 3 基于机器学习——异质网络的药物靶标预测方法一般流程

Figure 3 The general process of machine learning (heterogeneous network)-based target prediction

除了 AE 及其变体外, 其他模型在药物靶标预测方面也展现出不俗的预测效果。Manoochchri 等^[22] 利用更简单的输入 (仅考虑药物-药物相似性和靶标-靶标相似性信息) 和 FNN 模型进行学习预测, 但将更多的精力放在输入数据的处理上, 提出了独特的特征提取和构建方法。他们利用异质网络的拓扑结构来预测药物的未知靶标, 通过药物-药物相似性和靶标-靶标相似性信息把药物-靶标异质网络抽象成半二部图, 并从中提取出多个封闭子图, 然后采用 Weisfeiler-Lehman 算法对每个子图中的节点进行排序标记, 以表征药物-靶标对的拓扑结构。最后使用这种特殊的输入来训练 FNN 模型, 同时进行了 10 折交叉验证。结果显示, 该方法比 BLMNII、CMF、HNM 等同类模型预测能力更强。此外, GNN 模型也被用来处理这些异质网络, 进行药物靶标的预测。Huang 等^[23] 提出了 SkipGNN 模型, 并认为异质网络中直接相连的 2 个节点并不一定有很强的相似性, 反而是间接的或跳跃的节点间的相似性可能更加必要。根据这种思想, 他们以药物-药物、靶标-靶标、药物-靶标、基因-疾病相关

测能力, ROC-AUC 达 0.963。也有研究人员通过将 AE 和其他网络模型结合, 尝试发展了新的网络模型。如 Peng 等^[21] 提出了 DTI-CNN 模型, 特点在于使用 Jaccard 相似性系数结合重启随机游走算法 (random walk with restart, RWR) 来提取药物特征和靶标特征, 且经过 DAE 层后添加了 CNN 模块来预测最终结果, 训练后 ROC-AUC 达 0.9416, 与 deepDTnet 效果相当。

信息构建了异质网络, 从中提取跳跃相似性信息并构建跳跃相互作用图, 同时结合原始图输入至 GNN 模型中, 最后经由解码器输出药物与靶标相互作用概率。实验结果表明 SkipGNN 模型优于其他模型, 如 DeepWalk、图卷积神经网络 (graph convolutional neural network, GCN) 和 node2vec 模型等。

4 基于深度学习的合成路线设计

药物研发离不开合成路线设计, 设计高效的合成路线可大幅度降低药物研发成本、缩短生产周期、提高药物研发效率。传统的计算机辅助合成路线设计的方法主要是基于大量“专家”规则和逆合成分析方法来规划合成路线, 但其存在设计速度较慢、设计的合成路线往往不太合理等问题^[24]。随着 DL 算法在化合物性质预测和生物活性预测等领域中展现出巨大的潜力, 其也逐渐被应用于合成路线的设计并取得了一定的进展。

Waller 团队于 2018 年报道了一种 AI 工具 3N-MCTS, 通过使用 3 种不同的深度神经网络 (分别是拓展策略网络、筛选网络和展示策略网络)

和蒙特卡罗树搜索算法来设计目标化合物的合成路线^[3]。他们首先利用拓展策略网络对目标分子进行逆向化学转换, 搜索当前节点可能的变换路径, 然后使用筛选网络分析判断反应是否可行, 过滤不合理的反应路线, 最后通过展示策略网络多次随机采样对搜索节点进行评价打分。研究人员利用来自 Reaxys 数据库的 1 240 万条反应数据训练这些网络, 学习化学转化规则。与其他方法相比, 3N-MCTS 在合成路线的搜索速度、质量等方面均有显著提升, 能在短时间内生成数百个化合物的合成路线, 且双盲实验结果表明 3N-MCTS 预测分子合成路线水平接近合成化学家水平。这种方法的优势体现在无需专家自定义规则, DL 模型就可以学习到已知反应所蕴含的转化规则, 然后根据学习到的规则快速选择出最佳合成路线。

