



基于人工智能的蛋白质热力学稳定性预测*

陶林节¹⁾ 徐凡丁¹⁾ 郭宇²⁾ 龙建纲^{1)**} 鲁卓阳^{1)**}

⁽¹⁾ 西安交通大学生命科学与技术学院线粒体生物医学研究所, 西安 710049;

⁽²⁾ 西安交通大学人工智能与机器人研究所, 视觉信息与应用国家工程研究中心, 人机混合增强智能全国重点实验室, 西安 710049)

摘要 人工智能技术在生物学领域的应用在近几年取得了突飞猛进的发展, 其中最显著的成果为蛋白质结构预测和设计, 该成果于2024年荣获诺贝尔化学奖。可以预见, 对蛋白质各类物理和化学属性的精准预测将是蛋白质预测领域下一阶段的重要发展方向。蛋白质热力学稳定性在深入了解生命活动机制、药物研发、疾病诊断和治疗, 以及生物技术产业中酶制剂的生产、生物传感器研发以及蛋白质药物制备等方面均具有重要意义。借助人工智能技术进行蛋白质热力学稳定性的精准预测将大幅提升蛋白质相关的科学研究能力和产业发展效率。本文综述了蛋白质热力学稳定性预测技术的发展历程, 梳理了从生物实验测定方法、传统能量函数计算方法到现代机器学习预测方法。重点讨论了基于机器学习的预测模型, 尤其是深度神经网络、图神经网络和注意力机制等前沿算法在蛋白质热力学稳定性预测中的突破。深入讨论了突变稳定性预测的核心挑战, 如数据集质量与数量不平衡、模型过拟合及蛋白质动态性的建模等难题。旨在为研究人员提供一个全面的参考框架, 助力突变蛋白质热力学稳定性预测技术的发展。

关键词 机器学习, 蛋白质热力学稳定性, 突变

中图分类号 Q68, Q816, TP18

DOI: 10.16476/j.pibb.2024.0530

CSTR: 32369.14.pibb.20240530

在众多生物大分子中, 蛋白质展示了最多样的功能和结构变化, 几乎参与了所有生物过程。这种多样性不仅源于蛋白质独特的三维结构, 还与其在生理条件下的稳定性密切相关。需要指出的是, 蛋白质的功能实现并非完全依赖于折叠状态: 自然界中, 从低等到高等生物普遍存在一定比例的无序蛋白质 (intrinsically disordered proteins, IDPs) 或全蛋白质的无序区 (intrinsically disordered regions, IDRs), 这些区域在信号转导、分子识别等关键生命活动中发挥着重要作用, 其结构的动态无序性对于功能实现至关重要。

对于多数具有稳定三维结构的蛋白质而言, 其功能完整性仍然主要依赖于其热力学稳定性^[1], 即折叠自由能, 亦即维持折叠状态的能力。在这一背景下, Anfinsen^[2] 提出的“蛋白质折叠热力学假说” (Anfinsen's Thermodynamic Hypothesis) 奠定了现代蛋白质稳定性研究基础。该假说指出, 蛋白质的天然三维结构由其一级氨基酸序列唯一决定, 并且蛋白质的折叠过程是由热力学驱动的, 即蛋白

质会自发地折叠成使吉布斯自由能最小化的构象。这个理论基础揭示了蛋白质热力学稳定性与其折叠自由能之间的直接关系, 成为理解蛋白质折叠机制和稳定性预测的关键。

尽管如此, 大多数天然蛋白质的构象稳定性相对较低, 容易在温和条件下发生变性或降解, 难以满足实际应用需求^[3]。为解决这一问题, 研究人员通常采用定点突变 (site-directed mutation) 策略, 通过改变特定位点的氨基酸残基以优化其三维结构, 从而提升整体热力学稳定性。蛋白质热力学稳定性的调控不仅对于基础研究具有重要意义, 更在实际应用中发挥关键作用。尤其在蛋白质工程中, 稳定性是决定蛋白质功能表现的关键参数, 直接影响其催化效率、表达水平及环境适应性。例

* 国家自然科学基金 (32271281) 资助项目。

** 通讯联系人。

鲁卓阳 Tel: 029-82665849, E-mail: luzhuoyang@xjtu.edu.cn

龙建纲 Tel: 029-82665849, E-mail: jglong@xjtu.edu.cn

收稿日期: 2024-12-25, 接受日期: 2025-06-04

如,在生物医药领域,蛋白质的稳定性是确保生物治疗药物在整个开发过程中的安全性、有效性和可控性的基础^[4]。在酶工程中,酶的结构稳定性决定了其活性位点是否能够维持精确构象,从而实现高效的底物结合与催化转化^[5]。此外,蛋白质热力学稳定性的丧失或错误折叠与多种疾病密切相关,如神经退行性疾病和遗传性疾病^[6]。因此,提升蛋白质热力学稳定性被认为是改善相关疾病的一个重要策略。

随着技术的不断演进,测定蛋白质热力学稳定性的方法也在不断迭代升级。早期研究主要依赖实验方法来测定蛋白质的稳定性,如圆二色性(circular dichroism, CD)光谱法^[7]、差示扫描量热法(differential scanning calorimetry, DSC)^[8],以及荧光光谱法^[9]等。这些方法通常需要大量蛋白质样品,并需多次反复实验,耗时且成本较高。随后,基于物理建模的计算方法逐渐兴起,如FoldX^[10]、Rosetta^[11]和分子动力学(molecular dynamics, MD)模拟^[12]等。这些方法通过高精度

的能量函数与构象采样策略,不仅提供了高精度的预测结果,还揭示了蛋白质结构与动力学的详细信息,现仍广泛用于蛋白质热力学稳定性研究。

近年来,基于机器学习的蛋白质热力学稳定性预测算法受到广泛关注(图1)。大量的研究致力于开发高效且可靠的计算工具,以预测突变对蛋白质热力学稳定性所造成的影响。这类方法通过分析大规模的蛋白质结构和稳定性数据库,提取包括氨基酸替代类型、局部二级结构、溶剂可及性(solvent accessibility)、相邻残基的氨基酸组成等特征,并结合从传统线性回归到复杂深度学习算法的多种机器学习方法,进行建模与预测。当前已有多种算法成功应用于该任务,例如支持向量机(support vector machine, SVM)^[13]、人工神经网络(artificial neural network, ANN)^[14]、随机森林(random forest, RF)^[15]和3D卷积神经网络(3D convolutional neural network, 3D-CNN)^[16]等算法,均在预测突变导致的蛋白质稳定性变化方面展现出良好的性能。

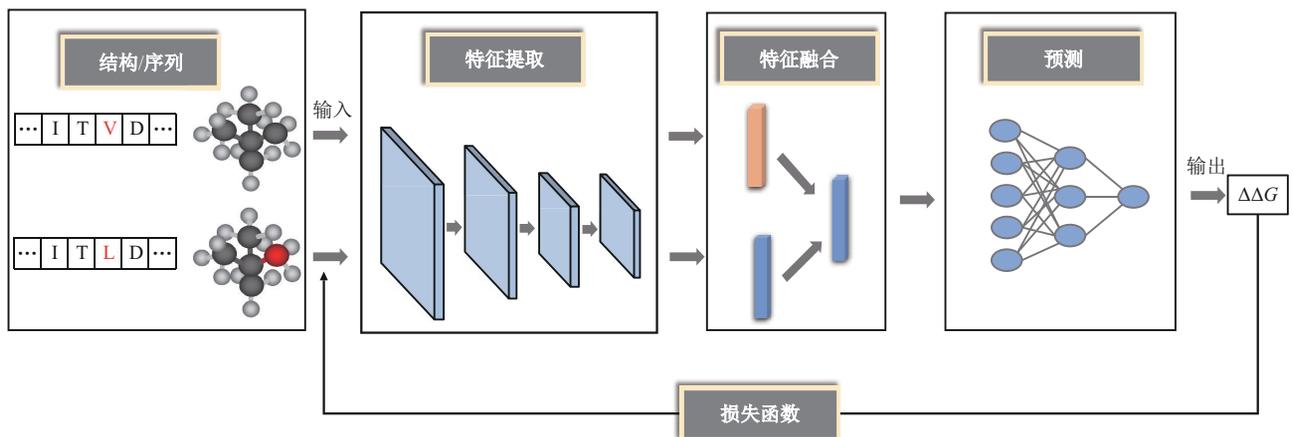


Fig. 1 Prediction of protein stability based on artificial intelligence mutations

图1 基于人工智能的突变蛋白质稳定性预测

蛋白质的稳定性是一个多维概念,其功能实现通常依赖于特定的三维结构,而热力学稳定性(即折叠自由能(ΔG))描述了蛋白质维持其折叠状态相对于变性状态的能力。蛋白质的稳定性不仅仅由其折叠状态决定,还涉及诸如寡聚、无定形沉淀和降解性等其他方面,这些特性在蛋白质的生物学功能中发挥着至关重要的作用。然而,目前计算方法仍主要聚焦于预测突变对蛋白质折叠自由能的影响,以评估其热力学稳定性。对于大多数具有稳定

