

# 理论多相催化计算中的机器学习方法: 现状与挑战

刘照清<sup>†</sup>, 邓哲<sup>†</sup>, 蒋鸿<sup>\*</sup>

北京大学化学与分子工程学院, 北京 100871

<sup>†</sup> 同等贡献

\* 联系人, E-mail: [jianghchem@pku.edu.cn](mailto:jianghchem@pku.edu.cn)

2024-11-11 收稿, 2025-02-11 修回, 2025-02-12 接受, 2025-02-13 网络版发表

国家重点研发计划(2022YFB4101401)资助

**摘要** 多相催化反应在化工领域具有重要地位, 是合成氨、煤制油等工业部门的核心。一直以来, 为了更好地理解活性位点、反应路径以及微观动力学等信息, 基于密度泛函理论的计算方法在这一领域发挥着独特的作用。但受限于计算能力, 这样的研究手段难以处理涉及大时间与空间尺度的问题。近年来, 作为人工智能重要分支的机器学习方法逐渐兴起, 为理论多相催化研究带来了全新的研究范式。本文首先介绍了针对特定催化问题的统计学习模型、应用策略与局限性; 然后重点探讨了机器学习势函数及其在多相催化理论模拟中应用的现状与挑战, 并强调了近期兴起的通用机器学习势模型及其预训练-微调部署方式在多相催化模拟中的潜力; 最后讨论了新兴的包括大语言模型在内的生成式模型在多相催化研究中的前景。在可预期的未来, 机器学习方法与多相催化将有更深入的交叉, 为多相催化剂理性设计工作带来更多便利, 这既是相关研究者的一大机遇, 也是新的挑战。

**关键词** 理论多相催化, 催化剂理性设计, 机器学习, 机器学习势, 生成式模型

作为一门加速化学反应、调控反应产物分布的科学, 催化已经与现代人类生活紧密地连接在一起。在具体研究中, 催化一般被分为均相催化与多相催化两大分支, 其中多相催化是指在不同相交界处发生的催化反应过程, 定义上包括了固-液、固-气、气-液等情形<sup>[1]</sup>。考虑到工业生产实践中通常将流动相气液组分通过固态催化剂以完成反应, 本文主要讨论的是发生在固体表界面的反应。过去一个世纪以来, 多相催化反应构成了化学工业必不可少的一部分, 但随之而来的也是巨大的资源消耗和碳排放<sup>[2]</sup>。有报道表明, 当前全球约2%的能源消耗在了合成氨相关的工业部门<sup>[3]</sup>。若能根据理论预测针对特定反应设计出高效的催化剂, 不仅可以提高反应速率与产物选择性, 还能降低反应的能耗, 以达到经济发展与环境保护双赢的目标。近些年, 更加系统的催化剂理性设计策略受到了学界的关

注<sup>[4]</sup>, 由此对催化剂具体作用机制与构效关系认识提出了较高的要求。在这样的形势下, 基于经验试错的传统催化研究范式已经难以适用于新时代的研究。

具体来看, 除了将各类表征技术综合应用的实验手段外, 多相催化研究中理论计算的贡献也是不容忽视的<sup>[5]</sup>: 其通过自底而上的跨尺度模拟, 从催化剂材料、表面与催化反应的建模开始, 再基于密度泛函理论(density functional theory, DFT)计算获取表面反应的热力学数据(包括基元反应能变、活化能垒以及气体组分的热力学信息), 随后通过微观动力学求解以从理论上预测催化反应的反应机制与决速步, 最后结合动力学信息与催化剂性质相关描述符, 即可在理论指导下更好地实现催化剂理性设计。以在合成气转化中发挥关键作用的费托合成反应为例, 通过基于第一性原理的DFT计算可以确定反应的大量信息, 具体结果包

引用格式: 刘照清, 邓哲, 蒋鸿. 理论多相催化计算中的机器学习方法: 现状与挑战. 科学通报, 2025, 70: 4081–4097

Liu Z-Q, Deng Z, Jiang H. Machine learning methods for theoretical heterogeneous catalysis: current status and challenges (in Chinese). Chin Sci Bull, 2025, 70: 4081–4097, doi: [10.1360/TB-2024-1207](https://doi.org/10.1360/TB-2024-1207)

括活性相、活性表面<sup>[6]</sup>与表面反应位点<sup>[7]</sup>与反应动力学数据<sup>[8]</sup>。在过去20年间,由于理论计算方法的发展和计算机软硬件水平的快速提高,相关的理论计算已经广泛应用于新型多相催化剂的设计与优化工作中。但由于密度泛函理论相对高的计算成本,目前仍然难以将这样的计算推广到更加复杂从而更接近真实反应条件的多相催化体系。在空间尺度方面,部分多相催化反应在实验上使用小粒径的纳米粒子催化剂<sup>[9]</sup>,其结构特征是简单的周期性平板(slab)模型不易描述的。若是针对纳米粒子建模,通常情况下直径为数纳米的粒子具有数千甚至上万的原子数,常规的第一性原理计算手段很难处理。在时间尺度方面,多相催化中涉及分子动力学的任务有时依赖于长时间(通常达到纳秒甚至更长的尺度)的采样<sup>[10]</sup>,其在百原子级别的体系也超过了第一性原理方法的计算能力。为了解决类似的问题,理论多相催化领域迫切需要一种计算量小但仍保持较高精度的研究策略:或是直接通过物理描述符预测目标性质,或是通过高精度力场代替第一性原理方法进行采样计算。这就引出了本文讨论的机器学习方法。

机器学习(machine learning, ML)这一概念的含义相当广泛,其涵盖范围从经典的线性回归模型,到21世纪才逐步兴起的深度神经网络<sup>[11]</sup>,更包括近年来异常火爆的生成式模型和大语言模型。这一概念从计算机视觉和自然语言处理等计算机领域兴起,并在近年来与自然科学形成了许多交叉,形成了AI for Science的研究新范式<sup>[12]</sup>。2024年的诺贝尔物理学奖与化学奖都颁发给了相关研究工作,充分表明机器学习方法在物质科学基础和应用研究的各个领域都大有可为。

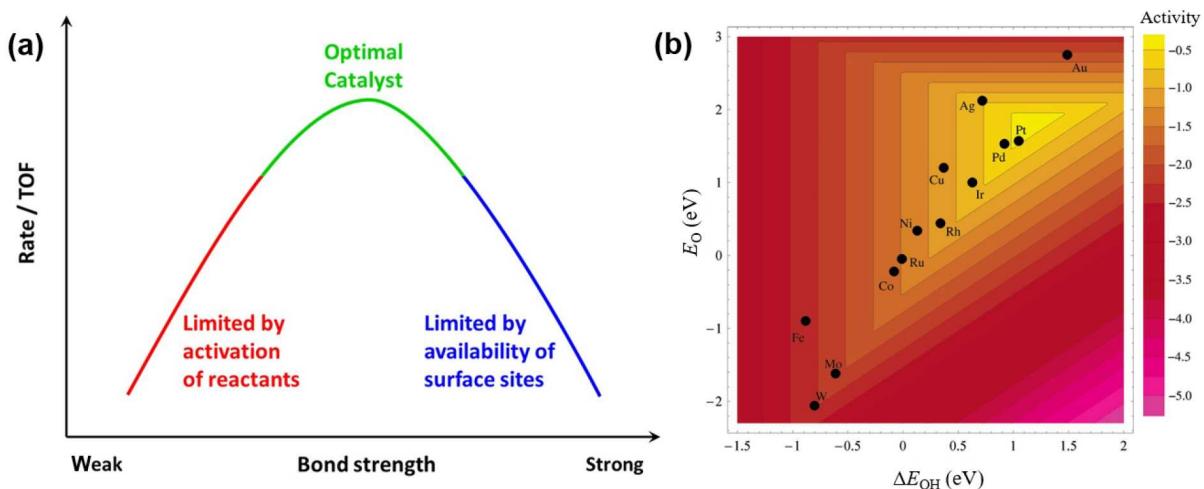
基于以上背景,本文对近年来理论多相催化领域内的机器学习方法研究进展进行了总结。首先讨论了针对特定催化性质,基于研究者自行构建的描述符形成的小机器学习模型。这类模型通常在训练中仅需要较小的数据量,一般具有较好的可解释性,但不易做外推应用。随后介绍了在多相催化模拟领域广泛应用的机器学习势。这类方法通过机器学习的方式表示能量对结构的依赖关系,兼顾了经典力场的计算效率和DFT的计算精度,大大拓展了力场方法的应用边界。我们着重介绍了近年来新兴的通用机器学习势,它打破了传统机器学习力场对特定元素种类的依赖,可应用于元素种类和结构类型非常多样的场合,且可通过预训练-微调的方式,实现针对特定问题的快速部署。对于人工智能领域新兴的包括大语言模型在内的生成式

模型,我们结合最新工作讨论了其在多相催化领域的应用现状与特点,以及这一方向的发展趋势。

## 1 针对特定催化问题的统计学习模型

在现代意义上的机器学习方法出现之前,多相催化研究就已经开始尝试基于物理描述符建立不同量之间的关系了。1912年,诺贝尔奖得主Paul Sabatier提出了著名的Sabatier原理:理想的催化剂应该具有适中的吸附能力,既能使反应物在活性位点上被有效吸附和活化,又能让产物及时脱附以重新释放活性位点<sup>[13]</sup>。这一论断相当于建立了物种吸附能与催化反应活性之间的关系,在实际研究中常常通过图1(a)所示的火山形曲线(volcano curve)来表达<sup>[14]</sup>。而在火山形曲线的顶点,也即反应物种吸附能适中时,催化剂具有最高的总体反应性能。Sabatier原理在理论多相催化领域应用广泛,直到近年还发挥着不可或缺的作用。如2023年, Jin等人<sup>[15]</sup>研究了不同浓度的Ir单原子催化剂对芳香化合物选择性加氢反应的催化活性,发现这一活性在中等浓度(0.7 atoms/nm<sup>2</sup>)时达到峰值,进一步的机理研究表明此时表面H<sup>\*</sup>的吸附与解吸强度均处于适中水平。另外, Sabatier原理还可以扩展到多个描述符的情形。例如, Nørskov等人<sup>[16]</sup>曾经研究了电催化氧还原反应(oxygen reduction reaction, ORR)在不同电极材料上的活性(图1(b)),并使用O与OH两种基团在金属表面的吸附能作为描述符以预测总反应活性。在其计算的十多种金属中,最接近火山图顶点的为Pt,由此从理论层面说明了为什么Pt是ORR过程最常用的电极材料之一。