随后, Coley 等^[4]推出了一个基于 AI 的自动化合成平台, 首先利用前馈神经网络生成目标分子

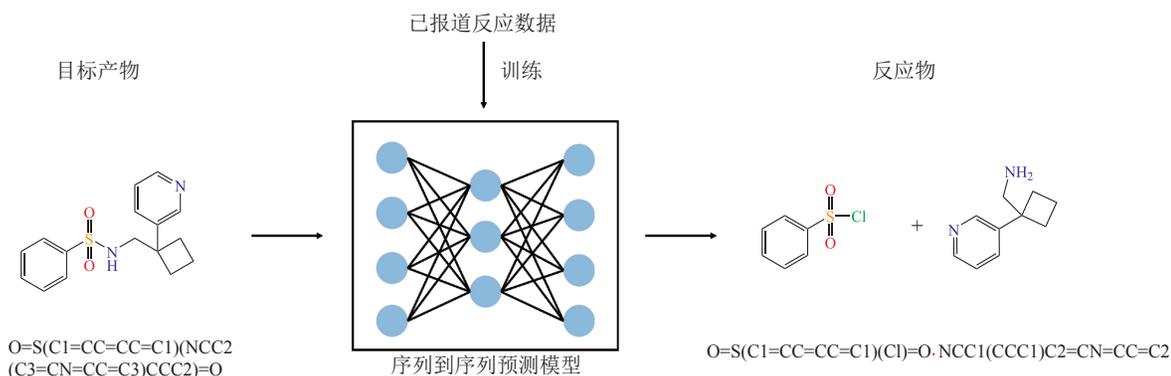


图 4 基于序列到序列模型进行合成路线预测

Figure 4 Prediction of synthetic route based on the seq2seq model

Liu 等^[25]率先将 seq2seq 模型应用到逆合成预测任务中, 使用的 seq2seq 模型是基于 RNN 的编码器-解码器结构, 并在包含 5 万个专利反应的数据集上训练, 并初步达到了与基于规则的基准方法效果相当的水平。该方法在一定程度上突破了专家规则的限制, 并表现出良好可扩展性的优势。随后 seq2seq 模型经过发展, 得到了较为流行基于注意力机制的 Transformer 模型。Zheng 等^[26]开发了一种无模板的自校正逆合成路线预测工具 SCROP, 通过使用基于多头注意力机制的 Transformer 网络模型预测逆合成路线, 同时引入了基于 Transformer 的语法校正器,

的合成路线, 然后机器人根据合成方案执行一系列具体的制备过程, 实现自动化合成。研究人员使用 Reaxys 和 USPTO 数据库中的反应数据训练网络模型, 学习反应转换规则, 为目标化合物设计出可行的合成路线, 包括给出反应条件, 同时根据合成路线中的反应类型是否容易实现以及中间产物是否多样化等条件进一步筛选得到最优合成路线。最后, 他们通过该平台成功完成了 15 种小分子药物合成路线设计并实现了自动化合成。同时, 基于 DL 的序列到序列 (sequence-to-sequence, seq2seq) 模型 (如 Transformer 模型等) 的发展给不依赖模板的逆合成预测任务提供了一种新的解决思路 (见图 4): 可将该任务看成自然语言处理 (natural language processing, NLP) 领域内机器翻译任务, 输入目标分子的 SMILES 序列, 不依赖反应规则, 就能输出对应单步的反应物 SMILES 序列。

对预测模型产生的不合理候选前体分子 SMILES 进行修正。SCROP 在基准数据集上预测准确率达 59%, 比基于模板的方法提高了 6%; 同时实验结果表明语法校正器的加入提高了模型预测质量, 使无效的候选前体分子比例从 12.1% 降至 0.7%。此外, Guo 等^[27]结合 Transformer 模型和贝叶斯推理算法进行逆向合成预测。他们将该任务视为组合优化问题, 即在所有可用的反应物组合中找到一组最佳的反应物对, 用来合成目标产物。他们首先通过训练好的 Molecular Transformer 模型对给定反应物组合进行高精度正向预测, 然后基于贝叶斯定理将正向

预测模型反演为逆向合成模型,同时使用蒙特卡罗搜索算法探索得到最佳的反应物组合。正向和逆向预测模型的组合提高了合成路线的可行性,同时改善了逆合成问题的不适应性。

这类序列模型一般利用分子的 SMILES 字符串作为输入,未能有效刻画出分子中各原子间复杂关系。为此,Shi 等^[28]提出了一种基于图神经网络的无模版逆合成预测框架 G2G (graph to graph framework),利用图表征分子,将任务转化为图到图的翻译问题,即将目标分子图转化为一组反应物分子图。研究人员首先基于 GCN 识别目标分子的反应中心,将目标分子拆分为一组合成子。然后,通过图 VAE 将每个合成子转换为最终的反应物分子图。实验结果表明 G2G 在 Top-1 准确率指标上明显优于其他无模版的基准模型(如 seq2seq 模型、transformer 模型等),并与最先进的基于模板的方法相比水平相当,如条件图逻辑网络(conditional graph logic network, GLN)模型。