三维结构的蛋白质而言, ΔG 作为热力学稳定性的关键参数,可通过其折叠态与去折叠态之间的吉布斯自由能差即 ΔG 表征,其计算公式如下:

$$\Delta G = G_{\text{folded}} - G_{\text{unfolded}} \quad (1)$$

氨基酸突变会引起蛋白质的 ΔG 的变化,而突变后稳定性的变化可以用 ΔG 的变化来量化,即折叠自由能差($\Delta\Delta G$)表示:

$$\Delta\Delta G = \Delta G_{\text{mutant}} - \Delta G_{\text{wild}} \quad (2)$$

$\Delta\Delta G$ 通常以kcal/mol为单位,其符号用于指示突变对蛋白质稳定性的影响方向。若 $\Delta\Delta G$ 为正值,

表示突变导致折叠自由能增加,即蛋白质稳定性降低(不稳定突变体);若 $\Delta\Delta G$ 为负值,则表示突变降低了折叠自由能,即增强了蛋白质的稳定性(稳定突变体)。

$\Delta\Delta G$ 的一个重要性质是反对称性(anti-symmetry),其对应的吉布斯自由能变化为逆变化(即从野生型到突变型的折叠自由能)是相等的,并且与直接变异的变化有相反的符号:

$$\Delta\Delta G_{\text{forward}} = \Delta G_{\text{mutant}} - \Delta G_{\text{wild}} = -(\Delta G_{\text{wild}} - \Delta G_{\text{mutant}}) = -\Delta\Delta G_{\text{reverse}} \quad (3)$$

虽然机器学习在自动提取数据中的潜在特征与复杂模式方面展现出卓越的能力,但其性能易受到数据噪声和虚假相关性的影响,尤其是在训练数据存在显著偏差或样本量有限的情况下。目前,蛋白质热力学稳定性相关的高质量实验数据仍为稀缺。现有数据库中,不仅蛋白质之间存在较高的序列相似性,实验条件的变化也导致 $\Delta\Delta G$ 测量结果存在固有误差,这些因素均为预测模型的泛化能力构成挑战。此外,由于忽略突变体与天然蛋白质形式之间 $\Delta\Delta G$ 值的非对称性,进一步放大了模型预测中的系统性偏差。目前,多数现有模型依赖于浅层机器学习算法,这在一定程度上与实验训练数据的稀缺性有关,限制了更复杂更深度模型训练于应用。因此,发展鲁棒性更强、能有效缓解数据偏倚影响的预测方法显得尤为关键。

本综述介绍了突变蛋白质热力学稳定性测量方法的发展历程,以及人工智能在蛋白质工程的应用,并重点介绍和探讨了当前用于预测突变所致蛋白质热力学稳定性变化的前沿技术与方法。希望本综述能为相关领域的研究人员提供有益的参考,并为蛋白质热力学稳定性预测领域的进一步研究提供启发性思考。

1 蛋白质热力学稳定性测量的发展

蛋白质的热力学稳定性直接影响其在表达、纯化和储存过程中的最佳条件^[17],因此准确测量蛋白质热力学稳定性是蛋白质工程中的一项重要任务。早期研究主要依赖生物实验手段评估突变体与野生型蛋白质之间稳定性差异。尽管实验方法在精确性方面具有显著优势,但其测量过程通常耗时且成本较高。为了克服实验方法的局限,研究人员提出基于物理原理的计算方法进行预测。虽然该方法对计算资源的需求较高,且模型预测精度依赖于力场参数和初始条件,但仍然得到了广泛应用。随着

人工智能(artificial intelligence, AI)技术的发展,越来越多的研究致力于将机器学习与深度学习方法引入蛋白质热力学稳定性预测领域。AI驱动的预测模型在效率与性能方面展现出显著潜力,正逐步弥补传统方法在准确性、可扩展性与适用性方面的不足,成为该领域的重要研究方向之一。

1.1 生物实验测定

早期生物实验方法测定蛋白质热力学稳定性主要通过物理化学技术,这些方法的核心原理是利用蛋白质在变性过程中的物理化学性质变化来推断其稳定性。性质变化包括热力学参数、光学特性和荧光响应等。常用的实验方法包括CD光谱法、DSC和热变性展开法等。虽然这些方法在灵敏度和操作复杂度上有所差异,但均为研究蛋白质的稳定性提供了可靠的数据基础。

CD光谱术是一种利用手性分子对左、右圆偏振光吸收差异的光谱技术^[18]。在远紫外波段(180~250 nm),蛋白质的二级结构(如 α 螺旋和 β 折叠)对圆偏振光表现特征吸收,产生独特的光谱,由此可推断蛋白质构象状态。通过在可逆热展开或化学展开过程中进行CD测量,CD可用于绘制蛋白质折叠和去折叠的转变曲线,进一步计算 $\Delta\Delta G$ 。然而,CD技术在解析三级或四级结构以及特定的氨基酸残基识别方面的分辨率有限,且对溶剂条件的变化不够敏感。

DSC是一种定量测量样品热熔变化的热分析技术,能够获取包括蛋白质折叠过程的自由能变化(ΔG_{fold})、焓变(ΔH_{fold})、熵变(ΔS_{fold})及熔化温度(T_m)和热容(C_p)在内的多项热力学参数。其工作原理通过对比含有蛋白质的样品池与不含蛋白质的参比池的温差,进而获得热吸收信息。DSC的优势在于无需标记蛋白质即可从单一实验中获取全面的热力学参数。尽管目前在提高通量方面存在一定挑战,例如减少样本量、并行读出和确保样品平衡。目前,DSC仍是研究蛋白质热稳定性和结构完整性的重要工具。

荧光光谱法则通过监测蛋白质内源荧光团(如色氨酸残基)在变性前后的激发与发射特性变化,反映其构象状态。由于蛋白质结构变化会改变荧光团所处的微环境,进而引起荧光强度和最大发射波长的变化。该方法具有高灵敏度、操作简便及非破坏性的优点,适用于多种蛋白质样品状态。然而,荧光光谱测量也存在一些缺陷,如对某些元素(例如氢、碳等轻元素)的敏感度较低,易受元素相互

干扰和重叠峰的影响, 以及定量分析需要标准样品^[19]。

热变性展开法和化学变性法广泛用于测定蛋白质的热力学稳定性, 即折叠态与去折叠态之间的 ΔG (公式 (1))。热变性展开法通过测定升温过程中的蛋白质溶解度或活性变化来评估稳定性, 常与CD、荧光光谱或DSC结合使用, 但受限于温度控制精度和实验条件的稳定性^[20]。在热变性实验中, 蛋白质溶液在恒定升温速率下进行处理, 利用光谱或DSC监测构象变化或热效应。关键参数包括 T_m 、 $\Delta H(T_m)$ 和 ΔC_p , 可用于计算蛋白质的稳定性函数($\Delta G(T)$)。DSC作为蛋白质热力学稳定性测量的金标准^[21], 虽不如CD光谱法那样能提供详细的结构变化信息, 但DSC本身也存在灵敏度低和通量有限等局限性^[22]。

目前, DSC因其高精度被认为是最主流的方法, 但由于成本高和通量限制, 其在高通量应用中的适用性受限。因此, CD光谱法和荧光光谱等技术在高通量筛选中的应用更为广泛。近年来, 随着高通量自动化实验设备、微流控技术和机器学习辅助数据分析的发展, 蛋白质稳定性实验的效率和数据处理能力得到了显著提升。例如, 基于纳米DSC (nanoDSC) 的小样本测量和基于微孔板的高通量荧光筛选正在成为主流趋势。未来, 结合计算模拟 (如MD模拟) 和机器学习预测, 有望进一步提高实验数据的利用率, 降低实验成本, 同时提升蛋白质稳定性研究的效率和准确性。

1.2 基于物理方法的计算预测

基于物理的计算方法通过计算机模拟和能量计算, 有效减少了对实际蛋白质样品、实验设备和试剂的需求, 从而降低了研究成本。这些方法不仅显著缩短了实验周期, 还为对蛋白质热力学稳定性的快速预测与分析提供了高效手段。传统的蛋白质热力学稳定性评估方法依赖于基于物理的建模工具、统计分析方法与传统机器学习技术。期望在准确性和计算效率之间取得平衡。目前常用的基于物理的预测工具包括FoldX、Rosetta和MD模拟, 它们依赖于不同的能量函数和统计模式来评估突变对蛋白质稳定性的影响。