Sabatier原理的背后实际上是线性标度关系(linear scaling relationship),在具体应用中一般包括线性自由能关系以及催化表面性质的描述符与其他物理量的线性关系等<sup>[17]</sup>。这实际上启发了更普适的线性模型,也即多元线性回归拟合的引入。这一研究方法由于所需数据量与计算量较小,且能方便地反映结果对各描述符的敏感程度,在领域内得到了较多的应用。如Schmack等人<sup>[18]</sup>研究了甲烷氧化偶联(oxidative coupling of methane, OCM)反应,从一系列文献数据中提取出催化剂理化性质描述符与反应性之间的联系。Hong等人<sup>[19]</sup>建立了一批钙钛矿氧化物(ABO<sub>3</sub>)中几何和电子结构描述符与析氧反应的线性关系,并发现其中最关键的指标是B位元素d电子数目与电荷转移能。不仅如此,现有的多元线性模型还可以进一步扩展,以引入不同描述符间的非线性关联。考虑到非线性函数组合的复杂



**图 1** (网络版彩色)用于描述Sabatier原理的火山形曲线. (a) 火山形曲线的一般形式, 在中等吸附强度时相应的反应速率达到最大<sup>[14]</sup>, Copyright © 2021 American Chemical Society; (b) ORR过程的三维火山形曲线, 横纵坐标分别代表OH与O两种基团在金属上的吸附能. 图中所示金属中Pt最为接近活性顶点<sup>[16]</sup>, Copyright © 2004 American Chemical Society

**Figure 1** (Color online) Volcano curves about the Sabatier principle. (a) The general form of the volcano curve, where the reaction rate reaches its maximum at moderate adsorption strength<sup>[14]</sup>, Copyright © 2021 American Chemical Society; (b) 3D volcano curve of the ORR process, with the horizontal and vertical axes representing the adsorption energies of OH and O species on metals, respectively. Among these metals shown in this figure, platinum is closest to the activity peak<sup>[16]</sup>, Copyright © 2004 American Chemical Society

程度, 我们可以根据领域知识先验地导出针对特定问题的物理方程, 再对已有数据拟合以得到方程中的未知系数<sup>[20]</sup>. 但在复杂的多相催化场合, 这样的物理方程很难预先得到, 如何通过自动化的手段完成类似的工作? 这方面比较有代表性的成果是欧阳润海等人<sup>[21]</sup>在符号回归机器学习框架中提出的SISSO (sure independence screening and sparsifying operator)方法. 其首先通过经验或具体领域知识构建出一系列较为简单的初始描述符. 随后, 程序将通过特定的规则将初始描述符与一系列运算符( $+$ ,  $-$ ,  $\times$ ,  $\div$ ,  $\sqrt{\cdot}$ ,  $\exp$ ,  $\log$ ,  $| \cdot |$ ,  $\cdots \cdots$ )组合起来, 形成更加多样化的描述符, 然后基于LASSO回归(利用稀疏性)以得到具体拟合表达式. 如果待研究的问题足够复杂, 该步可能会组合出百亿量级甚至更多的候选复合描述符. 最后, 程序将按照拟合效果与模型复杂度输出一系列可能的拟合表达式与相应的关键描述符, 这些表达式在拟合精度上差别较小, 而研究者可以根据自己对问题的理解选出其中最符合要求的.

以Shu等人<sup>[22]</sup>研究金属催化剂结构敏感性的工作为例, 作者主要选取了一氧化碳与部分醇类分子在过渡金属表面的解离反应为训练集, 尝试应用简单的几何描述符结合反应能量变化信息预测反应能垒. 如图2所示, 应用SISSO方法, 他们成功获得了具有明确物理意义的拟合式, 结果可视为经典的BEP关系表达式中增

加了表面物种配位数的贡献项, 并达到了较好的预测精度和可迁移性. 另外, 由于SISSO方法在物理可解释性上的优势, 这一方法还可以用于从复杂化学体系提取简单明确的理论模型. 如李微雪及其合作者<sup>[23]</sup>讨论了金属催化剂与氧化物载体的相互作用, 并针对显著影响催化剂活性和选择性的包覆(encapsulation)现象进行了详细研究. 通过拟合金属-氧化物界面上黏附能(adhesion energy)、金属-金属相互作用与金属-氧相互作用间的关系, 作者发现包覆现象的出现主要由强的金属-金属相互作用导致, 而不仅仅取决于金属对氧的亲和性, 从而提供了一个全新的研究视角. 需要指出的是, 在实际应用中, 这一方法不总是能够得到物理意义足够清晰的表达式, 特别是所计算的体系相对复杂时. Zhou等人<sup>[24]</sup>在研究金属助剂作用下铂催化剂表面发生的丙烷脱氢反应时, 应用SISSO方法得到的拟合式中出现了诸如电子亲和能与反应生成热的组合、电负性与原子半径的乘积等物理意义不够明确的项. 在这种情况下, 虽然仍然可以通过拟合表达式进行后续的催化剂优化, 但模型的物理可解释性受到了一定的限制. 值得注意的是, 通过恰当选取拟合中使用的描述符种类, 可以在一定程度上缓解这一问题的影响. 如Wang等人<sup>[25]</sup>提出在SISSO方法中使用谱学(如振动光谱、X射线吸收谱等)相关描述符, 所得结果往往可以将催化剂

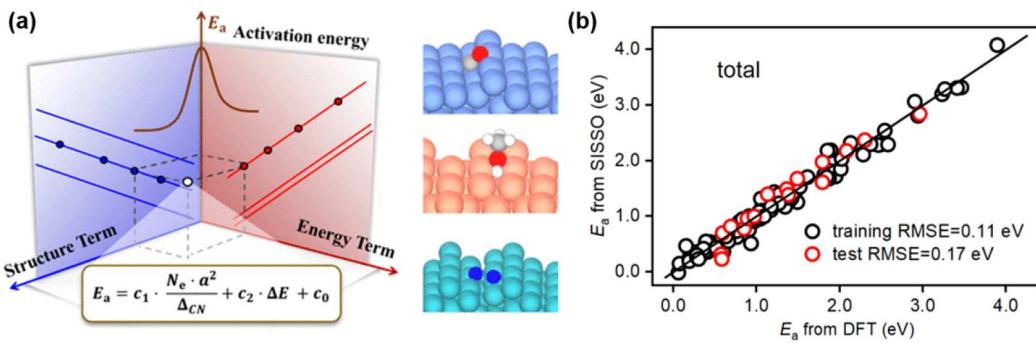


图 2 (网络版彩色)过渡金属表面小分子解离反应的SISSO预测.(a)拟合所得表达式,其中第一项与表面物种配位数有关,第二项对应BEP关系中反应能量<sup>[22]</sup>. (b)所得模型在训练集与测试集中的表现<sup>[22]</sup>. Copyright © 2024 American Chemical Society

**Figure 2** (Color online) SISSO results of small molecule dissociation reactions on transition metal surfaces. (a) The fitted formula, where the first term corresponds to the coordination number of surface species and the second term corresponds to reaction energy in the BEP relationship<sup>[22]</sup>. (b) Model performance on the training set and the test set<sup>[22]</sup>. Copyright © 2024 American Chemical Society

基底与吸附物的贡献分开,且所得模型具有较好的可迁移性,这可能与谱学描述符能够充分捕获催化体系的物理化学信息有关。从这一点来说,这类方法对研究者在描述符构造等环节中的领域知识提出了较高的要求。

视任务特点,其他在机器学习领域常见的模型也能应用于理论多相催化计算<sup>[26]</sup>:如无监督学习中的主成分分析、聚类;监督学习中的决策树、随机森林、高斯过程回归等(表1)。由于这类针对特定多相催化体系的统计模型主要基于人为构造的物理描述符,通常具有相对较好的可解释性,方便根据已有模型对催化剂开展进一步的优化。但这类统计模型的一个局限性是由于参数量较小且强烈依赖于具体描述符,通常难以迁移推广到离原始训练集较远的体系,因此在新的问题中往往需要重新训练。

理论多相催化研究的一大目标就是在实验上指导新催化剂的理性设计,目前这类统计模型已经广泛地应用于评价催化剂活性,从而减少后期实验筛选的工作量<sup>[11]</sup>。不仅如此,实际研究中还可以利用实验结果进一步优化原有模型,由此形成模型-实验-模型的回路以得到更好的表现。一个典型的例子来自江俊等人<sup>[34,35]</sup>近年的工作,他们将高通量DFT计算、机器学习模型与自动化实验数据充分结合,并通过自动化机器人实现了高效的催化剂开发。如基于火星陨石的产氧催化剂设计<sup>[35]</sup>,任务要求找到具有高析氧反应活性(由过电位表征)的化合物并完成实验。基于领域知识,作者选取了3个析氧反应中的物理描述符,其中两个与羟基和物种的吸附自由能有关、一个与活性位点吸附羟基前

后的电荷变化量有关。由于通过DFT计算这些描述符较为耗时,他们基于两个较小的两层神经网络,拟合了催化剂不同金属组分含量到物理描述符再到反应过电位间的关系,这一方案的好处是可以分别使用理论计算及实验结果对所得模型进行校准。最后,基于贝叶斯优化能够确定最优的催化剂组成。这样的流程可以反复进行,以主动学习新的实验和计算数据。

值得说明的是,本部分讨论的统计模型更多关注于建立不同描述符之间的联系。但在实际研究中,要对一个足够复杂的多相催化体系完成建模并计算特定物理量本身就是一个复杂的问题。因此,相当一部分研究还关注于如何应用机器学习方法提高模拟计算的效率,从而以远超第一性原理方法的计算速度实现相似的势能面精度,这也是下文重点讨论的对象。

## 2 针对复杂多相催化体系的机器学习势方法

机器学习势(machine learning potential, MLP),又称机器学习力场(machine learning force field, MLFF)或机器学习原子间势(machine learning interatomic potential, MLIP),是一类通过机器学习模型拟合第一性原理势能面,构造兼具第一性原理精度和经典经验性力场计算效率的机器学习势能面求解器的方法。势能面(potential energy surface, PES)是从原子尺度理解真实体系的核心,它包含了有关体系稳态与亚稳态结构与能量关系、有限温度下驱动化学动力学的原子间作用力,以及结构转变与化学反应之间的过渡态与能垒等重要化学物理信息<sup>[36]</sup>。在机器学习和第一性原理方法发展成熟之前,研究者们就在尝试以数学物理模型拟合实验数据

**表 1** 常见机器学习模型在理论多相催化中的应用举例**Table 1** Applications of common machine learning models in theoretical heterogeneous catalysis

应用体系	模型名称	模型作用	来源文献
二氧化碳加氢制甲醇	主成分分析	衡量催化描述符贡献大小	[27]
甲烷氧化偶联反应	K-means聚类	筛选可能具有高活性的催化剂	[28]
铑基合成气转化反应	高斯过程回归	预测反应中间体生成自由能	[29]
铂基水煤气变换反应	支持向量机	预测一氧化碳转化率	[30]
沸石类化合物选择性催化还原NO <sub>x</sub>	决策树	提出高效催化剂的合成策略	[31]
掺杂Ni <sub>2</sub> P催化剂的析氢反应	随机森林	预测热力学能并衡量描述符重要性	[32]
金属-沸石催化剂的CO <sub>2</sub> 还原反应	XGBoost	预测热力学能与产物选择性数据	[33]