5 基于深度学习的从头药物分子设计

近年来在从头药物分子设计领域,DL 方法因部分解决了传统方法的组合爆炸、多目标优化等问题而受到越来越多的关注。许多相关研究都证明了 DL 方法在从头药物分子设计的可行性,目前关于 DL 在这方面的应用已经被总结报道^[29-31],在此笔者将对最新的研究进展进行进一步介绍。Born 等^[32]构建了一种混合的 VAE 模型,用来生成具有抗癌药物特性的候选分子。值得注意的是,他们不仅使用分子 SMILES 作为输入,还首次加入疾病相关的基因表达数据,同时使用抗癌药物敏感性预测模型作为奖励函数。混合的 VAE 模型由 2 个并列的 VAE 组成,一个用于接收小分子 SMILES 以学习其语法规则,另一个 VAE 用于接收基因表达数据以学习其特征表示,然后将这 2 个 VAE 编码器的输出结果一并输入到同一解码器,生成新分子,最后用抗癌药物敏感性预测模型预测生成分子对靶细胞的活性值。应用在 4 种不同癌症类型的实例表明,该模型能够针对特定疾病生成具有较强抑制效果的分子,且生成的分子在结构、可合成性以及溶解性等方面均与

现有药物相似。然而,VAE 也存在一定局限,它只会最大限度地“模仿”训练数据,尽可能生成与训练数据在结构上相似的分子,因此生成分子的结构新颖性较低。

AAE 在 VAE 基础上增加了判别模型,对采样分子和真实样本进行区分,基于对抗的思想训练生成模型和判别模型,扩展了分子的生成空间,一定程度上弥补了 VAE 在生成分子时结构新颖性方面的缺陷。Polykovskiy 等^[33]构建了一种新的 AAE 模型,即条件 AAE,其能够基于指定条件(如药物分子的靶标特异性、溶解性、可合成性等)生成相应的分子。其中,基于长短时记忆网络(long short term memory, LSTM)分别构建编码器和解码器,同时使用多层的 FNN 作为判别模型,用来判断采样分子是否符合真实数据分布以及是否具备所需的理化性质,并基于半监督学习方法优化模型。他们利用该模型成功发现了一种新型的 Janus 激酶 3 (Janus kinase 3, JAK3) 抑制剂。

Bagal 等^[34]受生成式预训练新型神经网络模型(generative pre-training transformer, GPT) Transformer 在生成文本任务中取得突破性进展的启发,基于 GPT 构建了一个新的生成模型 MolGPT,能够根据给定条件(输入 SMILES 字符串、脂水分配系数、可合成性分数以及拓扑极性表面积等目标属性值)生成具有所需骨架和理想特性的分子。MolGPT 由多个堆叠的解码器模块组成,每个解码器包含一层掩码自注意力层和多层全连接网络,能够捕获 SMILES 字符串中字符间远距离依赖关系。与 VAE、AAE 等其他 DL 模型相比,MolGPT 在生成分子的有效性、独特性以及新颖性方面表现较好,打分为 0.981、0.998 和 1.0。

Goel 等^[35]结合 RNN 和强化学习,提出了一个分子生成模型 MoleGuLAR,其能够对分子的类药性、结合亲和力等方面进行多目标优化。尤其是,他们提出一种新的交替奖励策略,奖励函数随着生成不同分子的过程中在动态地改变,使得模型能够交替探索不同的化学区间,采样得到更加合理的分子。区别于以往大多数 DL 模型只能生成一维或二维的分子,Li 等^[36]将 DL 与基于结构的从头药物