FoldX是一种快速而广泛应用的能量计算工具, 其核心是将范德华相互作用、静电势、氢键和溶剂化贡献整合至能量函数中。FoldX通过在一组实验数据上拟合模型参数, 从而实现了对突变引起的 $\Delta\Delta G$ 的预测^[23]。在模拟突变时, FoldX采用侧链

旋转体方法, 允许侧链构象的调整而保持主链不变。在计算效率方面, FoldX计算单点突变的速度约为Rosetta的一半^[24]。但是, FoldX在不同基准测试中, 其与实验数据的相关性差异显著, 预测性能的皮尔逊相关系数在0.19~0.81之间不等^[25-26]。这种差异既受到FoldX模型误差的影响, 也与数据范围和分布有关。

Rosetta是一款功能强大的蛋白质建模与设计平台, 其特色在于采用蒙特卡洛模拟退火采样策略, 能够在庞大的构象空间中有效搜索并评估蛋白质构象。Rosetta的核心是基于能量函数的优化, 使用统计势来预测和设计蛋白质的结构和功能。这个能量函数综合了多种因素, 包括蛋白质内部的化学键、氢键、范德华力等相互作用, 以及蛋白质在溶剂中的行为。相较于传统力场方法, 时间效率被认为是Rosetta的一大优势, 这种在时间效率上的优势主要体现在其能够快速评估能量和预测蛋白质结构, 尤其是在处理大规模构象空间时。然而, 这种高效性能往往以较大的计算资源消耗和更长的模拟时间为代价, 尤其是在模拟复杂结构或进行高通量预测时, 其性能弱势较为突出。

MD模拟则提供了一种时间维度上的蛋白质行为洞察手段。该方法基于牛顿运动定律, 通过数值积分方式模拟分子系统中各原子的运动轨迹, 从而获取蛋白质系统在不同时间尺度 (从纳秒到微秒) 上的动态特性。MD模拟能够揭示构象变化、折叠路径、配体结合与溶剂相互作用等重要生物物理过程。通过对比突变体与野生型蛋白质在多种条件 (如温度、pH、溶剂环境) 下的动力学行为, 研究人员可进一步评估突变对蛋白质稳定性的影响。虽然MD模拟在精度与解释性方面具有显著优势, 但其对计算资源的需求极高, 尤其是在大规模、多样性突变分析中。

2 人工智能促进生物医学工程研究

AI是一门通过开发计算机系统以模拟和执行人类智能任务的科学。AI的核心能力包括感知、理解、推理和学习, 使其具备处理复杂问题的强大潜力。近年来, AI技术已广泛应用于生物医学、材料科学以及工业优化等领域, 尤其在生物学和医学领域, 通过快速处理和分析大规模生物数据, AI有效促进了个性化医疗、精准医学、新药开发以及蛋白质功能预测等前沿方向的发展。

机器学习 (machine learning, ML) 作为AI的

核心分支, 通过构建可以从数据中学习的算法, 使系统能够适应新输入并改进预测能力^[27]。在语音识别、自然语言处理以及基于移动设备、智能手表等可穿戴设备所采集的传感器数据进行人类行为模式识别等应用领域中, 各类机器学习算法均发挥着关键作用。根据学习机制不同, ML 通常分为监督学习、无监督学习、半监督学习及强化学习^[28-31]。其中, 监督学习通过拟合线性方程 (如线性回归) 预测连续值, 或通过逻辑回归处理二分类问题 (判断突变是否破坏稳定性)。例如, 决策树基于特征阈值构建树形结构进行分类, 梯度提升树 (如 XGBoost/LightGBM) 通过集成多棵弱决策树优化预测精度。无监督学习则不依赖标签, 通过识别数据中的潜在结构来提取信息。例如, K 均值聚类可用于将蛋白质突变数据划分为相似群组, 主成分分析 (principal component analysis, PCA) 可用于降维并保留主要特征。强化学习通过与环境交互并优化长期回报, 实现策略优化。近年来在药物开发和医学影像分析等领域取得了广泛应用。半监督学习结合少量标记数据和大量未标记数据进行学习, 其中标签传播已在医学影像分类与小样本稳定性预测任务中表现出色。

深度学习 (deep learning, DL) 是 ML 的重要发展方向, 其通过多层神经网络结构在无需显式特征工程的条件下自动学习数据中的复杂模式。在大规模数据和计算资源支持下, DL 广泛应用于图像识别^[32]、图像分类^[33]和运动识别^[34]等任务。典型架构如卷积神经网络 (convolutional neural network, CNN) 和循环神经网络 (recurrent neural network, RNN), 已在医学图像处理、序列建模和语义理解展现出优异性能。后续模型如长短期记忆 (long short-term memory, LSTM) 等 DL 架构, 在处理复杂、非结构化数据时表现卓越。

随着 AI 技术的不断发展, 一系列新型算法为蛋白质稳定性预测注入新动力。扩散模型^[35] (diffusion models) 通过正向噪声添加与反向去噪过程生成高质量数据, 能够模拟蛋白质构象动态变化, 辅助设计具备更高稳定性的蛋白质变体。语言模型^[36] (如 Transformer 和进化尺度蛋白质语言模型 (evolutionary scale protein language model, ESM)) 基于注意力机制捕捉序列上下文关系, 通过大规模预训练后微调至特定任务, 从蛋白质序列中提取深层语义特征。几何深度学习^[37] (geometric deep learning, GDL) 则基于蛋白质结

构的对称性与物理约束 (如旋转不变性), 直接处理分子点云或三维结构, 分析突变所引起的局部几何扰动, 评估稳定性变化。表 1 中详细列出了这些方法的应用及其具体描述。

值得注意的是, 随着 AI 方法在蛋白质稳定性建模中的深入应用, 蛋白质适应性预测也逐渐成为突变效应研究的重要分支。其核心目标是判断特定突变在自然选择或生理压力下是否具有进化优势, 重点评估突变是否破坏蛋白质功能或改变其在特定条件下的适应能力。与热稳定性预测 (聚焦于 $\Delta\Delta G$ 或 ΔT_m 等物理变化) 相比, 适应性预测更侧重于突变对生物功能的影响, 通常依赖于多序列比对 (multiple sequence alignment, MSA) 和进化保守性等信息, 常以功能保持情况或细胞表型变化 (如细胞生长率) 作为监督标签, 建模逻辑更偏向“功能敏感性”。在方法选择上, 适应性预测多采用统计模型或生成模型框架。例如: EVmutation^[38] 基于最大熵模型, 捕捉残基间的共进化关系; DeepSequence^[39] 利用变分自编码器 (variational autoencoder, VAE) 增强多突变场景下的泛化能力; AlphaMissense^[40] 模型则融合 AlphaFold2 的结构信息与蛋白质语言模型嵌入, 借助神经网络实现对突变致病性的精确评估, 展现出强大的综合建模能力。

适应性预测与热稳定性预测在输入形式、建模逻辑和监督标签等方面存在差异, 但两者在底层架构上高度相关, 常共享特征提取模块。适应性预测通常基于野生型序列建模, 评估突变在“进化语义空间”中的扰动程度来判断其可能的功能影响, 属于“序列敏感性驱动”的建模范式; 而热稳定性预测则侧重于突变前后状态的显式差异, 基于序列或结构输入估算 $\Delta\Delta G$ 或 ΔT_m , 体现“状态变化驱动”的建模逻辑。从生物学角度看, 两者密切相关但不等价。结构稳定性是蛋白质实现功能的前提, 热稳定性常被视为适应性的物理基础。然而, 部分突变即便降低热稳定性, 仍可能通过增强酶活性或调控特异性获得适应优势, 从而被正向选择。因此, 在适应性评估中, 热稳定性既是重要参考, 也存在例外情形。

尽管适应性预测与热稳定性预测在建模目标和监督信号上各有侧重, 但在实际应用中, 二者往往存在交叉需求。例如, 在疾病突变筛查或蛋白质设计任务中, 既需要评估突变对功能影响, 也需要考量其对结构稳定性的干扰。然而, 当前一些先进模

型并未在训练阶段显式整合稳定性标签, 更多是通过结构建模或进化表征间接反映蛋白稳定性。如 AlphaMissense^[41] 基于语言模型的表示与群体遗传标签, 虽可间接反映进化适应性, 却未显式引入稳定性指标, ESM-IF^[42] 则聚焦于结构适配性, 稳定

性评估则依赖后处理, 而非端到端学习。因此, 将功能性与稳定性信息整合至统一建模框架, 仍是一个值得深入探索的方向, 有望在提升预测性能的同时, 增强模型对突变机制的生物学解释能力。