来描述势能面信息，搭建了大量的经典经验力场模型，包括以Amber<sup>[37]</sup>、CHARMM<sup>[38]</sup>为代表的分子力学力场和以ReaxFF<sup>[39]</sup>为代表的反应性力场。基于机器学习方法尤其是人工神经网络(*artificial neural network*, ANN)拟合势能面的尝试早在20世纪90年代就有一定的开展<sup>[40,41]</sup>。而在近10年来ResNet<sup>[42]</sup>、Transformer<sup>[43]</sup>等深度神经网络(*deep neural network*, DNN)及其残差连接算法、注意力机制等机器学习算法和相关计算机算力迅猛发展的大背景下，机器学习势及其相关方法也随之不断更新换代，并广泛应用于多相催化理论模拟工作当中。它们显著地扩展了第一性原理精度计算模拟的空间尺度和时间尺度，能实现此前许多应用传统方法难以完成的催化模拟任务，使得研究者能够对各类催化问题形成全新的理解。

基于在拟合势能面任务中使用的不同机器学习模型架构，机器学习势主要分为基于核方法的机器学习势和基于人工神经网络的机器学习势两类。基于核方法的机器学习势利用核函数将原子结构在高维空间中表示，通过这种表示和特定的高斯回归策略对化学结构开展模式识别，构建特定结构模式到对应能量的映射，从而拟合势能面。这一方法包括但不限于高斯近似势(*Gaussian approximation potential*, GAP)<sup>[44]</sup>、谱邻域分析势(*spectral neighbor analysis potential*, SNAP)<sup>[45]</sup>，以及梯度域机器学习(*gradient domain machine learning*, GDML)<sup>[46]</sup>方法及其对称化改进<sup>[47]</sup>等。这些基于核方法的机器学习势具有相对较小的参数量和相对清晰的数学表达，能在几百到数千帧的少量训练集上充分训练，并在测试集上表现出较高的精度，因而适用于特定小规模问题的机器学习势构建，并在这些问题上体现出较强的表示能力。然而，其模型训练与推理成本随训练集规模增长的趋势在O( $N^2$ )~O( $N^3$ )量级，难以在针对较大规模问题的训练集上训练，可扩展性相对有限，缺乏

描述复杂化学体系的能力。与之对应的，基于人工神经网络的机器学习势，又称神经网络势(*neural network potentials*, NNPs)，则是通过神经网络来表示体系中各原子能量贡献对局域化学环境的依赖性。这一方法名列国际纯粹与应用化学联合会(IUPAC)认定的2024年度化学领域十大新兴技术榜单，它在相关概念和方法上具有最广的多样性，是目前发展最迅速且应用最广泛的机器学习势方法。相较于核方法，神经网络方法往往需要更多的训练数据量才能表现相似的精度，但这一方法具备更好的可扩展性，能在几千至上万帧的数据下训练，并获得对更广泛的化学空间的表示和推理能力<sup>[48]</sup>。

神经网络势的开创性工作是Behler和Parrinello于2007年提出的BPNN势，及其对应的高维神经网络势(*high-dimensional neural network potentials*, HDNNPs)架构<sup>[49]</sup>。这一模型架构如图3(a)所示<sup>[50]</sup>，包含输入层、隐藏层和输出层三部分。其输入层为由一系列以原子坐标为输入的、短程截断的代数表达式构成的描述符，它们被称作原子中心对称函数(*atom-centered symmetry functions*, ACSFs)<sup>[51]</sup>，由以体系各个原子为中心的、考虑其与周围原子间几何关系的径向函数(两体项)、角向函数(三体项)等类似于力场描述符的多体函数组合而成，能解析地表示体系各个原子的局域化学环境特征。这些特征经过后续隐藏层中的原子神经网络处理后，在最终的输出层输出各个原子对体系总能量的贡献(又称为原子能量)，并加和为体系的(短程的)总能量。为确保模型对体系原子置换的不变性，并保证模型支持任意原子数目的体系，模型中同一元素的各个原子共享同一套原子神经网络。该机器学习势的关键即在其描述符ACSFs，它的解析形式使体系的原子受力(即势能面梯度)和晶格应力可以解析求解，满足势能面应有的保守性；同时其表达式也能满足物理体系应有的

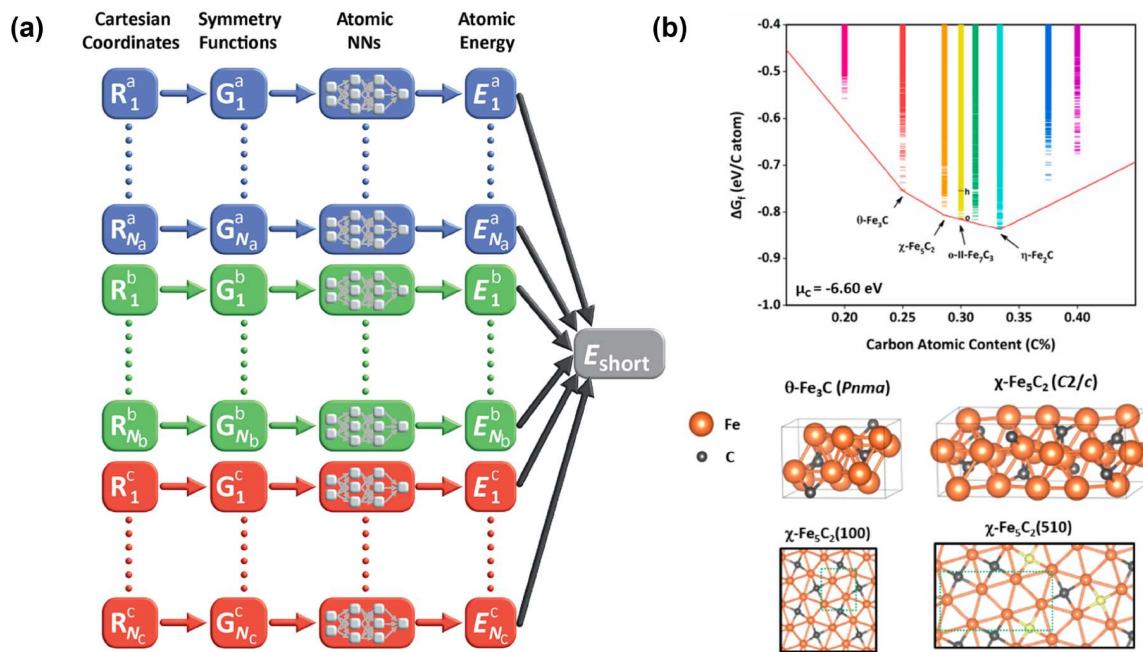


图 3 (网络版彩色)高维神经网络势架构及其在多相催化模拟中的应用. (a) 高维神经网络势架构示意, 图中三种颜色代表三种不同的元素类型<sup>[50]</sup>, Copyright © 2021 American Chemical Society; (b) 基于SSW-NN方法搜索得到的费托合成气氛下常见铁碳化合物物相的热力学凸包图, 以及它们的最概然体相结构和表面结构<sup>[56]</sup>, Copyright © 2021 American Chemical Society

**Figure 3** (Color online) HDNNPs and their applications in heterogeneous catalysis simulations. (a) Architecture of HDNNPs, where the three colors represent three different types of elements<sup>[50]</sup>, Copyright © 2021 American Chemical Society; (b) convex-hull and most probable structures for the bulk and surface of iron carbide under gas-phase environments of Fischer-Tropsch synthesis searched by the SSW-NN method<sup>[56]</sup>, Copyright © 2021 American Chemical Society

对称性, 即对体系平移、旋转和原子置换的不变性. 这两个特点确保了该机器学习势的表示能力和开展稳定可靠的动力学模拟的能力. 在这一架构设计的基础上, 刘智攀及其合作者<sup>[52,53]</sup>进一步将ACSFs描述符从指数形式改进为结合球谐函数的多项式形式, 并将描述符扩展到四体项, 构建高维神经网络势. 这一机器学习势与他们开发的随机势能面行走(stochastic surfaces walking, SSW)全局搜索方法<sup>[54]</sup>共同集成在其LASP软件包中<sup>[55]</sup>, 形成SSW-NN高效势能面搜索平台. 这一平台已被广泛用于多相催化理论模拟工作中, 如针对铁基费托合成, 该团队模拟了典型实验条件下铁碳化合物的结构和反应活性, 基于广泛的全局优化计算结果绘制了热力学凸包图, 细致阐述了反应气氛下各铁碳物相的热力学稳定性. 研究者还基于搜索所得结构进行切面, 并对切面后结构进行全局结构优化, 进一步揭示了反应气氛下的热力学最稳定表面结构、从切面结构到最稳定结构过程中的表面重构行为, 以及重构后表面的活性中心性质(图3(b))<sup>[6,56,57]</sup>. 类似方法也被用于CO/H<sub>2</sub>合成气催化转化到甲醇的ZnCrO二元金属氧化物催

化剂全局相图计算与活性相结构搜索<sup>[58]</sup>以及Ag氧化物表面乙烯环氧过程的活性位点全局搜索<sup>[59,60]</sup>等热力学最概然催化活性相研究工作. 进一步的, 基于LASP平台, 刘智攀及合作者还开发了一套微观动力学引导的机器学习路径搜索方法, 这一方法通过SSW-NN方法高效搜索表面基元反应的初末态和过渡态信息, 存储在反应对数据库中, 并通过微观动力学方法筛选低能量反应路径, 实现机器学习势加持的高效表面催化反应搜索. 该方法被成功应用于Cu-Zn催化剂上CO/CO<sub>2</sub>混合气体氢化合成甲醇的反应网络和反应动力学自动化计算, 并基于所得结果讨论了催化剂上Zn在甲醇合成过程中的关键作用<sup>[61]</sup>.