设计策略相结合, 发展了一种新的从头分子生成模型 DeepLigBuilder, 其能够直接生成具有高结合亲和力和力类药分子的三维结构。DeepLigBuilder 首先利用一种图生成模型即配体神经网络 (ligand neural network, L-Net) 实现生成类药分子的三维结构, 然后结合蒙特卡洛树搜索方法将靶标的结构信息引入到模型中, 在靶标活性位点搜索、优化分子的结合构象, 从而得到具有高结合亲和力的新分子。通过将其应用于严重急性呼吸综合征冠状病毒 2 (severe acute respiratory syndrome coronavirus 2, SARS-CoV-2) 抑制剂的从头设计, 他们得到了 3 种新型具有高预测结合亲和力且与已知抑制剂结构类似的 SARS-CoV-2 潜在抑制剂, 证明了 DeepLigBuilder 在从头药物设计以及先导物优化方面的实用性。

为了解决 DL 在小规模训练数据集上表现较差等问题, Krishnan 等^[37]设计了一套基于 RNN 的生成模型和迁移学习的药物从头设计流程, 生成的分子不仅具有所需类药特性, 同时还具有靶标特异性。他们首先利用 ChEMBL 数据库中的活性分子 SMILES 数据预先训练 RNN 生成模型, 以学习 SMILES 语法规则; 然后, 通过对接得到具有靶标选择性的分子并进行迁移学习, 生成作用于特定靶标的分子; 同时, 再建立另一个基于 RNN 的预测模型, 作为奖励函数评价生成的分子与靶标的结合亲和力。另外, Moret 等^[38]将 RNN 生成模型与数据增强、温度采样和迁移学习这 3 种优化方法结合起来, 也能够具有少量数据情况下生成所需特性的新分子。

6 基于深度学习的药物吸收、分布、代谢、排泄和毒性预测

药物的 ADMET 性质研究对于药物研发也是至关重要的。据统计, 将近 50% 的候选药物在临床试验阶段因 ADMET 性质不符合要求而宣告失败。因此, 在早期药物发现和药物设计阶段, 研究人员应提前对药物分子的 ADMET 性质进行预测评估, 以降低后续临床试验失败的风险。相较于耗时耗力的实验方法, 精确可靠的 ADMET 预测方法能极大地缩短时间花费、减少实验成本, 提高候选药物的筛选效率, 基于 DL 的 ADMET 预测方法则恰逢其会, 并逐渐

成为预测药物 ADMET 性质的重要手段。

近几年来, 利用 DL 方法来预测小分子性质已经较为普遍, 其中基于 GNN 模型的方法受到了学界的广泛认可, 预测结果相较其他 DL 方法更为可靠。2018 年, Wu 等^[39]基于 DeepChem 平台构建了一个用于分子性质预测的 DL 框架, 称为 MoleculeNet。他们通过这个框架为同行提供了一个基准, 可以用于比较各种不同模型的效果和可靠程度。该框架涵盖了不同的数据集拆分方法, 包括基于骨架、随机拆分等; 以及不同的特征构建方法, 处理为 ECFP、图结构等; 和不同的网络模型, 例如 GCN、MPNN、weave、随机森林 (random forest, RF)、核岭回归 (Kernel ridge regression, KRR) 等; 并针对各种 ADMET 性质相关的数据库 (如 QM8、Clintox、Lipophilicity、BBBP 等) 进行训练和测试。通过一系列基准测试, 他们发现在应用量子力学性质、物理化学性质、生理学性质相关的数据集时, 最佳的 GNN 模型比最佳的传统模型更为有效, 如应用 QM8 数据集训练模型并预测小分子量量子力学性质时, 以平均绝对误差 (mean absolute error, MAE) 为评价指标, 表现最佳的传统模型是 KRR 模型, 该模型 MAE 达 0.015, 而基于 GNN 的网络模型中表现最佳的是 MPNN 模型, 其测试结果 MAE 为 0.014 3, 误差低于 KRR 模型测试结果。随后研究人员从不同角度出发, 建立了一系列各具特色的 GNN 模型。Feinberg 等^[40]构建了一种新型 GNN 网络模型 PotentialNet, 其核心思想是在更新原子状态过程中考虑距离因素, 比常用的邻接矩阵更能描述药物分子结构。该方法相较于传统的机器学习方法和一些常见的 GNN 模型性能更佳, 仍以 QM8 数据集进行测试, 在基于此数据集预测小分子量量子化学性质任务中, MPNN 在测试集上 MAE 达 0.013 9, 而 PotentialNet 则提升明显, MAE 在 0.011 8 左右。后续研究中, 他们又进一步在 PotentialNet 模型基础上进行了改进, 设计出多任务 PotentialNet 模型, 同时采用 31 项 ADMET 性质进行训练, 最终同时预测这 31 项性质^[41], 例如电压门控钾离子通道 (human ether-à-go-go-related gene encoded potassium ion channel, hERG) 抑制性、人肝细胞清除率、半