Table 1 Machine learning methods for predicting mutation effects

表1 预测突变效应的机器学习方法

方法	描述	文献
卷积神经网络 (CNN)	一种神经网络架构, 通过卷积滤波器处理蛋白质结构数据, 类似于常见的图像处理网络	[43]
图神经网络 (GNN)	一种神经网络架构, 它将蛋白质中的原子或残基作为节点, 它们之间的关系作为边, 构建成图结构	[44]
支持向量机 (SVM)	一种传统的监督式机器学习方法, 学习在特征空间中划分不同类别的边界。这是应用于预测突变效应的最古老的机器学习方法之一	[45]
迁移学习	一种机器学习技术, 它允许模型将从一个任务学到的知识应用到另一个相关任务上, 以提高学习效率和性能	[46]
语言模型 (Transformer)	一种非常强大的神经网络架构, 学习特征嵌入空间, 并结合注意力机制从序列数据中进行预测	[47]
随机森林 (RF)	一种集成学习方法, 通过构建多个决策树并结合它们的预测结果来提高模型的准确性和鲁棒性	[48]
扩散模型 (Diffusion)	一种生成模型, 通过逐步向数据添加噪声并学习去除噪声来恢复原始数据, 类似于从模糊图像中还原清晰细节的过程	[35]
几何深度学习 (GDL)	一种处理图结构数据的神经网络架构, 通过图卷积操作聚合节点信息, 类似于在复杂网络中传递和更新信息, 适用于蛋白质结构和生物网络分析	[37]

CNN: 卷积神经网络 (convolutional neural network); GNN: 图神经网络 (graph neural network); SVM: 支持向量机 (support vector machine); GDL: 几何深度学习 (geometric deep learning); Transformer: 一种语言模型, 基于注意力机制的序列建模架构; RF: 随机森林 (random forest)。

3 蛋白质热力学稳定性预测计算方法

3.1 传统机器学习方法的进展与挑战

ML 模型已广泛应用于蛋白质热力学稳定性变化 (如 $\Delta\Delta G$) 的预测, 这一趋势主要由两个独立因素驱动: 一是大规模蛋白质序列和结构数据集的可获得性显著提升, 二是计算机视觉和自然语言处理算法的重大突破被有效迁移至蛋白质生物化学问题的建模中。随着 ProTherm 数据库的建立, 该数据集系统性地整合了大量突变蛋白质的热力学实验数据^[49], 研究者们开始尝试使用各类 ML 方法来替代或增强传统基于物理建模的预测策略。这些方法能够充分利用大规模突变- $\Delta\Delta G$ 对的数据, 融合序列、结构及其他辅助信息, 构建统一的多模态特征空间, 以支持更全面且高精度的预测任务。

早期研究提出了许多基于监督学习的模型, 广泛采用分类与回归技术以预测蛋白质突变所引起的稳定性变化方向 (增强或削弱) 及其具体数值。例如, Capriotti 等^[50] 提出了一种基于 ANN 的预测方法, 在预测蛋白质变体的稳定性 (稳定/不稳定)

方面达到了 81% 的准确率, 显著优于当时基于能量函数的物理模型。Frenz 等^[51] 构建的前馈神经网络模型在特定数据集上达到了 92.8% 的准确率。此外, SVM^[52]、神经网络^[50] 和基于决策树的随机梯度提升算法 (stochastic gradient boosting based on decision trees, SDB-DT)^[53] 等 ML 模型, 也被成功应用于突变蛋白质热力学稳定性的预测, 并取得了良好的性能表现。

这些方法在预测性能上看似表现优异, 但实际上其高准确率往往是基于低质量、重复性较高且误差较大的数据集验证的结果, 这导致了预测结果的虚假优越性, 且存在过拟合的风险, 影响模型的可靠性和泛化能力。有研究指出, 许多 ML 方法在预测中存在显著偏差^[54], 表现为对降低稳定性的突变预测效果优于对增强稳定性的突变预测。这种偏差削弱了模型在稳定突变情形下的预测能力, 从而导致整体预测结果与实验数据之间虽呈现较高的线性相关性, 但该相关性未能在稳定突变数据上得到充分体现。造成这种偏差的主要原因之一在于数据集的质量问题, 这限制了模型的泛化能力和准确

性。2018年的一项研究表明, ProTherm数据库中大量关于序列、结构和稳定性数据的错误^[55]。该数据库标注为 $\Delta\Delta G$ 的4 148个条目中, 仅有1 197个(约占29%)被认为是可靠的。为了改善数据质量, 后续研究者开发了多个经手工筛选和严格质量控制的数据集, 如 Ssym^[56]、ThermoMutDB^[57]和 S669^[58], 这些新数据集在实验来源、测定条件和数据一致性方面均显著优于早期数据集。然而, 这些数据集在规模上相对有限, 且人工整理过程中可能引入新的主观偏差, 仍对算法的训练与评估带来挑战。2020年 Fang^[59]提出使用假设性反向突变策略来评估5种主流ML算法的鲁棒性, 发现许多模型存在明显的过拟合现象。Pancotti等^[58]在 S669数据集上评估了21种蛋白质热力学稳定性预测器的反对称性(anti-symmetry)表现, 结果表明, 仅有5个预测器符合反对称性要求。因此, 设计具有反对称性并能在有限数据条件下有效克服过拟合问题的模型, 成为亟待解决的难题。

3.2 深度学习模型的突破

DL的兴起为蛋白质热力学稳定性预测领域带来了新的研究活力, 表2列举了近年来具有代表性的预测模型, 展现了深度学习在该领域的应用成果。与传统ML相比, DL通过自动化特征提取显著减少了对手工特征的依赖, 能够学习更复杂的非线性映射关系, 从而在预测精度方面实现了进一步提升。然而, 数据规模有限和实验数据的噪声仍是当前面临的挑战, 促使研究者们不断探索新型神经网络架构。例如, SCONES模型^[60]通过引入蛋白质对称性和传递性原则, 计算残基对的相互作用能量, 有效提升了突变对蛋白质热力学稳定性的影响的预测能力, DeepDDG^[61]模型则设计了“共享残基对(shared residual pair, SRP)”网络结构, 以共享残基对之间权重来学习突变的影响, 充分利用了蛋白质的局部结构信息, 在多个数据集中的性能评估结果优于其他方法。

图神经网络(graph neural network, GNN)因其强大的处理蛋白质三维结构信息的能力, 成为近年来最受关注的技术路线之一。GNN将蛋白质的原子或残基表示为图中的节点, 原子间键合或空间邻近关系建模为图的边, 从而有效捕捉分子内部的空间关系, 提升预测性能。例如, ProS-GNN模型^[62]采用门控GNN捕捉蛋白质结构与属性之间的复杂关系, 从而提高模型性能, ABYSSAL^[63]

模型融合光注意力机制, 提高了突变前后氨基酸的识别能力, Pythia^[64]则结合自监督训练模块与GNN, 在预测精度和计算速度(提升105倍)上均实现突破。Li等^[65]提出的3D-CNN进一步增强了局部结构环境的空间建模能力, 有效缓解传统模型在预测不稳定突变时的偏倚问题。mutDDG-SSM^[66]则通过图注意力网络(graph attention network, GAT)和自监督学习结合XGBoost, 进一步提高了稳定性预测的准确度, 并有效缓解了过拟合现象。

此外, 基于单序列的结构预测模型(如ESMFold和OmegaFold)以及多模态融合模型(例如将序列、结构、进化和理化特征结合的模型)正在逐步推动蛋白质稳定性预测向更广泛和更精细的方向发展。然而UniProt数据库中只有不到0.2%的蛋白质拥有三维结构, 因此蛋白质结构数据的稀缺也限制着模型的潜力。随着蛋白质结构预测技术(如AlphaFold3)的迅速发展, 该问题正逐步得到缓解。这些方法借助深度3D-CNN, 可从预测结构中提取原子类型、相邻原子、氢原子的数量等特征, 从而构建更具表达力的结构表示。OmeDDG^[67]首次结合OmegaFold预测的结构信息进行蛋白质热力学稳定性变化预测, 且在多个盲测数据集上表现良好, 并展示了良好的反对称性; PROST^[68]则基于AlphaFold2生成的结构描述符, 通过集成学习策略实现了无结构条件下的高精度预测。尽管基于单序列的结构预测方法(如ESMFold、OmegaFold)可在几秒的时间尺度内完成对泛型蛋白质的结构预测, 速度超过从MSA中提取的进化信息来预测蛋白质结构的AlphaFold的数量级, 但在预测精度上, AlphaFold结果更接近实验结果。而近期提出的一种新的基于单序列结构预测模型SPIRED^[69], 通过优化推理效率和建模精度, 在蛋白质热力学稳定性预测任务中表现出与AlphaFold相近预测精度的同时, 推理速度提升约5倍, 训练资源消耗降低一个数量级。同时, SPIRED在蛋白质热力学稳定性预测任务中展现出先进性能, 成为兼顾效率与准确性的重要模型。

GAT、Transformer架构以及多模态融合方法被广泛应用于 $\Delta\Delta G$ 预测中, 分别在结构建模、序列建模和多源信息整合方面展现出关键作用, 它们正在成为推动这一领域发展的关键算法。GAT通过加权节点间的信息传递机制, 能够更加精确地捕捉蛋白质结构中的关键依赖关系, 从而提升模型对

