

专题: AI 物质科学

目标性质导向的材料生成: 迈向按需构筑的材料逆向设计*

刘章赫¹⁾ 陈新宇¹⁾ 周楚桦^{1)2)†} 王金兰^{1)2)‡}

1) (东南大学物理学院, 量子材料与信息器件教育部重点实验室, 南京 211189)

2) (苏州实验室, 苏州 215004)

(2025年7月24日收到; 2025年9月10日收到修改稿)

近年来, 机器学习在材料科学中的应用显著加快了新材料的发现, 特别是在结合第一性原理计算等传统方法后, 能够高效筛选已有数据库中的潜在高性能材料. 然而, 此类方法大多局限于已有化学空间, 难以实现对全新材料结构的主动设计. 为突破这一瓶颈, 基于生成模型的材料逆向设计方法逐渐兴起, 成为探索未知结构与性质空间的重要手段. 尽管当前生成模型在晶体结构生成方面取得了初步进展, 但如何实现目标性质导向的材料生成仍面临显著挑战. 本文首先介绍了近年来在材料生成领域中具有代表性的生成模型, 包括 CDVAE, MatGAN 以及 MatterGen, 分析其在结构生成上的基本能力与局限. 随后重点探讨如何将目标性质有效引入生成模型, 实现性质导向的结构生成, 具体包括基于目标性质向量的 Con-CDVAE、融合结构约束与引导机制的 SCIGEN、通过适配器实现性质调控的微调版 MatterGen 以及结合隐空间搜索优化的 CDVAE 隐变量优化策略. 最后总结当前性质导向生成机制面临的挑战, 并展望其未来的发展方向. 本文旨在为研究者深入理解和拓展性质驱动的材料生成方法提供系统性参考和启发.

关键词: 机器学习, 生成模型, 逆向设计, 性质导向**PACS:** 07.05.Mh, 91.60.Ed, 81.05.Zx**DOI:** 10.7498/aps.74.20250989**CSTR:** 32037.14.aps.74.20250989

1 引言

材料科学正在迈入以数据驱动 (data-driven) 为核心特征的“第四范式”时代, 其中, 人工智能 (artificial intelligence, AI) 作为数据挖掘与模型构建的重要工具, 正在重塑材料设计的流程与范式^[1-3]. 生成模型 (generative models) 作为 AI 领域的关键方法之一, 因其具备从数据中学习分布并生成新材料结构的能力, 在材料设计领域具有巨大潜力. 传统的数据驱动材料设计通常遵循“正向设计”范式,

即从已有结构出发, 结合第一性原理计算 (first-principles calculation) 和机器学习方法 (machine learning, ML) 预测其物理化学性质, 再通过筛选得到具有目标性质的候选材料^[4-6]. 这种以结构为起点的设计流程虽然在过去取得了大量成果, 但其探索路径往往受限于已有材料数据库 (如 ICSD^[7], OQMD^[8], Materials Project^[9], C2DB^[10]) 中的样本结构, 难以跳出既有材料组合的范畴. 因而, 其对全新材料结构的主动发现能力有限, 难以覆盖广阔的化学空间. 这一局限在追求高性能、低成本乃至绿色可持续材料的当下, 已逐渐成为潜在高性能

* 国家重点研发计划 (批准号: 2021YFA1500703)、国家自然科学基金 (批准号: 22033002, T2321002, 22373013) 和江苏省科技计划专项资金前沿引领技术基础研究重大项目 (批准号: BK20222007, BK20232012) 资助的课题.

† 通信作者. E-mail: qh.zhou@seu.edu.cn

‡ 通信作者. E-mail: jlwang@seu.edu.cn

新材料发现的瓶颈^[11].

近年来,生成模型的发展^[12-16]为材料科学提供了一种范式跃迁的可能路径.相较于正向筛选流程,生成模型通过学习已有材料的结构分布,可以实现从目标性质反推出结构的逆向设计,即根据预设的性质需求直接生成满足该性质的潜在材料结构.这种方法具备探索未知材料空间的潜力,能够突破正向设计的搜索边界,更高效地发现新型高性能材料(见图1).特别是在面对性质标注数据有限、高维复杂性强的结构-性质空间时,生成模型可以通过在隐空间中进行连续表示与优化,生成多样化的候选材料结构,为实现性质导向的材料设计提供了新的思路与工具.