衰期、脂溶性等, 并与 RF 模型进行了比较。对于绝大部分性质而言, 多任务 PotentialNet 模型预测的相关系数 (R^2) 与 RF 模型相比都有不同程度的提高, 例如以时序拆分方法拆分数据集时, 多任务 PotentialNet 模型较 RF 模型, 在 31 项性质预测中 R^2 平均高出 64%。

Yang 等^[42] 则开发了一种有向信息传递网络 (directed message passing neural network, D-MPNN), 与往常的 GNN 模型做法不同, 在表征药物分子结构时, 他们将原子间的键考虑为有方向的边, 而非常规的无向的边, 且通过边的方向来对原子的状态进行更新, 减少了无效冗余的原子状态更新。预测结果表明, 在所有数据集上 D-MPNN 都比 RF 模型、FNN 模型等性能更好或者相当; 例如, 在血脑屏障透过能力预测方面, D-MPNN 模型的 ROC-AUC 高达 0.925, 而 RF 模型和 FNN 模型分别仅为 0.788 和 0.899。Li 等^[43] 提出了基于多头三联注意力机制的 MPNN 模型 TrimNet, 通过给定的邻接矩阵、边特征矩阵、节点特征矩阵, 分析周围原子对当前原子的影响, 从而实现高效地从图结构表征的药物分子结构中学习潜在信息, 并大幅度减少模型参数数量、降低计算成本, 最终在多个数据集上取得良好的预测结果, 如在 ClinTox 数据集上 ROC-AUC 高达 0.948。

除了 GNN 相关模型, 研究人员也尝试了其他类型的 DL 模型, 并获得一定成果。Kim 等^[44] 开发了首个基于自注意力机制具有可解释性的 DNN 模型, 用于预测药物是否存在 hERG 毒性。尽管只是采用了较为简单的 ECFP 描述符和 FNN 网络模型, 但在测试集上 ROC-AUC 依旧高达 0.893, 较传统的定量构效关系 (quantitative structure-activity relationship, QSAR) 模型, 有明显的改善。Wang 等^[45] 基于概念新颖的胶囊网络模型 (capsule neural network, CapsNet), 并结合 CNN、受限波尔兹曼机 (restricted boltzmann machine, RBM) 等网络模型构建了一系列衍生网络, 用于预测药物 hERG 毒性, 训练得到的最佳模型 ROC-AUC 达 0.944。也有研究团队通过 DL 模型直接学习实验数据并预测给药后患者体内药物的药效学 (pharmacodynamics,

PD)/药动学 (pharmacokinetics, PK) 性质变化曲线。例如, 最近 Lu 等^[46] 基于 RNN 模型和神经微分方程 (Neural-ODE) 提出了 Neural-PK/PD 模型, 其创新之处在于设计网络框架时, 保留了 PK/PD 的一些基本原理, 如药物的体内效应与给药剂量、体内浓度直接相关等, 从而提升了 PK/PD 性质的预测准确度。

7 结语与展望

DL 技术在药物发现多个环节中取得了惊人的预测能力, 正在改变着药物研发进程, 将有可能降低药物发现成本、提高药物研发效率。然而, 现有 DL 技术仍面临着诸多挑战。首先, 大多数 DL 技术严重依赖大量的计算资源, 一定程度上限制了 DL 方法的发展及应用。如何在保持模型预测准确率的前提下, 降低 DL 模型对计算资源的依赖已成为 DL 领域的一个研究热点^[47]。其中一个主流思路是通过修剪 DL 模型或者改善 DL 模型结构以减少网络参数数量和运算量, 从而降低对计算资源的需求。目前已有一些新型的轻量级 DL 模型被开发和应用^[14], 如 SqueezeNet、ThiNet、ShuffleNet。其次, 数据样本量、来源、质量等参差不齐, 也限制了 DL 技术建立和优化。DL 模型的训练依赖于大规模且高质量的数据样本。如何有效进行小样本学习是未来 DL 重要的发展方向^[48], 目前已有一些针对小样本学习的方法, 如采用数据增强技术、迁移学习、多任务学习策略等。同时, 数据集的质量也决定着 DL 模型预测性能的好坏。药物研发相关原始数据的提取、特征构建等方法尚存在不足, 影响着高质量 DL 模型的发展。近年来, 图神经网络的发展, 蕴含更多结构信息的图被逐渐用来表征分子并应用于药物发现领域, 已取得一些研究进展。此外, DL 模型中超参数搜索、内部机制的不可解释性等, 也一定程度上阻碍了该技术的发展。总而言之, 以上 DL 技术面临的种种不足和挑战都在提示我们, 需要更多不同背景的研究人员加入到这一领域, 来提出更多精湛的 DL 算法, 并且要充分结合传统的药物设计方法, 才能逐步解决药物研发过程中各个环节的具体问题, 从而能助力创新药物发现, 进一步推动药物研发领域迈向智能时代。