复杂三维结构信息的理解和稳定性预测的准确性。GAT在蛋白质稳定性预测中的应用,尤其在识别局部环境对稳定性变化的影响方面,表现出了明显优势。Transformer架构凭借其自注意力机制能够高效处理序列数据,特别是在捕捉长程依赖关系时展现出独特的优势,可以精确地提取蛋白质序列中

的特征;多模态模型则融合序列、结构及进化信息,实现多源数据的互补,增强模型的泛化能力。上述方法在一定程度上克服了传统方法对特定工程和线性关系建模的依赖,为蛋白质热力学稳定性建模开辟了更具有前景的新方向。

Table 2 An overview of AI-based prediction models for protein thermodynamic stability changes ($\Delta\Delta G$)

表2 基于人工智能的蛋白质热力学稳定性变化 ($\Delta\Delta G$) 预测模型概述

方法 (年份)	三维 结构	特征类型	RMSE/ (kcal·mol ⁻¹)	<i>r</i>	人工智能方法	文献
BayeStab (2022)	√	突变残基邻接矩阵、突变点坐标	/	0.53	贝叶斯神经网络	[70]
ABYSSAL (2023)		野生型和突变型蛋白质嵌入表示	1.32	0.37	ESM-2模型+光注意力机制	[63]
mutDDG-SSM (2024)	√	原子类型、原子所属的氨基酸类型、DSSP二级结构类型、溶剂可及表面积	1.69	0.49	图注意力网络(GAT)+极端梯度增强(XGBoost)	[66]
PSP-GNM (2022)	√	序列特征、结构特征(突变位点、溶剂可及性、可达表面积、疏水性、和二级结构)	1.53	0.59	高斯网络模型	[71]
PROST (2022)		突变残基邻域、相对溶剂可及性、净体积、净疏水性、突变疏水性以及氨基酸评分等	1.46	0.64	极端梯度提升决策树+额外树回归器	[68]
ThermoMPNN (2024)	√	蛋白质序列嵌入、结构特征(原子距离)、突变点坐标	1.52	0.43	卷积神经网络+迁移学习	[72]
THPLM (2024)		蛋白质嵌入表示	1.63	0.53	ESM-2+卷积神经网络	[73]
Stability Oracle (2024)	√	原子坐标邻接矩阵、野生和突变型的残基微环境	1.43	0.52	图神经网络+变压器	[74]
ProSTAGE (2024)	√	空间邻接矩阵,蛋白质嵌入,野生型残基微环境	1.37	0.7	图卷积网络+蛋白质语言模型	[75]

RMSE: 均方根误差 (root mean square error); *r*: 皮尔逊相关系数 (Pearson correlation coefficient); DSSP: 蛋白质二级结构字典 (dictionary of secondary structure of proteins); XGBoost: 极端梯度提升 (extreme gradient boosting); GAT: 图注意力网络 (graph attention network); √: 表示模型使用了蛋白质结构信息; /: 表示模型未提供RMSE结果; 所有RMSE和*r*值均基于“S669数据集的测试结果”。

3.3 蛋白质语言模型 (PLMs) 与未来发展

随着蛋白质结构预测技术的发展,蛋白质语言模型 (protein language models, PLMs) 在结构与功能预测任务中展现出卓越的潜力。PLMs能够从一级序列中捕获原子级别的结构特征,有助于深入理解单点突变对蛋白质稳定性和功能的影响。通过在数亿个天然蛋白质序列上进行无监督预训练,PLMs能够生成高质量的蛋白质序列嵌入,这些嵌入在多个蛋白质设计任务中表现出良好的泛化能力和预测性能。在蛋白质热力学稳定性预测任务中,PLMs的拟合倾向较低,并且其测试集性能上过拟合或高估的风险也较低。然而,在训练过程中,PLMs不可避免地会产生与现有模型不同的偏差,这种偏差在某些情况下可能是互补的^[76]。ThermoMPNN^[72]通过迁移学习策略,将预训练模

型ProteinMPNN提取的序列嵌入特征迁移至蛋白质稳定性优化任务,实现了知识的有效迁移;DDMut^[77]模型结合局部三维环境表示、CNN和Transformer编码器,更加精准捕捉短距离与长距离相互作用,在稳定性预测中达到了0.70的Pearson相关性,并展现出良好的反对称性。

虽然PLMs在蛋白质设计中的应用取得了显著进展,但数据稀缺和偏差问题依然是限制其进一步发展的重要难题。随着蛋白质数据库的不断扩展,PLMs在大规模数据上的训练的表现不断优化,但在处理小样本学习任务 and 显著类别不平衡的数据集中,模型泛化能力和稳定性仍不理想。为了应对这一难题,越来越多研究开始引入数据增强技术和跨模态学习策略。例如,Stability Oracle^[74]模型通过热力学排列 (thermodynamic ranking) 生成新的伪

标签样本进行数据增强,显著扩大了训练数据集规模,提升了模型在独立测试集上的泛化能力。在T2837测试集上的测试结果显示,该模型能正确识别出48%的稳定突变($\Delta\Delta G < -0.5$ kcal/mol),其中74%的预测确实是稳定的。

同时,随着蛋白质结构预测技术(如AlphaFold2)的不断成熟,结构信息的可获得性逐渐提升,使得基于序列的蛋白质稳定性预测模型得以融合更多空间特征以增强表达力。例如,ProSTAGE^[75]融合蛋白质结构和序列嵌入,引入一个空间节点特征来捕获突变附近的残基相互作用。该模型利用蛋白质语言模型的序列嵌入作为图卷积网络(graph convolutional network, GCN)的节点输入,进一步验证了PLMs在突变效应预测中的潜力。

面向未来,PLMs的发展需在计算效率与模型可解释性两个维度实现突破。一方面,尽管现有模型在预测精度方面已取得一定的提升,但如何实现大规模预测仍是技术难点。随着计算资源的优化,未来PLMs需在保持高精度的同时提升推理速度,以支持蛋白质工程与药物发现中的高效应用。另一方面,增强模型可解释性(如可视化突变对稳定性的影响路径)对提升临床与工业界的信任度至关重要。在技术融合方面,PLMs与GNN的结合可同步捕捉蛋白质三维结构的局部环境与全局拓扑信息,而自监督学习的引入则能进一步挖掘序列中的隐含特征。多模态数据整合将成为另一核心方向:通过融合序列、结构、生物物理等多源信息,辅助PLMs突破单一数据源的局限性,例如,利用AlphaFold预测的结构数据构建多模态训练集,或通过生成对抗网络(generative adversarial networks, GANs)模拟高质量突变数据,可有效缓解数据稀缺问题。随着这些技术的协同优化,PLMs将在精准蛋白质设计、稳定性预测及靶向药物开发中释放更大潜力,最终为疾病研究与生物工程提供兼具高效性、可解释性与实用性的工具。

4 总结与展望

蛋白质的稳定性在研究由不稳定蛋白质引发的疾病机制以及开发具有特定功能的工程蛋白质中具有不可或缺的作用。因此,提升蛋白质热力学稳定性不仅对生物医学研究意义重大,也为工业应用提供了广阔前景。近年来,蛋白质热力学稳定性预测领域取得了显著进展,已形成了实验方法与计算模

型协同发展的技术体系。在数据基础层面,ProTherm、FireProtDB等数据库的建立为算法开发提供了关键支撑;在方法学层面,基于物理的算法(如Rosetta、FoldX)与ML模型(如DeepDDG、ThermoNet)实现了从序列到结构的稳定性预测;在技术融合方面,AlphaFold等结构预测工具的突破有效填补了未知蛋白质结构信息的空白。这些发展为疾病机制解析、工程蛋白质设计等应用场景奠定了务实的基础。尽管该领域在数据集建设、算法创新和实验验证方面均取得了诸多成果,但仍存在若干需解决的关键挑战。本文综述了蛋白质热力学稳定性预测领域中实验方法和计算方法的最新进展,并对未来发展方向进行了探讨。

当前,蛋白质热力学稳定性预测领域有望从以下3个方面取得突破。

首先,数据局限性层面。现有实验数据集(如ProTherm)不仅规模有限,且存在显著偏差——仅29%的数据可靠性经严格验证,约85%的突变类型局限于单点变异。这种数据失衡直接导致预测模型对复杂场景(如多点协同突变)的泛化能力不足。Fang^[59]的研究进一步指出,现有算法对突变位点空间动态变化的表征能力薄弱,其根源在于过度依赖突变后已知结构数据。然而,约70%的人类蛋白质缺乏实验解析结构,这一缺口虽因AlphaFold的突破性进展得到部分填补,但其局限性依然突出:一方面,AlphaFold预测的突变体与野生型结构差异较小^[60],难以捕捉折叠缺陷;另一方面,其输出的pLDDT指标与实验 $\Delta\Delta G$ 值相关性较弱或没有^[61],无法直接关联热力学稳定性。因此,未来研究需要重点关注数据整合与增强技术。通过整合多源数据,特别是结合实验数据与计算生成的模拟数据,或许可以显著弥补单一数据来源的不足。例如,结合蛋白质热力学数据、动力学参数及多态性结构信息的多模态数据融合技术,可能成为突破数据瓶颈的关键途径。此外,开发基于结构预测模型生成的虚拟数据集并结合迁移学习,也为模型训练提供了新的思路。