高维神经网络势在势能面拟合与实际催化计算任务中取得了不小的成功, 其架构所满足的保守性和物理对称性等性质也为后来的神经网络势架构设计与优化奠定了关键基调, 直接基于高维神经网络势基础架构进行优化是神经网络势发展的一个重要方向. 比如, 传统的高维神经网络势基于ACSFs描述符构建, 这一表示针对局域化学环境采用显式区分的两体项、三体

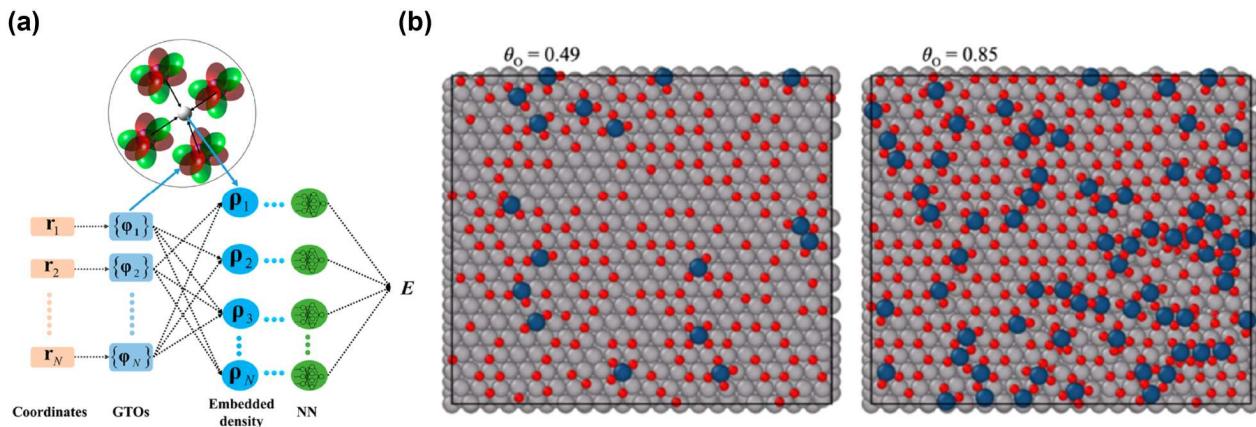


图 4 (网络版彩色) (a) EANN的模型架构<sup>[62]</sup>, Copyright © 2019 American Chemical Society; (b) 基于EANN的大尺度GCMC揭示具有不同氧覆盖度的Pt(111)表面结构<sup>[63]</sup>, Copyright © 2022 The Authors, Open Access

**Figure 4** (Color online) (a) Model architecture of EANN potential<sup>[62]</sup>, Copyright © 2019 American Chemical Society; (b) large-scale GCMC calculation based on EANN potential revealed the structure of Pt(111) surface under different oxygen coverage<sup>[63]</sup>, Copyright © 2022 The Authors, Open Access

项等多体函数表达，将带来较大的编程难度和计算成本。蒋彬及其合作者<sup>[62]</sup>开发的嵌入原子神经网络(embedding atom neural network, EANN)给出了一种改进方案：如图4(a)所示，EANN从嵌入原子模型这一经典力场的基础上发展而来，其以高斯轨道基函数线性组合的平方构建嵌入原子电荷密度，作为原子局域化学环境的表示，物理意义清晰且无需区分描述符多体作用项，并在这一基础上采用神经网络构建这一表示与原子局域能量的泛函。EANN势在模型概念、代码实现和数值计算上都相对简单，近年来在多相催化理论研究中获得了很多关注。如胡培君与合作者<sup>[63]</sup>通过结合遗传算法(genetic algorithm, GA)、分子动力学(molecular dynamics, MD)和巨正则系综蒙特卡罗(grand canonical Monte Carlo, GCMC)等势能面搜索方法的主动学习策略获取数据集训练Pt-O二元EANN势，并基于GCMC方法在300 K和1 mbar的氧气气氛下开展平整和台阶形貌的Pt表面微观氧化机制的计算。通过这一方法，研究者揭示了Pt表面在不同氧覆盖度下的稳定结构(图4(b))，识别了表面上的关键氧化物结构(如方形平面PtO<sub>4</sub>、最小条纹Pt<sub>2</sub>O<sub>6</sub>等)，并讨论了这些结构的微观形成机制，为催化表面活性位点的动态形成与演化机制研究提供了参考。类似方法还被他们用于锌氧化物催化合成气甲醇转化过程中氧空位和金属掺杂剂的微观作用机制<sup>[64,65]</sup>，以及ZnCrO催化合成气甲醇转化过程中H-CO键断裂过程活性位点的全局搜索等<sup>[66]</sup>。

上述基于描述符的神经网络势具有一个共同的问

题：模型输入，即体系的原子类别与坐标，需要先经过专家构造的代数描述符处理，再通过神经网络映射到势能面信息，模型的质量将显著依赖于描述符的表示能力。相比之下，端到端(end-to-end)神经网络势具有直接从原子类型和坐标中自主学习合适的原子局域化学环境表示的能力，可以避免基于人工先验的描述符处理，因此在近年来得到了更多的关注和发展<sup>[48]</sup>。

端到端神经网络势可以通过深度神经网络架构构建，其中一类典型代表是由张林峰、王涵等人<sup>[67~70]</sup>开发的，集成在深度势能分子动力学(deep potential molecular dynamics, DeePMD)方法中的深度势能(deep potential, DP)模型。该模型的架构如图5(a)所示<sup>[70]</sup>，它直接应用深度神经网络架建立三维坐标到能量的泛函，其描述符并没有一个具体的代数表达式，而是直接以体系笛卡儿坐标为输入，对体系每个原子在一定截断半径内定义一套局域内坐标 $R^i$ ，并通过基于这一内坐标的嵌入矩阵(embedding matrix)变换和矩阵转置相乘，将三维欧氏空间的原子局域环境表示成一个高维的特征矩阵(feature matrix)。这一特征矩阵确保了体系原子坐标的平移、旋转和置换的不变性，使模型能够对物理意义相同的微观系统给出一致的表示和预测结果；它再通过一套以神经网络构成的拟合网络(fitting-net)映射到体系各原子能量，形成一个完整的端到端编码器-解码器架构。DP势具备较强的可迁移性与可推广能力，在多相催化模拟中也有广泛的应用。如Bonati等人<sup>[71]</sup>基于增强采样分子动力学方法，强化键解离过程

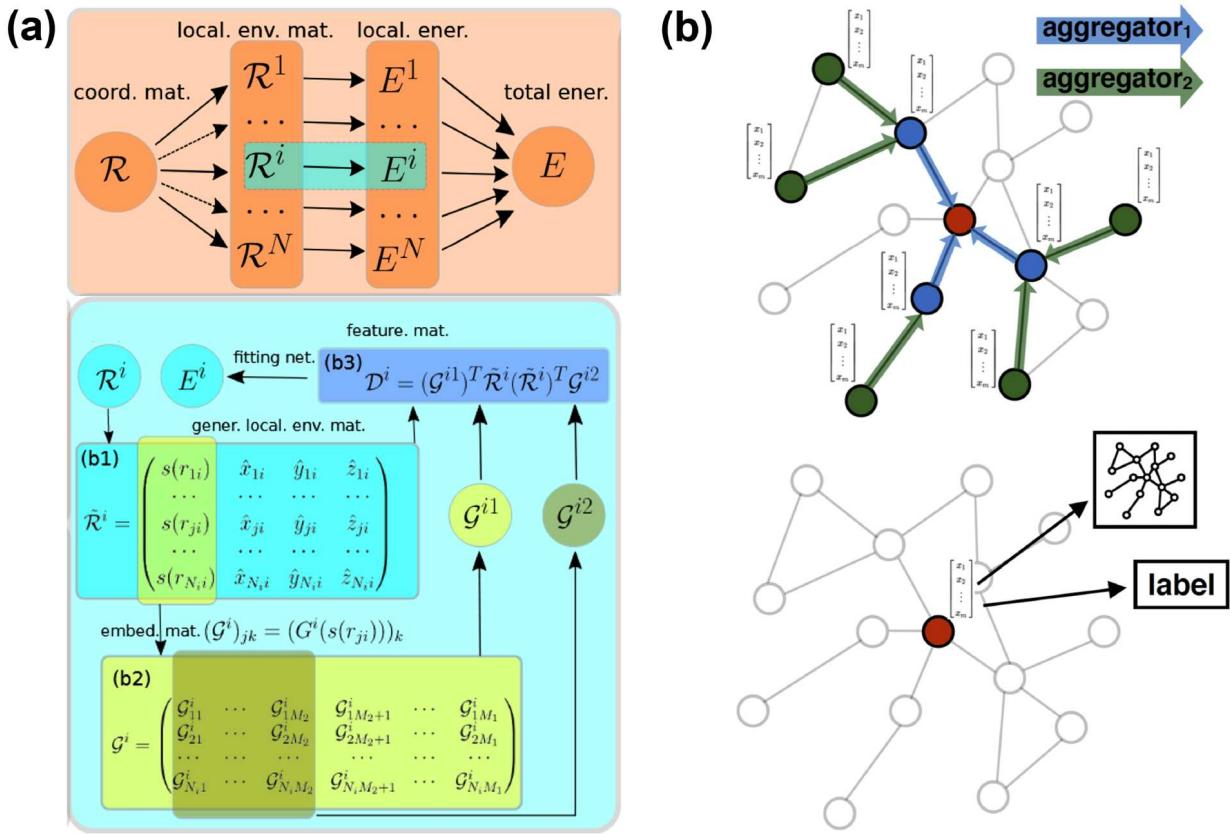


图 5 (网络版彩色)端到端神经网络势. (a) DP势的模型架构, 其描述符通过基于局域内坐标的矩阵变换提取多体相互作用特征<sup>[70]</sup>; (b) 图神经网络势的描述符, 图顶点包含原子特征信息, 通过邻接边消息传递机制提取多体相互作用特征<sup>[75]</sup>

**Figure 5** (Color online) End-to-end NNPs. (a) Architecture of DP potential, where the descriptor extracts features of many-body interactions based on matrix transformation of local internal coordinates<sup>[70]</sup>; (b) descriptors of GNN potentials, the nodes save atomic features and the edges extract features of interactions<sup>[75]</sup>

的稀有事件采样, 以训练适合用于描述化学反应过程的DP势, 并用于Fe(111)表面合成氨过程的研究, 包括不同温度下的表面活性中心形貌变化, 以及其上N<sub>2</sub>的吸附与解离过程详细机制. 这种基于分子动力学的动态反应模拟方法能超越传统过渡态理论的谐振近似, 揭示反应条件下活性中心的动态效应, 以及动态效应影响下的反应过渡态和自由能垒, 且这一方法计算成本高得不足可被高效率的DP势所克服. DP势也被应用于不同大小的Au团簇在CeO<sub>2</sub>氧化物载体表面上烧结过程的MD模拟<sup>[72]</sup>; 不同覆盖度下Pd(111)面上CO氧化过程中的活性中心实时结构演化机制探究<sup>[73]</sup>; 通过Cu<sub>55</sub>团簇上CO<sub>2</sub>解离过程的约束MD模拟揭示团簇预熔化过程带来的异常熵效应<sup>[10]</sup>, 以及SnO<sub>2</sub>(110)/水界面上的水介导质子跳跃机制研究<sup>[74]</sup>等多相催化计算模拟工作.