【参考文献】

- [1] Sengupta S, Basak S, Saikia P, *et al.* A review of deep learning with special emphasis on architectures, applications and recent trends[J/OL]. *Knowl-Based Syst*, 2020, 194: 105596[2021-10-01]. <https://doi.org/10.1016/j.knosys.2020.105596>.
- [2] Fawaz H I, Forestier G, Weber J, *et al.* Deep learning for time series classification: a review[J]. *Data Min Knowl Discov*, 2019, 33(4): 917-963.
- [3] Segler M H S, Preuss M, Waller M P. Planning chemical syntheses with deep neural networks and symbolic AI[J]. *Nature*, 2018, 555(7698): 604-610.
- [4] Coley C W, Thomas D A, Lummiss J A M, *et al.* A robotic platform for flow synthesis of organic compounds informed by AI planning[J/OL]. *Science*, 2019, 365(6453): eaax1566[2021-10-01]. <https://doi/10.1126/science.aax1566>.
- [5] Jumper J, Evans R, Pritzel A, *et al.* Highly accurate protein structure prediction with AlphaFold[J]. *Nature*, 2021, 596(7873): 583-589.
- [6] Baek M, Dimaio F, Anishchenko I, *et al.* Accurate prediction of protein structures and interactions using a 3-track network[J]. *Science*, 2021, 373(6557): 871-876.
- [7] Bohr H, Bohr J, Brunak S, *et al.* A novel approach to prediction of the 3-dimensional structures of protein backbones by neural networks[J]. *FEBS Lett*, 1990, 261(1): 43-46.
- [8] Fariselli P, Olmea O, Valencia A, *et al.* Prediction of contact maps with neural networks and correlated mutations[J]. *Protein Eng*, 2001, 14(11): 835-843.
- [9] Rahman J, Newton M A H, Islam M K B, *et al.* Enhancing protein inter-residue real distance prediction by scrutinising deep learning models[J/OL]. *Sci Rep*, 2022, 12(1): 787[2021-10-01]. <https://doi/10.1038/s41598-021-04441-y>.
- [10] Yang H, Wang M H, Yu Z H, *et al.* GANcon: Protein contact map prediction with deep generative adversarial network[J/OL]. *IEEE Access*, 2020, 8: 80899-80907[2021-10-01]. <https://doi/10.1109/ACCESS.2020.2991605>.
- [11] Öztürk H, Özgür A, Ozkirimli E. DeepDTA: deep drug-target binding affinity prediction[J/OL]. *Bioinformatics*, 2018, 34(17): i821-i829[2021-10-01]. <https://doi/10.1093/bioinformatics/bty593>.
- [12] Karimi M, Wu D, Wang Z Y, *et al.* DeepAffinity: interpretable deep learning of compound protein affinity through unified recurrent and convolutional neural networks[J]. *Bioinformatics*, 2019, 35(18): 3329-3338.
- [13] Li Y J, Rezaei M A, Li C L, *et al.* DeepAtom: a framework for protein-ligand binding affinity prediction[J/OL]. 2019[2021-10-01]. <https://doi/10.1109/BIBM47256.2019.8982964>.