其次,解决蛋白质结构预测动态性问题。蛋白质的 $\Delta\Delta G$ 值不仅取决于静态结构的差异,还与蛋白质的动态特性密切相关。尽管蛋白质结构及其构象预测都取得了重大突破,但动态结构的研究仍然相对滞后。以AlphaFold为例,其能够准确预测蛋白质的三维结构,但只能预测蛋白质在一个瞬间的静态结构,尚无法实现动态变化的预测。因此,如

何将蛋白质的动态变化纳入稳定性预测中, 仍然是当前研究的重大挑战。其中, 复旦大学的研究团队提出了4D扩散模型AlphaFolding^[78], 该模型结合了MD模拟数据, 能够预测蛋白质在多个时间步长上的动态结构变化。通过MD与DL技术相结合, 从而捕捉多突变组合对蛋白质热力学稳定性的协同影响。此外, 模型还可以融入生物物理约束条件, 通过模拟溶液条件、温度变化或其他环境因素对蛋白质热力学稳定性的影响, 也有望提供更加精确的突变影响评估。

最后, 算法需要多元化的创新和突破。在蛋白质热稳定性预测领域, 当前算法面临同质化的困境。许多预测模型仍然依赖于单一的DL架构, 尤其是GNN和CNN。虽然这些模型在某些任务中取得了一定的成果, 但它们在如何有效结合不同技术范式、提升预测性能方面仍缺乏创新性。因此, 推动多范式协同的创新性框架成为提升该领域预测能力的关键。当前的算法大多侧重于单点突变的优化, 而在面对多点突变时, 表现较差。多点突变引起的协同效应对蛋白质热稳定性有着深远的影响, 而现有算法通常仅关注突变位点的局部特征, 忽略了突变之间的交互作用及其对稳定性的整体影响。因此, 如何在算法中有效地融合多点突变的协同效应, 仍然是一个待解的难题。此外, 尽管许多算法将蛋白质的原子坐标作为节点特征, 但这些算法往往未能充分考虑蛋白质结构的几何本质属性。蛋白质的三维结构不仅仅是原子之间的连接关系, 还包括旋转、平移等几何对称性, 这些几何性质对蛋白质的稳定性和折叠行为至关重要, 但现有算法忽视了这些因素, 导致其在构象变化的敏感性上存在局限性。为了解决这一问题, GDL作为一种新兴的架构, 能够显式地嵌入物理对称性先验, 提升模型对蛋白质几何约束(如旋转、平移和手性守恒等)的理解, 从而提高对构象变化的预测能力。通过将分子力场参数(如Lennard-Jones势能)作为几何网络的约束项, GDL可实现“几何-物理联合优化”, 将蛋白质的几何拓扑信息与热力学能量结合, 构建端到端的稳定性预测模型。这一方法有望突破现有模型的限制, 在稳定性预测中实现更高的准确性。

蛋白质热稳定性预测领域研究已取得许多重要进展, 但仍然存在许多进步空间。未来的研究应重点围绕数据集的扩展与质量提升, 开发更加精准的动态预测方法, 并通过多范式架构的协同创新推进

算法的发展。随着新的实验数据、计算模型的出现以及跨学科技术的融合, 蛋白质稳定性预测的精度和广泛性将不断提高, 进而推动在疾病研究、药物开发, 以及蛋白质工程等领域的应用。总之, 蛋白质热稳定性预测的研究正迈向新的突破, 数据的整合与增强、动态性结构的研究, 以及算法的多元化创新, 将是推动该领域发展的关键因素。

参 考 文 献

- [1] Jarzab A, Kurzawa N, Hopf T, *et al.* Meltome atlas-thermal proteome stability across the tree of life. *Nat Methods*, 2020, **17**(5): 495-503
- [2] Anfinsen C B. Principles that govern the folding of protein chains. *Science*, 1973, **181**(4096): 223-230
- [3] Pace C N. Measuring and increasing protein stability. *Trends Biotechnol*, 1990, **8**(4): 93-98
- [4] Tokuriki N, Stricher F, Serrano L, *et al.* How protein stability and new functions trade off. *PLoS Comput Biol*, 2008, **4**(2): e1000002
- [5] Ming Y, Wang W, Yin R, *et al.* A review of enzyme design in catalytic stability by artificial intelligence. *Brief Bioinform*, 2023, **24**(3): bbad065
- [6] Chen F, Chen X, Jiang F, *et al.* Computational analysis of androgen receptor (AR) variants to decipher the relationship between protein stability and related-diseases. *Sci Rep*, 2020, **10**: 12101
- [7] Mairbäurl H, Weber R E. Oxygen transport by hemoglobin. *Compr Physiol*, 2012, **2**(2): 1463-1489
- [8] Matthews C R, Crisanti M M, Gepner G L, *et al.* Effect of single amino acid substitutions on the thermal stability of the alpha subunit of tryptophan synthase. *Biochemistry*, 1980, **19**(7): 1290-1293
- [9] Sharma A, Boelens J J, Cancio M, *et al.* CRISPR-Cas9 editing of the *HBG1* and *HBG2* promoters to treat sickle cell disease. *N Engl J Med*, 2023, **389**(9): 820-832
- [10] Schymkowitz J, Borg J, Stricher F, *et al.* The FoldX web server: an online force field. *Nucleic Acids Res*, 2005, **33**(web server issue): W382-W388
- [11] Kellogg E H, Leaver-Fay A, Baker D. Role of conformational sampling in computing mutation-induced changes in protein structure and stability. *Proteins*, 2011, **79**(3): 830-838
- [12] Gapsys V, Michielssens S, Seeliger D, *et al.* Accurate and rigorous prediction of the changes in protein free energies in a large-scale mutation scan. *Angew Chem Int Ed*, 2016, **55**(26): 7364-7368
- [13] Pires DE, Ascher D B, Blundell T L. DUET: a server for predicting effects of mutations on protein stability using an integrated computational approach. *Nucleic Acids Res*, 2014, **42**(web server issue): W314-W319
- [14] Benevenuto S, Pancotti C, Fariselli P, *et al.* An antisymmetric neural network to predict free energy changes in protein variants. *J Phys D Appl Phys*, 2021, **54**(24): 245403
- [15] Chen Y, Lu H, Zhang N, *et al.* PremPS: predicting the impact of missense mutations on protein stability. *PLoS Comput Biol*, 2020,