端到端神经网络势的另一套发展路线是基于消息传递机制(message-passing)的图神经网络(graph neural network, GNN). 在这类模型中, 分子、材料等物质体系以图的结构表示, 其顶点表示原子, 边则表示化学键等原子间相互作用, 图的结构自然保证了表示的物理对称性. 如图5(b)所示<sup>[75]</sup>, 这一图表示在GNN中通过消息传递机制进行自适应更新, 即节点通过邻接边接收邻居节点的信息, 并结合自身已有特征更新其特征. 这一机制可以捕捉到图结构中的局部和全局信息, 使图表示编码器提取到节点间的复杂关联, 以提高模型预测性能. 图神经网络势百花齐放, 包括CGCNN<sup>[76]</sup>, NequIP<sup>[77]</sup>, DimeNet<sup>[78]</sup>, GemNet-OC<sup>[79]</sup>, MACE<sup>[80]</sup>, EquiformerV2<sup>[81]</sup>, M3GNet<sup>[82]</sup>, CHGNET<sup>[83]</sup>, GPTFF<sup>[84]</sup>等. 它们在表面吸附能预测、表面动力学和表面催化反应模拟中有着广泛的应用<sup>[85~89]</sup>.

上述机器学习势模型通过不同的描述符设计策略，保证了模型以遵循物理系统对称不变性的方式表示原子局域化学环境信息与原子间关联信息的能力，从算法层面确保了模型具备一定的外推能力，能够在仅含数十或一两百原子的数据集上训练完善之后，应用于几百、上千乃至上万原子的空间尺度和较长时间尺度的计算模拟任务。然而，仅凭这些精心设计的机器学习势算法还不足以针对多相催化等特定问题构建机器学习势，要想训练得到能用于特定多相催化体系理论模拟的高质量机器学习势，就离不开针对对应问题的、具有足够特征多样性和足够质量的势能面数据集。一个简单的案例是，仅以固体体相结构及其势能面信息作为数据集训练得到的机器学习势在开展固体表面的计算模拟时往往无法得到正确的结果，因为机器学习势未能学习到关于表面物相特征的任何信息，不具备将固体体相特征信息自行外推到表面物相特征信息的能力。类似的，基于结构优化后的平衡结构数据集训练得到的机器学习势因缺乏化学反应过程中非平衡结构的特征信息，一般也难以在化学反应模拟和分子动力学模拟中表现出足够的精度<sup>[90]</sup>。要想使机器学习势在这些特定领域的模拟任务上具备足够好的表现，相关的数据集必不可少。这些数据集可以直接在对应问题的原子尺度模型上，通过第一性原理方法或特定力场方法开展的化学键扫描、结构优化、GA、MD、MC、SSW等势能面搜索方法进行采样之后加以足够精度的第一性原理计算打标签得到。这种做法直接有效，但一般是不够高效且自动化的。主动学习(active learning)是所有监督学习模型常用的训练集更新策略，其核心思想是：从基于当前(未收敛)模型或某个简化替代模型所产生的采样数据中，以某种预设的方案抽取最有助于提升模型性能的数据；通过循环执行采样与数据选取、数据打标、模型训练三个主要步骤，完成这一自动化工作流。这种策略能同时提升机器学习数据集质量及其对应的模型性能，且能避免过于庞大与冗余的数据集的生成，降低用于数据集生成的计算消耗。针对机器学习势开展主动学习训练的平台包括张与之等人<sup>[91,92]</sup>开发的专用于DP势主动学习的DP-GEN软件，以及许嘉琰等人<sup>[93]</sup>开发的GDPy软件等。在此基础上，为进一步提升主动学习过程中数据选取的针对性，基于稀有事件采样和数据特征不确定度指标的主动学习方法正得到持续的发展<sup>[66,85,94~98]</sup>。与此同时，随着物质科学各领域计算模拟

的发展，针对分子、材料、催化等不同领域的开放性数据库也逐步被建立了起来。其中包括QM9<sup>[99]</sup>、MD17<sup>[46]</sup>等分子领域数据库，Materials Project<sup>[100]</sup>、Materials Project Trajectory(MPtrj)<sup>[83]</sup>、Alexandria<sup>[101]</sup>、OMat24<sup>[102]</sup>等材料领域的数据库，针对小分子反应过渡态的数据库Transition1x<sup>[103]</sup>，以及针对多相催化领域的数据库OC20<sup>[104]</sup>、OC22<sup>[88]</sup>等，它们极大便利了针对特定问题的机器学习势的快速部署。然而，以OC22为代表的多相催化领域数据库大多只包括表面吸附构型，缺少反应过渡态构型，且这些数据库各有自己的一套第一性原理计算参数，在针对特定表面反应模拟需求构建机器学习力场时，这些数据库中的数据往往难以直接使用，研究者往往仍需自行构造相应训练数据。

近年来，随着势函数模型和势能面数据库的快速发展，具备强表示能力和全元素周期表编码能力的通用描述符、能从海量数据集中将元素周期表中各元素间复杂相互作用的信息提取到单个模型中，并能同时应用于多领域不同模拟任务的通用机器学习势吸引了极大的关注。一般来说，传统机器学习势需要针对特定的元素类型及元素类型组合构建特定的描述符，导致描述符的数目随元素类型数目呈幂指数趋势增长，且每加入一个新的元素类别，相应的描述符都要重新构建，所需要的训练数据量也会随之增加，依然没有真正克服“维数灾难”，从而导致传统机器学习势往往局限于只包含较少给定元素类别的具体问题的模拟。而通用机器学习势则为克服以上限制提供了有效的解决方案。同时，通用机器学习势往往会在模型架构上做出针对性改进，以强化模型的跨不同类型结构特征表示能力和在大量数据集上开展训练的能力。由于图的结构能自然嵌入元素类别信息，因此以GemNet-OC、M3GNet为代表的的消息传递图神经网络势能直接用于通用机器学习势的构建。其中一个突出案例是微软AI for Science研究院Yang等人<sup>[105]</sup>基于改进的M3GNet架构在约1700万材料数据集上构建的MatterSim通用机器学习势，其训练数据集涵盖数十种元素类别、0~5000 K的温度区间和0~1000 GPa的压力区间，且同时包括近平衡态和非平衡态结构，多样性远超现有开源数据集，使得模型能在未经进一步训练的情况下，直接应用于不同领域的多种模拟任务，并在结构能量与原子受力预测、材料稳定性和力学性质预测等任务上具有显著优于其他机器学习势的表现。另一种构建

通用机器学习势的方法是在描述符中针对性地构建元素类别的表示，并通过深度神经网络策略强化模型表示能力，如北京科学智能研究院等机构的张锋等人通过在DP模型描述符中引入类型嵌入网络构建元素类别信息的表示，并结合包括门控注意力机制在内的Transformer架构，发展出了包括DPA-1<sup>[106]</sup>，DPA-2<sup>[107]</sup>在内的系列通用机器学习势。其中，DPA-2能通过多任务训练策略，同时在分子、材料、催化等领域的囊括数十种元素且具备不同结构特征和计算参数设置的开源数据库上进行训练，并将这些数据中的化学物理信息提取到一套通用的编码器中，构建周期表通用的“大原子模型”。测试表明，在训练数据相同的情况下，DPA-2模型在铁电材料、固态电解质、半导体、有机分子等领域均具备比传统DP模型更高的精度<sup>[107]</sup>。在预先训练好的上游通用机器学习势的基础上，研究者还可以通过微调的方式，将通用机器学习势快速部署在目标领域的小量势能面数据集上，这一策略被称为预训练-微调，属于迁移学习的一种。以DPA-2为例，相关研究者公开发布了基于这一架构经由多任务训练机制业已训练完成的预训练通用模型，供其他领域研究者直接使用。当这些领域研究者有针对目标问题自行构建的数据集时，可以通过程序将预训练模型的拟合网络输出与已有数据集的势能面标签对齐，并在这之后以预训练模型参数为基础训练自己的模型，此即为微调方法。微调所得的模型能保留上游预训练模型的通用编码器架构及其在预训练过程中提取到的化学物理信息，确保了模型的精度和泛用性。在实际测试中，基于预训练-微调策略训练的模型相比从头训练的模型往往在测试集上具有更高的精度，且在几十个训练数据上开展微调得到的模型就已经有着和从头训练模型在上千数据下训练所得模型相同的精度，说明了预训练-微调策略的有效性<sup>[107]</sup>。这一预训练-微调策略已经开始在多相催化研究中得到应用。如刘锦程及其合作者<sup>[108]</sup>基于GemNet-OC和DPA-1等模型，通过针对性的数据采样工作流，从头构建了催化大原子模型，并通过在线局部微调的方式，在催化模拟任务中自动化地完善数据集稀有事件采样和模型性能。多相催化体系复杂度的一大来源即是其囊括了固体、表面和吸附分子三种化学特征，它们之间的复杂相互作用，以及化学键解离等势能面稀有事件，在一个模型内描述好这些问题是非常具有挑战的。可以预见的是，由于通用机器学习势出色的跨化学体系训练和预

测性能，它在理论多相催化模拟中具有广阔的发展和应用前景。当然，这一前景依赖于更多高质量数据的产生与利用。

### 3 生成式人工智能

近年来，生成式人工智能的迅猛发展为物质科学领域研究创造了大量的机遇。其中，最引人注目的就是以ChatGPT为代表的大语言模型(large language models, LLMs)，它们基于海量文本信息训练，并能根据特定提问回答对应问题，已经成为日常生活和科学研究中心重要的助手。大语言模型在多相催化等化学领域具有广阔的应用前景，在简单开展提示词工程(prompt engineering)之后，它们就可能具有帮助实验和理论研究者自动化工作流程与解决一些繁琐问题的能力，并在针对化学领域理解性问题的问答中体现出一定的性能<sup>[109]</sup>。若是能进行针对性的训练与调整，大语言模型将具有更为惊艳的表现。Zhao等人<sup>[110]</sup>开发了一个针对化学问题的基座大语言模型ChemDFM。这一模型基于已有LLaMa-13B开源大语言模型框架<sup>[111]</sup>，通过先在大量科研文献、书籍和相关网络资料上预训练，再在实验和理论计算得到的数据库上开展指导性微调的方式构建，具备强大的化学领域自然语言处理与知识推理能力，在分子识别、分子设计、分子性质预测、化学反应预测与逆合成分析等任务中名列前茅。Sprucill等人<sup>[112]</sup>开发了一套基于蒙特卡罗树搜索算法的大语言模型查询方法，强化模型复杂科学推理能力，并用于固体催化剂的设计。Ock等人<sup>[113]</sup>基于BERT大模型架构<sup>[114]</sup>开发了一套用于吸附能预测的大语言模型CatBERTa，它无需原子坐标输入，能直接基于文本信息自动提取表面吸附特征相关的关键词开展吸附能预测，在OC20等数据集上具有比传统的GNN机器学习势更灵活的数据处理方式和更广泛的适用性。他们还基于ChatGPT-4开发了一套表面吸附构型预测的智能代理程序Adsorb-Agent，并基于这一程序开展了氮还原反应中CuPd<sub>3</sub>和MoPd<sub>3</sub>合金上氮物种最优吸附结构的自动化识别，其具有明显超出传统全局结构搜索方法的性能<sup>[115]</sup>。