- [14] Cheng J, Wang P S, Li G, *et al.* Recent advances in efficient computation of deep convolutional neural networks[J]. *Front Inform Technol Electro*, 2018, 19(1): 64-77.
- [15] Zheng S, Li Y, Chen S, *et al.* Predicting drug-protein interaction using quasi-visual question answering system[J]. *Nat Mach Intell*, 2020, 2: 134-140[2021-10-01]. <https://doi.org/10.1038/s42256-020-0152-y>.
- [16] Cho H, Lee E K, Choi I S. Layer-wise relevance propagation of InteractionNet explains protein-ligand interactions at the atom level[J/OL]. *Sci Rep*, 2020, 10(1): 21155[2021-10-01]. <https://doi/10.1038/s41598-020-78169-6>.
- [17] Zeng Y N, Chen X R, Luo Y J, *et al.* Deep drug-target binding affinity prediction with multiple attention blocks[J]. *Brief Bioinform*, 2021, 22(5): bbab117[2021-10-01]. <https://doi/10.1093/bib/bbab117>.
- [18] Sajadi S Z, Zare Chahooki M A, Gharaghani S, *et al.* AutoDTI++: deep unsupervised learning for DTI prediction by autoencoders[J]. *BMC Bioinformatics*, 2021, 22(1): 204[2021-10-01]. <https://doi/10.1186/s12859-021-04127-2>.
- [19] Zeng X X, Zhu S Y, Liu X R, *et al.* deepDR: a network-based deep learning approach to *in silico* drug repositioning[J]. *Bioinformatics*, 2019, 35(24):5191-5198.
- [20] Zeng X X, Zhu S Y, Lu W Q, *et al.* Target identification among known drugs by deep learning from heterogeneous network[J]. *Chem Sci*, 2020, 11(7): 1775-1797.
- [21] Peng J J, Li J Y, Shang X Q. A learning-based method for drug-target interaction prediction based on feature representation learning and deep neural network[J/OL]. *BMC Bioinformatics*, 2020, 21(13): 394[2021-10-01]. <https://doi/10.1186/s12859-020-03677-1>.
- [22] Manoochehri H E, Nourani M. Drug-target interaction prediction using semi-bipartite graph model and deep learning[J/OL]. *BMC Bioinformatics*, 2020, 21(4): 248[2021-10-01]. <https://doi/10.1186/s12859-020-3518-6>.
- [23] Huang K X, Xiao C, Glass L M, *et al.* SkipGNN: predicting molecular interactions with skip-graph networks[J/OL]. *Sci Rep*, 2020, 10(1): 21092[2021-10-01]. <https://doi/10.1038/s41598-020-77766-9>.
- [24] Liu Y D, Yang Q, Li Y, *et al.* Application of machine learning in organic chemistry[J]. *Chin J Org Chem*, 2020, 40(11): 3812-3827.
- [25] Liu B, Ramsundar B, Kawthekar P, *et al.* Retrosynthetic reaction prediction using neural sequence-to-sequence models[J]. *ACS Cent Sci*, 2017, 3(10): 1103-1113.
- [26] Zheng S J, Rao J H, Zhang Z Y, *et al.* Predicting retrosynthetic reactions using self-corrected transformer neural networks[J]. *J Chem Inf Model*, 2020, 60(1): 47-55.