- 16(12): e1008543
- [16] Isomura T, Kotani K, Jimbo Y, *et al.* Experimental validation of the free-energy principle with *in vitro* neural networks. *Nat Commun*, 2023, **14**(1): 4547
- [17] Deller M C, Kong L, Rupp B. Protein stability: a crystallographer's perspective. *Acta Crystallogr F Struct Biol Commun*, 2016, **72**(pt 2): 72-95
- [18] Johnson C M. Differential scanning calorimetry as a tool for protein folding and stability. *Arch Biochem Biophys*, 2013, **531**(1/2): 100-109
- [19] Yavşan E, Kalyoncu Uzunlar S. Fluorescence-based thermal stability screening is concentration-dependent and varies with protein size. *Front Life Sci Relat Technol*, 2023, **4**(2): 62-67
- [20] Giugliarelli A, Sassi P, Paolantoni M, *et al.* Heat-denatured lysozyme aggregation and gelation as revealed by combined dielectric relaxation spectroscopy and light scattering measurements. *J Phys Chem B*, 2012, **116**(35): 10779-10785
- [21] Miles A J, Wallace B A. Circular dichroism spectroscopy of membrane proteins. *Chem Soc Rev*, 2016, **45**(18): 4859-4872
- [22] Fiedler S, Cole L, Keller S. Automated circular dichroism spectroscopy for medium-throughput analysis of protein conformation. *Anal Chem*, 2013, **85**(3): 1868-1872
- [23] Guerois R, Nielsen J E, Serrano L. Predicting changes in the stability of proteins and protein complexes: a study of more than 1000 mutations. *J Mol Biol*, 2002, **320**(2): 369-387
- [24] Bednar D, Beerens K, Sebestova E, *et al.* FireProt: energy- and evolution-based computational design of thermostable multiple-point mutants. *PLoS Comput Biol*, 2015, **11**(11): e1004556
- [25] Song X, Wang Y, Shu Z, *et al.* Engineering a more thermostable blue light photo receptor *Bacillus subtilis* YtvA LOV domain by a computer aided rational design method. *PLoS Comput Biol*, 2013, **9**(7): e1003129
- [26] Sapozhnikov Y, Patel J S, Ytreberg F M, *et al.* Statistical modeling to quantify the uncertainty of FoldX-predicted protein folding and binding stability. *BMC Bioinformatics*, 2023, **24**(1): 426
- [27] Mitchell T M, Mitchell T M. *Machine Learning*. New York: McGraw-hill, 1997
- [28] Kotsiantis S B, Zaharakis I, Pintelas P. Supervised machine learning: a review of classification techniques. *Artificial Intelligence Review*, 2007, **160**(1): 3-24
- [29] Hu F, Xia G S, Wang Z, *et al.* Unsupervised feature learning via spectral clustering of multidimensional patches for remotely sensed scene classification. *IEEE J Sel Top Appl Earth Observations Remote Sensing*, 2015, **8**(5): 2015-2030
- [30] Sheikhpour R, Sarram M A, Gharaghani S, *et al.* A Survey on semi-supervised feature selection methods. *Pattern Recognit*, 2017, **64**: 141-158
- [31] Kiumarsi B, Vamvoudakis K G, Modares H, *et al.* Optimal and autonomous control using reinforcement learning: a survey. *IEEE Trans Neural Netw Learning Syst*, 2018, **29**(6): 2042-2062
- [32] Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition. *arXiv*, 2024. DOI: 10.48550/arXiv.1409.1556
- [33] Esteva A, Kuprel B, Novoa R A, *et al.* Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 2017, **542**(7639): 115-118
- [34] Göltz J, Kriener L, Baumbach A, *et al.* Fast and energy-efficient neuromorphic deep learning with first-spike times. *Nat Mach Intell*, 2021, **3**(9): 823-835
- [35] Guo Z, Liu J, Wang Y, *et al.* Diffusion models in bioinformatics and computational biology. *Nat Rev Bioeng*, 2024, **2**(2): 136-154
- [36] Bikman M. Using Transformers for Prediction of The Effect of Mutations [D]. Haifa: University of Haifa, 2024
- [37] Liu J, Guo Z, You H, *et al.* All-atom protein sequence design based on geometric deep learning. *Angew Chem Int Ed*, 2024, **63**(50): e202411461
- [38] Hopf T A, Ingraham J B, Poelwijk F J, *et al.* Mutation effects predicted from sequence co-variation. *Nat Biotechnol*, 2017, **35**(2): 128-135
- [39] Riesselman A J, Ingraham J B, Marks D S. Deep generative models of genetic variation capture the effects of mutations. *Nat Methods*, 2018, **15**(10): 816-822
- [40] Cheng J, Novati G, Pan J, *et al.* Accurate proteome-wide missense variant effect prediction with AlphaMissense. *Science*, 2023, **381**(6664): eadg7492
- [41] Yi M, Liu Y, Su Z. AlphaMissense, a groundbreaking advancement in artificial intelligence for predicting the effects of missense variants. *MedComm*, 2024, **3**(1): e70
- [42] Shanker V R, Bruun T U J, Hie B L, *et al.* Unsupervised evolution of protein and antibody complexes with a structure-informed language model. *Science*, 2024, **385**(6704): 46-53
- [43] Alzubaidi L, Zhang J, Humaidi A J, *et al.* Review of deep learning: concepts, CNN architectures, challenges, applications, future directions. *J Big Data*, 2021, **8**(1): 53
- [44] Cai H, Zhang H, Zhao D, *et al.* FP-GNN: a versatile deep learning architecture for enhanced molecular property prediction. *Brief Bioinform*, 2022, **23**(6): bbac408
- [45] Huang S, Cai N, Pacheco P P, *et al.* Applications of support vector machine (SVM) learning in cancer genomics. *Cancer Genomics Proteomics*, 2018, **15**(1): 41-51
- [46] Weiss K, Khoshgoftaar T M, Wang D. A survey of transfer learning. *J Big Data*, 2016, **3**(1): 9
- [47] Chowdhury R, Bouatta N, Biswas S, *et al.* Single-sequence protein structure prediction using a language model and deep learning. *Nat Biotechnol*, 2022, **40**(11): 1617-1623
- [48] Parmar A, Katariya R, Patel V. A review on random forest: an ensemble classifier//Hemanth J, Fernando X, Lafata P, *et al.* International Conference on Intelligent Data Communication Technologies and Internet of Things (ICICI) 2018. Cham: Springer International Publishing, 2019: 758-763
- [49] Bava K A, Gromiha M M, Uedaira H, *et al.* ProTherm, version 4.0: thermodynamic database for proteins and mutants. *Nucleic Acids Res*, 2004, **32**(database issue): D120-D121
- [50] Capriotti E, Fariselli P, Casadio R. A neural-network-based method for predicting protein stability changes upon single point mutations. *Bioinformatics*, 2004, **20**(Suppl 1): i63-i68

- [51] Frenz C M. Neural network-based prediction of mutation-induced protein stability changes in Staphylococcal nuclease at 20 residue positions. *Proteins*, 2005, **59**(2): 147-151
- [52] Capriotti E, Fariselli P, Calabrese R, *et al.* Predicting protein stability changes from sequences using support vector machines. *Bioinformatics*, 2005, **21**(Suppl 2): ii54-ii58
- [53] Berliner N, Teyra J, Colak R, *et al.* Combining structural modeling with ensemble machine learning to accurately predict protein fold stability and binding affinity effects upon mutation. *PLoS One*, 2014, **9**(9): e107353
- [54] Pucci F, Schwersensky M, Rooman M. Artificial intelligence challenges for predicting the impact of mutations on protein stability. *Curr Opin Struct Biol*, 2022, **72**: 161-168
- [55] Yang Y, Urolagin S, Niroula A, *et al.* PON-tstab: protein variant stability predictor. Importance of training data quality. *Int J Mol Sci*, 2018, **19**(4): E1009
- [56] Pucci F, Bernaerts K V, Kwasigroch J M, *et al.* Quantification of biases in predictions of protein stability changes upon mutations. *Bioinformatics*, 2018, **34**(21): 3659-3665
- [57] Xavier J S, Nguyen T B, Karmarkar M, *et al.* ThermoMutDB: a thermodynamic database for missense mutations. *Nucleic Acids Res*, 2021, **49**(d1): D475-D479
- [58] Pancotti C, Benevenuta S, Birolo G, *et al.* Predicting protein stability changes upon single-point mutation: a thorough comparison of the available tools on a new dataset. *Brief Bioinform*, 2022, **23**(2): bbab555
- [59] Fang J. A critical review of five machine learning-based algorithms for predicting protein stability changes upon mutation. *Brief Bioinform*, 2020, **21**(4): 1285-1292
- [60] Samaga Y B L, Raghunathan S, Priyakumar U D. SCONES: self-consistent neural network for protein stability prediction upon mutation. *J Phys Chem B*, 2021, **125**(38): 10657-10671
- [61] Cao H, Wang J, He L, *et al.* DeepDDG: predicting the stability change of protein point mutations using neural networks. *J Chem InfModel*, 2019, **59**(4): 1508-1514
- [62] Wang S, Tang H, Shan P, *et al.* ProS-GNN: predicting effects of mutations on protein stability using graph neural networks. *Comput Biol Chem*, 2023, **107**: 107952
- [63] Pak M A, Dovidchenko N V, Sharma S M, *et al.* New mega dataset combined with deep neural network makes a progress in predicting impact of mutation on protein stability. *BioRxiv*, 2023. DOI: 10.1101/2022.12.31.522396
- [64] Sun J, Zhu T, Cui Y, *et al.* Structure-based self-supervised learning enables ultrafast protein stability prediction upon mutation. *Innovation (Camb)*, 2025, **6**(1): 100750
- [65] Li B, Yang Y T, Capra J A, *et al.* Predicting changes in protein thermodynamic stability upon point mutation with deep 3D convolutional neural networks. *PLoS Comput Biol*, 2020, **16**(11): e1008291
- [66] Li S S, Liu Z M, Li J, *et al.* Prediction of mutation-induced protein stability changes based on the geometric representations learned by a self-supervised method. *BMC Bioinformatics*, 2024, **25**(1): 282
- [67] Liu B, Jiang Y, Yang Y, *et al.* OmeDDG: improved protein mutation stability prediction based on predicted 3D structures. *J Phys Chem B*, 2024, **128**(1): 67-76
- [68] Iqbal S, Ge F, Li F, *et al.* PROST: AlphaFold2-aware sequence-based predictor to estimate protein stability changes upon missense mutations. *J Chem InfModel*, 2022, **62**(17): 4270-4282
- [69] Chen Y, Xu Y, Liu D, *et al.* An end-to-end framework for the prediction of protein structure and fitness from single sequence. *Nat Commun*, 2024, **15**(1): 7400
- [70] Wang S, Tang H, Zhao Y, *et al.* BayeStab: predicting effects of mutations on protein stability with uncertainty quantification. *Protein Sci*, 2022, **31**(11): e4467
- [71] Mishra S K. PSP-GNM: predicting protein stability changes upon point mutations with a Gaussian network model. *Int J Mol Sci*, 2022, **23**(18): 10711
- [72] Dieckhaus H, Brocchiacono M, Randolph N Z, *et al.* Transfer learning to leverage larger datasets for improved prediction of protein stability changes. *Proc Natl Acad Sci USA*, 2024, **121**(6): e2314853121
- [73] Gong J, Jiang L, Chen Y, *et al.* THPLM: a sequence-based deep learning framework for protein stability changes prediction upon point variations using pretrained protein language model. *Bioinformatics*, 2023, **39**(11): btad646
- [74] Diaz D J, Gong C, Ouyang-Zhang J, *et al.* Stability Oracle: a structure-based graph-transformer framework for identifying stabilizing mutations. *Nat Commun*, 2024, **15**(1): 6170
- [75] Li G, Yao S, Fan L. ProSTAGE: predicting effects of mutations on protein stability by using protein embeddings and graph convolutional networks. *J Chem InfModel*, 2024, **64**(2): 340-347
- [76] Reeves S, Kalyaanamoorthy S. Zero-shot transfer of protein sequence likelihood models to thermostability prediction. *Nat Mach Intell*, 2024, **6**(9): 1063-1076
- [77] Zhou Y, Pan Q, Pires D E V, *et al.* DDMut: predicting effects of mutations on protein stability using deep learning. *Nucleic Acids Res*, 2023, **51**(w1): W122-W128
- [78] Cheng K, Liu C, Su Q, *et al.* 4D diffusion for dynamic protein structure prediction with reference guided motion alignment. *arXiv*. 2024. DOI: 10.48550/arXiv.2408.12419