除了基于文本生成的大语言模型，生成式人工智能的方法也具有基于针对特定科学问题的生成模型，直接用于催化体系的理论模拟的潜力，相关应用也正得到广泛的探索。Ishikawa<sup>[116]</sup>基于DFT精度的表面微观动力学模拟所得催化剂转化频率(turnover frequency,

TOF)训练生成对抗网络(generative adversarial network, GAN), 用于合成氨反应高效合金催化剂的快速生成. 在经过多次迭代后, GAN生成的催化剂表面逐渐表现出更高的TOF值, 特别是在第5次迭代时, 生成的 $\text{Rh}_8\text{Ru}_{78}$ 表面的TOF值达到了初始数据集最高TOF的10倍以上, 表明生成式模型具备针对目标催化性质通过外推生成新材料的能力, 为催化剂理性设计提供了新的思路. 然而, 通过GAN方法训练得到的生成器在多样性和生成效率方面都具有一定的不足. 而基于扩散模型(diffusion model)的应用方案在领域内受到了较多的关注. 如段辰儒等人<sup>[117]</sup>开发了一套用于化学反应快速生成的扩散生成模型OA-ReactDiff, 该模型能在几秒内快速生成小分子反应的精确过渡态结构, 并具有在广域化学空间内探索化学反应的能力, 为扩散生成模型助力高效催化剂开发提供了开创性思路. 微软AI for Science研究院Claudio等人<sup>[118]</sup>构建了一套包含多种生成策略的扩散生成模型MatterGen, 它能直接生成稳定、多样化且满足目标约束(如化学组成、对称性、磁密度、能带带隙和体积模量等属性)的材料. Hammer及其合作者<sup>[119]</sup>设计了一种基于去噪网络和旋转等变神经网络的扩散生成模型, 并通过引入基底固定和 $z$ 方向限制与通用机器学习势引导采样的策略训练这一模型, 所得模型在低能表面结构生成方面具有显著优势, 生成的 $\text{Ag}_{29}\text{O}_{22}$ 表面结构与实验上通过STM观察到的结构高度一致, 且能在诸如 $\text{SnO}/\text{Pt}_3\text{Sn}(111)$ 这类复杂表面系统中生成高度稳定的缺陷氧化物结构, 是扩散生成模型在表面催化领域应用的重要进展. 然而, 生成式模型在多相催化反应中的应用仍然鲜有报道, 一个主要的可能原因是缺乏针对多相催化反应体系的大规模高质量数据集, 以及针对性的模型设计思路, 这些既是生成式模型进一步应用的挑战, 也是生成式模型方法发展的机遇.

#### 4 总结与展望

本文围绕理论多相催化领域中常见的机器学习方法, 讨论涵盖了早期的Sabatier原理到基于物理描述符的统计模型、用于高效地描述原子间相互作用的机器学习势函数与应用前景广阔的生成式模型. 受制于多相催化反应的复杂性, 实际研究中往往难以使用第一性原理计算方案对大时空尺度的体系进行系统研究, 由此为机器学习方法留下了充分的发挥空间. 近十年来, 领域内的新机器学习策略层出不穷: SISSO拟合、

通用机器学习势、扩散生成式模型等方法为理论多相催化提供了更多可能性, 推动了计算对象由静态到动态、理想表面到重构表面、小体系到大体系、单个结构到统计系综的转变. 然而, 新的挑战也随之而来, 主要可分为以下三个方面.

(1) 模型物理可解释性与描述符可解释性的距离. 基于小数据集的统计和回归模型通常具有比较好的可解释性, 但随着统计回归方式变得复杂, 将具有清晰物理意义的一系列描述符通过这些复杂统计方式拟合得到能够描述具体问题的模型后, 所得模型并不一定具有确定的物理意义. 这一问题给复杂环境下催化剂理性设计带来了一定困难, 由此对研究者的领域知识与研究手段都提出了较高的要求.

(2) 多元素复杂体系机器学习势函数的训练和应用仍然存在挑战. 在通用机器学习势兴起之前, 由于模型设计和训练难度, 领域内很少得到能够描述四元素以上体系的势函数. 但即使有了通用机器学习势与预训练-微调策略, 高精度模型的训练仍然高度依赖于具体训练参数和高质量数据集的选取, 针对表面催化, 尤其是表面反应问题的机器学习势训练的数据采样方式和数据集管理依然是一大挑战. 同时, 这种高精度通用性机器学习势的训练和推理效率可能会比相对简单的势函数慢上1~2个数量级. 如何兼顾精度和计算效率, 将高精度通用机器学习势推广到更大时空尺度和更为自动化的催化过程模拟, 并让通用机器学习势在模拟的过程中在线训练, 自主更新, 依然是机器学习势架构和训练方法设计的一大难题.

(3) 生成式模型在理论多相催化的进一步应用. 当前大部分的应用停留在大语言模型本身, 即应用文献数据或其他文本信息生成相应的化学信息, 对多相催化领域具体问题涉及较少. 如何针对性实现原子尺度科学信息的直接高效生成, 是领域内将要克服的一大挑战, 这依赖于新的模型架构设计, 以及针对多相催化问题的大规模的高质量数据集构建. 同时, 让生成式模型所得到的催化剂结构和化学反应在实验中得到验证, 并结合实验信息进一步完善生成式模型的性能, 将是该领域发展的重要方向.

总的来说, 未来的机器学习模型还需朝着模拟真工况条件下多相催化表面形貌演化与详细反应机制的目标前进, 在模型部署成本、模型精度与泛用性等方面上取得进一步提升, 让人工智能驱动的催化剂理性设计时代早日到来.

## 参考文献

- 1 Fechete I, Wang Y, Védrine J C. The past, present and future of heterogeneous catalysis. *Catal Today*, 2012, 189: 2–27
- 2 Chen B W J, Xu L, Mavrikakis M. Computational methods in heterogeneous catalysis. *Chem Rev*, 2021, 121: 1007–1048
- 3 Ishaq H, Crawford C. Review of ammonia production and utilization: enabling clean energy transition and net-zero climate targets. *Energy Convers Manage*, 2024, 300: 117869
- 4 Wang Z, Hu P. Towards rational catalyst design: a general optimization framework. *Phil Trans R Soc A*, 2016, 374: 20150078
- 5 Xu J, Cao X M, Hu P. Perspective on computational reaction prediction using machine learning methods in heterogeneous catalysis. *Phys Chem Chem Phys*, 2021, 23: 11155–11179
- 6 Liu Q Y, Shang C, Liu Z P. *In situ* active site for Fe-catalyzed Fischer–Tropsch synthesis: recent progress and future challenges. *J Phys Chem Lett*, 2022, 13: 3342–3352
- 7 Chen B, Wang D, Duan X, et al. Charge-tuned CO activation over a  $\chi$ -Fe<sub>5</sub>C<sub>2</sub> Fischer–Tropsch catalyst. *ACS Catal*, 2018, 8: 2709–2714
- 8 Filot I A W, van Santen R A, Hensen E J M. The optimally performing Fischer–Tropsch catalyst. *Angew Chem Int Ed*, 2014, 53: 12746–12750
- 9 Astruc D. Introduction: nanoparticles in catalysis. *Chem Rev*, 2020, 120: 461–463
- 10 Gong F, Liu Y, Wang Y, et al. Machine learning molecular dynamics shows anomalous entropic effect on catalysis through surface pre - melting of naanoclusters. *Angew Chem Int Ed*, 2024, 63: e202405379
- 11 Erdem Günay M, Yıldırım R. Recent advances in knowledge discovery for heterogeneous catalysis using machine learning. *Catal Rev*, 2021, 63: 120–164
- 12 Wang F Y, Miao Q H. Novel paradigm of AI-driven scientific research: from AI4S to intelligent science (in Chinese). *Bull Chin Acad Sci*, 2023, 38: 536–540 [王飞跃, 缪青海. 人工智能驱动的科学新范式: 从AI4S到智能科学. 中国科学院院刊, 2023, 38: 536–540]
- 13 Medford A J, Vojvodic A, Hummelshøj J S, et al. From the Sabatier principle to a predictive theory of transition-metal heterogeneous catalysis. *J Catal*, 2015, 328: 36–42
- 14 Motagamwala A H, Dumesic J A. Microkinetic modeling: a tool for rational catalyst design. *Chem Rev*, 2021, 121: 1049–1076
- 15 Jin H, Zhao R, Cui P, et al. Sabatier phenomenon in hydrogenation reactions induced by single-atom density. *J Am Chem Soc*, 2023, 145: 12023–12032
- 16 Nørskov J K, Rossmeisl J, Logadottir A, et al. Origin of the overpotential for oxygen reduction at a fuel-cell cathode. *J Phys Chem B*, 2004, 108: 17886–17892
- 17 Greeley J. Theoretical heterogeneous catalysis: scaling relationships and computational catalyst design. *Annu Rev Chem Biomol Eng*, 2016, 7: 605–635
- 18 Schmack R, Friedrich A, Kondratenko E V, et al. A meta-analysis of catalytic literature data reveals property-performance correlations for the OCM reaction. *Nat Commun*, 2019, 10: 441
- 19 Hong W T, Welsch R E, Shao-Horn Y. Descriptors of oxygen-evolution activity for oxides: a statistical evaluation. *J Phys Chem C*, 2016, 120: 78–86
- 20 Pablo-García S, Sabadell-Rendón A, Saadun A J, et al. Generalizing performance equations in heterogeneous catalysis from hybrid data and statistical learning. *ACS Catal*, 2022, 12: 1581–1594
- 21 Ouyang R, Curtarolo S, Ahmetcik E, et al. SISSO: a compressed-sensing method for identifying the best low-dimensional descriptor in an immensity of offered candidates. *Phys Rev Mater*, 2018, 2: 083802
- 22 Shu W, Li J, Liu J X, et al. Structure sensitivity of metal catalysts revealed by interpretable machine learning and first-principles calculations. *J Am Chem Soc*, 2024, 146: 8737–8745
- 23 Wang T, Hu J, Ouyang R, et al. Nature of metal-support interaction for metal catalysts on oxide supports. *Science*, 2024, 386: 915–920
- 24 Zhou N, Liu W, Jan F, et al. Efficient screening of metal promoters of Pt catalysts for C–H bond activation in propane dehydrogenation from a combined first-principles calculations and machine-learning study. *ACS Omega*, 2023, 8: 23982–23990
- 25 Wang S, Jiang J. Interpretable catalysis models using machine learning with spectroscopic descriptors. *ACS Catal*, 2023, 13: 7428–7436
- 26 Guan Y, Chaffart D, Liu G, et al. Machine learning in solid heterogeneous catalysis: recent developments, challenges and perspectives. *Chem Eng Sci*, 2022, 248: 117224
- 27 Bhardwaj A, Ahluwalia A S, Pant K K, et al. A principal component analysis assisted machine learning modeling and validation of methanol formation over Cu-based catalysts in direct CO<sub>2</sub> hydrogenation. *Separation Purification Tech*, 2023, 324: 124576
- 28 Pirro L, Mendes P S F, Paret S, et al. Descriptor–property relationships in heterogeneous catalysis: exploiting synergies between statistics and fundamental kinetic modelling. *Catal Sci Technol*, 2022, 9: 3109–3125