- [27] Guo Z L, Wu S, Ohno M, *et al.* Bayesian algorithm for retrosynthesis[J]. *J Chem Inf Model*, 2020, 60(10): 4474–4486.
- [28] Shi C, Xu M, Guo H, *et al.* A graph to graphs framework for retrosynthesis prediction[EB/OL]. (2020-03-28)[2021-10-01]. <https://arxiv.org/abs/2003.12725v1>.
- [29] 梁礼, 邓成龙, 张艳敏, 等. 人工智能在药物发现中的应用与挑战[J]. *药学进展*, 2020, 44(1): 18–27.
- [30] Xu Y J, Lin K J, Wang S W, *et al.* Deep learning for molecular generation[J]. *Fut Med Chem*, 2019, 11(6): 567–597.
- [31] Wang M Y, Wang Z, Sun H Y, *et al.* Deep learning approaches for de novo drug design: an overview[J]. *Curr Opin Struct Biol*, 2021, 72: 135–144[2021-10-01]. <https://doi/10.1016/j.sbi.2021.10.001>.
- [32] Born J, Manica M, Oskooei A, *et al.* PacMannRL: *de novo* generation of hit-like anticancer molecules from transcriptomic data via reinforcement learning[J]. *iScience*, 2021, 24(4): 102269[2021-10-01]. <https://doi/10.1016/j.isci.2021.102269>.
- [33] Polykovskiy D, Zhebrak A, Vetrov D, *et al.* Entangled conditional adversarial autoencoder for *de novo* drug discovery[J]. *Mol Pharm*, 2018, 15(10): 4398–4405.
- [34] Bagal V, Aggarwal R, Vinod P K, *et al.* MolGPT: molecular generation using a transformer-decoder model[J]. *J Chem Inf Model*, 2021[2021-10-01]. <https://doi/10.1021/acs.jcim.1c00600>.
- [35] 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.
- [36] Li Y B, Pei J F, Lai L H. Structure-based *de novo* drug design using 3D deep generative models[J]. *Chem Sci*, 2021, 12(41): 13664–13675.
- [37] Krishnan S R, Bung N, Bulusu G, *et al.* Accelerating *de novo* drug design against novel proteins using deep learning[J]. *J Chem Inf Model*, 2021, 61(2): 621–630.
- [38] Moret M, Friedrich L, Grisoni F, *et al.* Generative molecular design in low data regimes[J]. *Nat Mach Intell*, 2020, 2(3): 171–180.
- [39] Wu Z Q, Ramsundar B, Feinberg E N, *et al.* MoleculeNet: a benchmark for molecular machine learning[J]. *Chem Sci*, 2018, 9(2): 513–530.
- [40] Feinberg E N, Sur D, Wu Z Q, *et al.* PotentialNet for molecular property prediction[J]. *ACS Cent Sci*, 2018, 4(11): 1520–1530.
- [41] Feinberg E N, Joshi E, Pande V S, *et al.* Improvement in ADMET prediction with multitask deep featurization[J]. *J Med Chem*, 2020, 63(16): 8835–8848.
- [42] Yang K, Swanson K, Jin W G, *et al.* Analyzing learned molecular representations for property prediction[J]. *J Chem Inf Model*, 2019, 59(8): 3370–3388.
- [43] Li P Y, Li Y Q, Hsieh C Y, *et al.* TrimNet: learning molecular representation from triplet messages for biomedicine[J]. *Brief Bioinform*, 2021, 22(4): bbaa266[2021-10-01]. <https://doi/10.1093/bib/bbaa266>.
- [44] Kim H, Nam H. hERG-Att: Self-attention-based deep neural network for predicting hERG blockers[J]. *Comput Biol Chem*, 2020, 87: 107286[2021-10-01]. <https://doi/10.1016/j.compbiolchem.2020.107286>.
- [45] Wang Y W, Huang L, Jiang S W, *et al.* Capsule networks showed excellent performance in the classification of hERG blockers/nonblockers[J]. *Front Pharmacol*, 2020, 10: 1631[2021-10-01]. <https://doi/10.3389/fphar.2019.01631>.
- [46] Lu J, Bender B, Jin J Y, *et al.* Deep learning prediction of patient response time course from early data via neural-pharmacokinetic/pharmacodynamic modelling[J]. *Nat Mach Intell*, 2021, 3(8): 696–704.
- [47] Wu J X, Gao B B, Wei X W, *et al.* Resource-constrained deep learning: challenges and practices[J]. *Sci China Chem*, 2018, 48(5): 501–510.
- [48] Wang Y, Yao Q, Kwok J T, *et al.* Generalizing from a few examples: a survey on few-shot learning[J]. *ACM Comput Surv*, 2020, 53(3): 1–34.



[专家介绍] 李国波: 四川大学华西药学院教授, 博士生导师, 课题组长。主要从事药物设计与药物化学方向, 聚焦于靶向金属酶药物设计与药物发现。主持了国家优秀青年基金项目、国家自然科学基金面上项目、四川省国际合作项目等多项科研项目。作为第一或通讯作者在 *Chem Sci*、*J Med Chem*、*Bioinformatics*、*Med Res Rev*、*Acta Pharm Sin B*、*Drug Discov Today*、*J Chem Inf Model* 等期刊发表 SCI 论文 40 余篇。申请国家发明专利 10 余项, 获授权专利 3 项, 获授权计算机软件著作权 5 项。曾获教育部自然科学一等奖、四川省科技进步奖自然科学类一等奖、四川大学学术新人奖、四川大学好未来优秀学者奖等。担任 *Eur J Med Chem* 期刊客座编辑, 担任《药学学报》中英文专刊、《中国化学快报》等期刊青年编委。