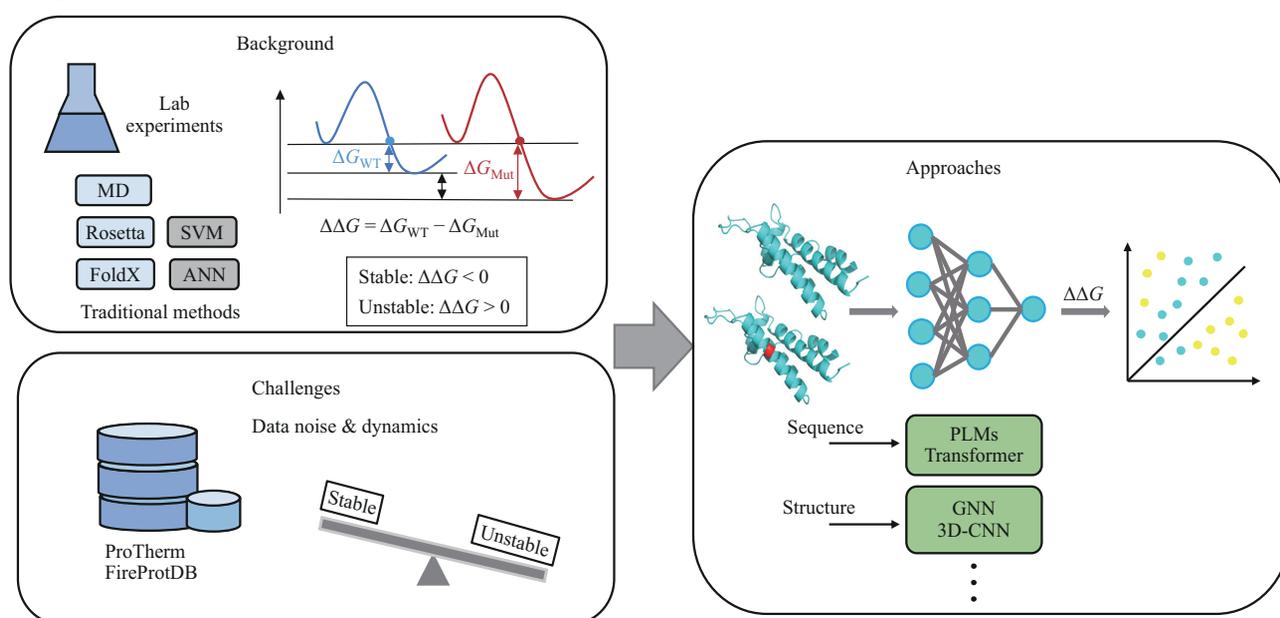
Prediction of Protein Thermodynamic Stability Based on Artificial Intelligence*

TAO Lin-Jie¹⁾, XU Fan-Ding¹⁾, GUO Yu²⁾, LONG Jian-Gang^{1)**}, LU Zhuo-Yang^{1)**}

¹⁾Institute of Mitochondrial Biomedicine, School of Life Sciences and Technology, Xi'an Jiaotong University, Xi'an 710049, China;

²⁾National Key Laboratory of Human-Machine Hybrid Augmented Intelligence, National Engineering Research Center for Visual Information and Applications, Institute of Artificial Intelligence and Robotics, Xi'an Jiaotong University, Xi'an 710049, China)

Graphical abstract



Abstract In recent years, the application of artificial intelligence (AI) in the field of biology has witnessed remarkable advancements. Among these, the most notable achievements have emerged in the domain of protein structure prediction and design, with AlphaFold and related innovations earning the 2024 Nobel Prize in Chemistry. These breakthroughs have transformed our ability to understand protein folding and molecular interactions, marking a pivotal milestone in computational biology. Looking ahead, it is foreseeable that the accurate prediction of various physicochemical properties of proteins—beyond static structure—will become the next critical frontier in this rapidly evolving field. One of the most important protein properties is thermodynamic stability, which refers to a protein's ability to maintain its native conformation under physiological or stress conditions. Accurate prediction of protein stability, especially upon single-point mutations, plays a vital role in numerous scientific and industrial domains. These include understanding the molecular basis of disease, rational

* This work was supported by a grant from The National Natural Science Foundation of China (32271281).

** Corresponding author.

LU Zhuo-Yang. Tel: 86-29-82665849, E-mail: luzhuoyang@xjtu.edu.cn

LONG Jian-Gang. Tel: 86-29-82665849, E-mail: jglong@xjtu.edu.cn

Received: December 25, 2024 Accepted: June 4, 2025

drug design, development of therapeutic proteins, design of more robust industrial enzymes, and engineering of biosensors. Consequently, the ability to reliably forecast the stability changes caused by mutations has broad and transformative implications across biomedical and biotechnological applications. Historically, protein stability was assessed via experimental methods such as differential scanning calorimetry (DSC) and circular dichroism (CD), which, while precise, are time-consuming and resource-intensive. This prompted the development of computational approaches, including empirical energy functions and physics-based simulations. However, these traditional models often fall short in capturing the complex, high-dimensional nature of protein conformational landscapes and mutational effects. Recent advances in machine learning (ML) have significantly improved predictive performance in this area. Early ML models used handcrafted features derived from sequence and structure, whereas modern deep learning models leverage massive datasets and learn representations directly from data. Deep neural networks (DNNs), graph neural networks (GNNs), and attention-based architectures such as transformers have shown particular promise. GNNs, in particular, excel at modeling spatial and topological relationships in molecular structures, making them well-suited for protein modeling tasks. Furthermore, attention mechanisms enable models to dynamically weigh the contribution of specific residues or regions, capturing long-range interactions and allosteric effects. Nevertheless, several key challenges remain. These include the imbalance and scarcity of high-quality experimental datasets, particularly for rare or functionally significant mutations, which can lead to biased or overfitted models. Additionally, the inherently dynamic nature of proteins—their conformational flexibility and context-dependent behavior—is difficult to encode in static structural representations. Current models often rely on a single structure or average conformation, which may overlook important aspects of stability modulation. Efforts are ongoing to incorporate multi-conformational ensembles, molecular dynamics simulations, and physics-informed learning frameworks into predictive models. This paper presents a comprehensive review of the evolution of protein thermodynamic stability prediction techniques, with emphasis on the recent progress enabled by machine learning. It highlights representative datasets, modeling strategies, evaluation benchmarks, and the integration of structural and biochemical features. The aim is to provide researchers with a structured and up-to-date reference, guiding the development of more robust, generalizable, and interpretable models for predicting protein stability changes upon mutation. As the field moves forward, the synergy between data-driven AI methods and domain-specific biological knowledge will be key to unlocking deeper understanding and broader applications of protein engineering.

Key words machine learning, protein thermodynamic stability, mutations

DOI: 10.16476/j.pibb.2024.0530

CSTR: 32369.14.pibb.20240530