- 29 Ulissi Z W, Medford A J, Bligaard T, et al. To address surface reaction network complexity using scaling relations machine learning and DFT calculations. *Nat Commun*, 2022, 8: 14621
- 30 Avşar E. Dimensionality reduction for predicting CO conversion in water gas shift reaction over Pt-based catalysts using support vector regression models. *Int J Hydrogen Energy*, 2017, 42: 23326–23333
- 31 Bae S, Lee H, Shin J, et al. Data-driven inference of synthesis guidelines for high-performance zeolite-based selective catalytic reduction catalysts at low temperatures. *Chem Mater*, 2022, 34: 7761–7773
- 32 Wexler R B, Martirez J M P, Rappe A M. Chemical pressure-driven enhancement of the hydrogen evolving activity of Ni<sub>2</sub>P from nonmetal surface doping interpreted via machine learning. *J Am Chem Soc*, 2018, 140: 4678–4683
- 33 Zhu Q, Gu Y, Liang X, et al. A machine learning model to predict CO<sub>2</sub> reduction reactivity and products transferred from metal-zeolites. *ACS Catal*, 2022, 12: 12336–12348
- 34 Xie Y, Feng S, Deng L, et al. Inverse design of chiral functional films by a robotic AI-guided system. *Nat Commun*, 2023, 14: 6177
- 35 Zhu Q, Huang Y, Zhou D, et al. Automated synthesis of oxygen-producing catalysts from Martian meteorites by a robotic AI chemist. *Nat Synth*, 2023, 3: 319–328
- 36 Kocer E, Ko T W, Behler J. Neural network potentials: a concise overview of methods. *Annu Rev Phys Chem*, 2022, 73: 163–186
- 37 Cornell W D, Cieplak P, Bayly C I, et al. A second generation force field for the simulation of proteins, nucleic acids, and organic molecules. *J Am. Chem. Soc.* 1995, 117, 5179–5197. *J Am Chem Soc*, 1996, 118: 2309
- 38 Brooks B R, Bruccoleri R E, Olafson B D, et al. CHARMM: a program for macromolecular energy, minimization, and dynamics calculations. *J Comput Chem*, 1983, 4: 187–217
- 39 van Duin A C T, Dasgupta S, Lorant F, et al. ReaxFF: a reactive force field for hydrocarbons. *J Phys Chem A*, 2001, 105: 9396–9409
- 40 Blank T B, Brown S D, Calhoun A W, et al. Neural network models of potential energy surfaces. *J Chem Phys*, 1995, 103: 4129–4137
- 41 Sumpter B G, Noid D W. Potential energy surfaces for macromolecules. A neural network technique. *Chem Phys Lett*, 1992, 192: 455–462
- 42 He K M, Zhang X Y, Ren S Q, et al. Deep residual learning for image recognition. 2015, arXiv: 1512.03385
- 43 Vaswani A, Shazeer N, Parmar N, et al. Attention is all you need. 2017, arXiv: 1706.03762
- 44 Bartók A P, Payne M C, Kondor R, et al. Gaussian approximation potentials: the accuracy of quantum mechanics, without the electrons. *Phys Rev Lett*, 2010, 104: 136403
- 45 Thompson A P, Swiler L P, Trott C R, et al. Spectral neighbor analysis method for automated generation of quantum-accurate interatomic potentials. *J Comput Phys*, 2015, 285: 316–330
- 46 Chmiela S, Tkatchenko A, Sauceda H E, et al. Machine learning of accurate energy-conserving molecular force fields. *Sci Adv*, 2017, 3: e1603015
- 47 Chmiela S, Sauceda H E, Poltavsky I, et al. sGML: constructing accurate and data efficient molecular force fields using machine learning. *Comput Phys Commun*, 2019, 240: 38–45
- 48 Unke O T, Chmiela S, Sauceda H E, et al. Machine learning force fields. *Chem Rev*, 2021, 121: 10142–10186
- 49 Behler J, Parrinello M. Generalized neural-network representation of high-dimensional potential-energy surfaces. *Phys Rev Lett*, 2007, 98: 146401
- 50 Behler J. Four generations of high-dimensional neural network potentials. *Chem Rev*, 2021, 121: 10037–10072
- 51 Behler J. Atom-centered symmetry functions for constructing high-dimensional neural network potentials. *J Chem Phys*, 2011, 134: 074106
- 52 Huang S D, Shang C, Zhang X J, et al. Material discovery by combining stochastic surface walking global optimization with a neural network. *Chem Sci*, 2017, 8: 6327–6337
- 53 Huang S D, Shang C, Kang P L, et al. Atomic structure of boron resolved using machine learning and global sampling. *Chem Sci*, 2018, 9: 8644–8655
- 54 Shang C, Liu Z P. Stochastic surface walking method for structure prediction and pathway searching. *J Chem Theor Comput*, 2013, 9: 1838–1845
- 55 Huang S, Shang C, Kang P, et al. LASP: fast global potential energy surface exploration. *WIREs Comput Mol Sci*, 2019, 9: e1415
- 56 Liu Q Y, Shang C, Liu Z P. *In situ* active site for CO activation in Fe-catalyzed Fischer–Tropsch synthesis from machine learning. *J Am Chem Soc*, 2021, 143: 11109–11120
- 57 Liu Q Y, Chen D, Shang C, et al. An optimal Fe–C coordination ensemble for hydrocarbon chain growth: a full Fischer–Tropsch synthesis mechanism from machine learning. *Chem Sci*, 2023, 14: 9461–9475
- 58 Ma S, Huang S D, Liu Z P. Dynamic coordination of cations and catalytic selectivity on zinc–chromium oxide alloys during syngas conversion. *Nat Catal*, 2019, 2: 671–677
- 59 Chen D, Shang C, Liu Z P. Automated search for optimal surface phases (ASOPs) in grand canonical ensemble powered by machine learning. *J*

- Chem Phys*, 2022, 156: 094104
- 60 Chen D, Chen L, Zhao Q C, et al. Square-pyramidal subsurface oxygen [ $\text{Ag}_4\text{OAg}$ ] drives selective ethene epoxidation on silver. *Nat Catal*, 2024, 7: 536–545
- 61 Shi Y F, Kang P L, Shang C, et al. Methanol synthesis from  $\text{CO}_2/\text{CO}$  mixture on Cu–Zn catalysts from microkinetics-guided machine learning pathway search. *J Am Chem Soc*, 2022, 144: 13401–13414
- 62 Zhang Y, Hu C, Jiang B. Embedded atom neural network potentials: efficient and accurate machine learning with a physically inspired representation. *J Phys Chem Lett*, 2019, 10: 4962–4967
- 63 Xu J, Xie W, Han Y, et al. Atomistic insights into the oxidation of flat and stepped platinum surfaces using large-scale machine learning potential-based grand-canonical Monte Carlo. *ACS Catal*, 2022, 12: 14812–14824
- 64 Han Y, Xu J, Xie W, et al. Comprehensive study of oxygen vacancies on the catalytic performance of  $\text{ZnO}$  for  $\text{CO}/\text{H}_2$  activation using machine learning-accelerated first-principles simulations. *ACS Catal*, 2023, 13: 5104–5113
- 65 Han Y, Xu J, Xie W, et al. Unravelling the impact of metal dopants and oxygen vacancies on syngas conversion over oxides: a machine learning-accelerated study of CO activation on Cr-doped  $\text{ZnO}$  surfaces. *ACS Catal*, 2023, 13: 15074–15086
- 66 Wu J W, Han Y L, Jia M L, et al. Identifying the active site on  $\text{Zn}_x\text{Cr}_y\text{O}_z$  for HC-O bond cleavage in syn-gas conversion. 2024, ChemRxiv: 2024-chrkx
- 67 Wang H, Zhang L, Han J, et al. DeePMD-kit: a deep learning package for many-body potential energy representation and molecular dynamics. *Comput Phys Commun*, 2018, 228: 178–184
- 68 Zeng J, Zhang D, Lu D, et al. DeePMD-kit v2: a software package for deep potential models. *J Chem Phys*, 2023, 159: 054801
- 69 Zhang L, Han J, Wang H, et al. Deep potential molecular dynamics: a scalable model with the accuracy of quantum mechanics. *Phys Rev Lett*, 2018, 120: 143001
- 70 Zhang L F, Han J Q, Wang H, et al. End-to-end symmetry preserving inter-atomic potential energy model for finite and extended systems. *Adv Neural Inform Process Syst*, 2018, 31: 4436–4446
- 71 Bonati L, Polino D, Pizzolitto C, et al. The role of dynamics in heterogeneous catalysis: surface diffusivity and  $\text{N}_2$  decomposition on Fe(111). *Proc Natl Acad Sci USA*, 2023, 120: e2313023120
- 72 Liu J C, Luo L, Xiao H, et al. Metal affinity of support dictates sintering of gold catalysts. *J Am Chem Soc*, 2022, 144: 20601–20609
- 73 Wu J, Chen D, Chen J, et al. Structural and composition evolution of palladium catalyst for CO oxidation under steady-state reaction conditions. *J Phys Chem C*, 2023, 127: 6262–6270
- 74 Jia M, Zhuang Y B, Wang F, et al. Water-mediated proton hopping mechanisms at the  $\text{SnO}_2(110)/\text{H}_2\text{O}$  interface from *ab initio* deep potential molecular dynamics. *Precision Chem*, 2024, 2: 644–654
- 75 Hamilton W L, Ying R, Leskovec J. Inductive representation learning on large graphs. *Adv Neural Inform Process Syst*, 2017, 30: 1024–1034
- 76 Xie T, Grossman J C. Crystal graph convolutional neural networks for an accurate and interpretable prediction of material properties. *Phys Rev Lett*, 2018, 120: 145301
- 77 Batzner S, Musaelian A, Sun L, et al. E(3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials. *Nat Commun*, 2022, 13: 2453
- 78 Gasteiger J, Groß J, Günemann S. Directional message passing for molecular graphs. 2020, arXiv: 2003.03123
- 79 Gasteiger J, Shuaibi M, Sriram A, et al. GemNet-OC: developing graph neural networks for large and diverse molecular simulation datasets. 2022, arXiv: 2204.02782
- 80 Batatia I, Kovács D P, Simm G N C, et al. MACE: Higher order equivariant message passing neural networks for fast and accurate force fields. 2022, arXiv: 2206.07697
- 81 Liao Y L, Wood B, Das A, et al. EquiformerV2: improved equivariant transformer for scaling to higher-degree representations. 2023, arXiv: 2306.12059
- 82 Chen C, Ong S P. A universal graph deep learning interatomic potential for the periodic table. *Nat Comput Sci*, 2022, 2: 718–728
- 83 Deng B, Zhong P, Jun K J, et al. CHGNet as a pretrained universal neural network potential for charge-informed atomistic modelling. *Nat Mach Intell*, 2023, 5: 1031–1041
- 84 Xie F, Lu T, Meng S, et al. GPTFF: a high-accuracy out-of-the-box universal AI force field for arbitrary inorganic materials. *Sci Bull*, 2024, 69: 3525–3532
- 85 Roongcharoen T, Conter G, Sementa L, et al. Machine-learning-accelerated DFT conformal sampling of catalytic processes. *J Chem Theor Comput*, 2024, 20: 9580–9591
- 86 Schwalbe-Koda D, Govindarajan N, Varley J. Comprehensive sampling of coverage effects in catalysis by leveraging generalization in neural

- network models. 2024, ChemRxiv: 2023-f6l23-v2
- 87 Stark W G, van der Oord C, Batatia I, et al. Benchmarking of machine learning interatomic potentials for reactive hydrogen dynamics at metal surfaces. 2024, arXiv: [2403.15334](#)
- 88 Tran R, Lan J, Shuaibi M, et al. The open catalyst 2022 (OC22) dataset and challenges for oxide electrocatalysts. *ACS Catal.*, 2023, 13: 3066–3084
- 89 Wander B, Shuaibi M, Kitchin J R, et al. CatTSunami: accelerating transition state energy calculations with pre-trained graph neural networks. 2024, arXiv: [2405.02078](#)
- 90 Kulichenko M, Nebgen B, Lubbers N, et al. Data generation for machine learning interatomic potentials and beyond. *Chem Rev.*, 2024, 124: 13681–13714
- 91 Zhang Y, Wang H, Chen W, et al. DP-GEN: a concurrent learning platform for the generation of reliable deep learning based potential energy models. *Comput Phys Commun.*, 2020, 253: 107206
- 92 Zhang L, Lin D Y, Wang H, et al. Active learning of uniformly accurate interatomic potentials for materials simulation. *Phys Rev Mater.*, 2019, 3: 023804
- 93 Cheng X, Wu C, Xu J, et al. Leveraging machine learning potentials for *in-situ* searching of active sites in heterogeneous catalysis. *Precision Chem.*, 2024, 2: 570–586
- 94 Busk J, Bjørn Jørgensen P, Bhowmik A, et al. Calibrated uncertainty for molecular property prediction using ensembles of message passing neural networks. *Mach Learn-Sci Technol.*, 2022, 3: 015012
- 95 Perego S, Bonati L. Data-efficient modeling of catalytic reactions via enhanced sampling and on-the-fly learning of machine learning potentials. 2024, ChemRxiv: 2024-nsp7n-v2
- 96 Schaaf L L, Fako E, De S, et al. Accurate energy barriers for catalytic reaction pathways: an automatic training protocol for machine learning force fields. *npj Comput Mater.*, 2023, 9: 180
- 97 Vandermause J, Torrisi S B, Batzner S, et al. On-the-fly active learning of interpretable Bayesian force fields for atomistic rare events. *npj Comput Mater.*, 2020, 6: 20
- 98 Blais C J, Xu C, West R H. Uncertainty quantification of linear scaling, machine learning, and density functional theory derived thermodynamics for the catalytic partial oxidation of methane on rhodium. *J Phys Chem C.*, 2024, 128: 17418–17433
- 99 Ramakrishnan R, Dral P O, Rupp M, et al. Quantum chemistry structures and properties of 134 kilo molecules. *Sci Data.*, 2014, 1: 140022
- 100 Jain A, Ong S P, Hautier G, et al. The Materials Project: a materials genome approach to accelerating materials innovation. *APL Mater.*, 2013, 1: 011002
- 101 Schmidt J, Hoffmann N, Wang H, et al. Machine-learning-assisted determination of the global zero-temperature phase diagram of materials. *Adv Mater.*, 2023, 35: 2210788
- 102 Barroso-Luque L, Shuaibi M, Fu X, et al. Open materials 2024 (OMat24) inorganic materials dataset and models. 2024, arXiv: [2410.12771](#)
- 103 Schreiner M, Bhowmik A, Vegge T, et al. Transition1x - a dataset for building generalizable reactive machine learning potentials. *Sci Data.*, 2022, 9: 779
- 104 Chanussot L, Das A, Goyal S, et al. Open catalyst 2020 (OC20) dataset and community challenges. *ACS Catal.*, 2021, 11: 6059–6072
- 105 Yang H, Hu C, Zhou Y, et al. MatterSim: a deep learning atomistic model across elements, temperatures and pressures. 2024, arXiv: [2405.04967](#)
- 106 Zhang D, Bi H, Dai F Z, et al. Pretraining of attention-based deep learning potential model for molecular simulation. *npj Comput Mater.*, 2024, 10: 94
- 107 Zhang D, Liu X, Zhang X, et al. DPA-2: a large atomic model as a multi-task learner. *npj Comput Mater.*, 2024, 10: 293
- 108 Wu Z H, Zhou L, Hou P F, et al. Catalytic large atomic model (CLAM): a machine-learning-based interatomic potential universal model. 2024, ChemRxiv: 2024-2xzct-v2
- 109 Castro Nascimento C M, Pimentel A S. Do large language models understand chemistry? A conversation with ChatGPT. *J Chem Inf Model.*, 2023, 63: 1649–1655
- 110 Zhao Z H, Ma D, Chen L, et al. ChemDFM: a large language foundation model for chemistry. 2024, arXiv: [2401.14818](#)
- 111 Touvron H, Lavril T, Izacard G, et al. LLaMA: open and efficient foundation language models. 2023, arXiv: [2302.13971](#)
- 112 Sprueill H W, Edwards C, Olarte M V, et al. Monte Carlo thought search: large language model querying for complex scientific reasoning in catalyst design. 2023, arXiv: [2310.14420](#)
- 113 Ock J, Guntuboina C, Barati Farimani A. Catalyst energy prediction with CatBERTa: unveiling feature exploration strategies through large language models. *ACS Catal.*, 2023, 13: 16032–16044
- 114 Devlin J, Chang M W, Lee K, et al. BERT: pre-training of deep bidirectional transformers for language understanding. 2018, arXiv: [1810.04805](#)
- 115 Ock J, Vinchurkar T, Jadhav Y, et al. Adsorb-agent: autonomous identification of stable adsorption configurations via large language model agent.

2024, arXiv: [2410.16658](https://arxiv.org/abs/2410.16658)

- 116 Ishikawa A. Heterogeneous catalyst design by generative adversarial network and first-principles based microkinetics. *Sci Rep*, 2022, 12: 11657
- 117 Duan C, Du Y, Jia H, et al. Accurate transition state generation with an object-aware equivariant elementary reaction diffusion model. *Nat Comput Sci*, 2023, 3: 1045–1055
- 118 Zeni C, Pinsler R, Zügner D, et al. MatterGen: a generative model for inorganic materials design. 2023, arXiv: [2312.03687](https://arxiv.org/abs/2312.03687)
- 119 Rønne N, Aspuru-Guzik A, Hammer B. Generative diffusion model for surface structure discovery. *Phys Rev B*, 2024, 110: 235427

Summary for “理论多相催化计算中的机器学习方法: 现状与挑战”

# Machine learning methods for theoretical heterogeneous catalysis: current status and challenges

Zhao-Qing Liu<sup>†</sup>, Zhe Deng<sup>†</sup> & Hong Jiang<sup>\*</sup>

*College of Chemistry and Molecular Engineering, Peking University, Beijing 100871, China*

<sup>†</sup> Equally contributed to this work

\* Corresponding author, E-mail: [jianghchem@pku.edu.cn](mailto:jianghchem@pku.edu.cn)

Catalysis, the science of accelerating chemical reactions and regulating the selectivity of reaction products, lies in the heart of modern chemical industry. Typically, catalysis can be categorized into heterogeneous catalysis and homogeneous catalysis. Over the past century, heterogeneous catalytic processes such as ammonia synthesis and Fischer-Tropsch synthesis have become an indispensable part of the chemical industry. In recent years, the rational design of catalysts has gained a lot of attention due to both economic and environmental concerns, which has raised new demands for a deep understanding of the mechanism of catalytic reactions. Machine learning methods as one of the most important areas of artificial intelligence (AI) are becoming more and more popular in practices of theoretical simulations. This review gives an overview on various machine learning methods and their applications in theoretical heterogeneous catalysis.

Starting with statistical learning models on small datasets based on specific physical descriptors, which play an important role in early theoretical heterogeneous research, this review analyzes how they connect linear scaling relationship and catalytic experiments. However, due to the complexity of heterogeneous catalytic systems, sometimes nonlinear models can be essential. In this way, the exact form of equations should be obtained before fitting, which is extremely challenging. This problem can be solved by methods based on symbolic regression like SISSO to obtain a set of expressions, and the researchers can select one by their domain knowledge. This review also shows the applications of some popular machine learning methods (such as SVM, PCA, RF, etc.) in research, most of which are based on specific descriptors, simple but have difficulty transferring to other systems.

First-principle calculations like DFT are crucial in catalysis rational design by computational simulations for surface structure determination and reaction channel exploration, but they are impractical for simulations on large and complex catalytic systems due to high computational costs. In recent years, machine learning potentials (MLPs), as tools to bridge the gap between first-principle accuracy and computational efficiency, have pushed forward the frontier of theoretical heterogeneous catalysis. Popular MLPs (LASP, EANN, DP, GemNet-OC, MACE, DPA, etc.) have numerous applications in this field, from global optimization and molecular dynamics of active surfaces under reaction conditions to automatic reaction network exploration. Recently, universal machine learning potentials have gained much attention and are very promising in heterogeneous catalysis simulations, for their promising capability to encapsulate huge chemical space of the whole periodic table in one pre-trained model, which is easy to fine-tune on a specific chemical system, efficiently deploying a specific MLP with high accuracy and transferability.

Recent advances in generative models have shown considerable promise for heterogeneous catalysis research, and the dramatic emergence of large language models like ChatGPT is especially important in this case. LLMs can incorporate information from massive natural language data, having a huge ability to push forward the understanding of topics in catalytic science. Furthermore, the concept of generative models themselves has much potential in generating novel catalytic information directly without natural language as the intermediate, but their application in heterogeneous catalysis is under limitation due to the lack of large datasets with high quality.

In a nutshell, machine learning methods play a key role in theoretical understanding of heterogeneous catalytic systems. Challenges remain due to the complexity of heterogeneous catalytic systems and lack of high-quality dataset, but the era of AI-driven rational design of surface catalysts is truly coming.

**theoretical heterogeneous catalysis, catalyst rational design, machine learning, machine learning potential, generative model**

doi: [10.1360/TB-2024-1207](https://doi.org/10.1360/TB-2024-1207)