尽管生成模型已在晶体结构生成^[17,18]、分子生成^[19]等任务中取得突破性进展,但要实现“面向目标性质”的材料生成仍面临一系列核心挑战.首先,材料数据通常具有性质分布高度不均、样本数量有限、结构与性质之间高度非线性等特征,导致生成模型在构建稳定的结构-性质映射关系和实现有效的目标性质引导方面面临一定的困难.其次,不同性质对材料结构的敏感性差异显著,例如带隙对原子排布的微小变化极为敏感^[20],而形成能则更依赖于整体构型与组分分布^[21],这使得“性质导向”的生成策略必须具备足够的精度与结构感知能力.此外,如何在保持生成结构物理合理性与晶体对称性的同时,实现性质的精确可调,也成为研究的关

键难题之一.

针对上述问题,研究者们提出了一系列可用于材料生成并引入了目标性质导向机制的生成模型框架.例如, Ye等^[22]提出的 Con-CDVAE通过在编码器和解码器引入目标性质向量来增强性质引导能力,实现隐空间的性质调控,使生成过程可根据所需性质进行有针对性的解码; Okabe等^[23]提出的 SCIGEN方法能够在扩散模型中引入特定几何构型来实现特定晶体结构的生成.而 Ye等^[24]最近提出的 PODGen方法则通过构建条件生成框架,将通用生成模型与多个性质预测模型相结合,在采样过程中直接引入目标性质的条件概率分布,从而显著提升生成满足目标性质晶体结构的效率和成功率.此外,诸如 FiLM 调控机制^[25]、结构-性质联合优化^[26]等技术也被广泛用于提升目标性质的控制精度与生成样本的多样性.这些方法不仅拓展了生成模型在材料领域的适用性,也推动了“从性质到结构”的设计范式加速落地.

然而,当前基于目标导向的生成方法仍存在一些亟待突破的技术瓶颈.一方面,目标性质与结构隐空间之间的耦合关系建模尚不充分,容易导致生成结果性质偏离目标,或牺牲结构的物理合理性.另一方面,在小样本甚至极少样本目标值区域,如何利用迁移学习^[27]、物理先验^[28]或主动学习^[29]等手段提升生成准确性,仍是需要深入研究的问题.同时,由于性质标签通常来源于计算或实验,其

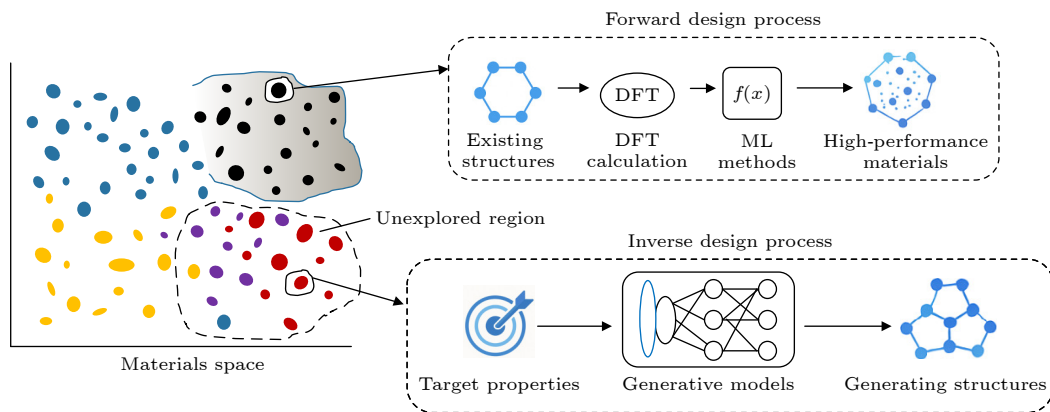


图1 数据驱动材料设计的两种范式:正向设计与逆向设计.图中黑色区域表示材料化学空间中已探索的区域,彩色点所在的区域表示未探索区域,黄色点所在的区域表示潜在的高性能材料区域;正向设计流程(上)从已知材料结构出发,经第一性原理与机器学习预测逐步筛选候选材料;逆向设计流程(下)从目标性质出发,通过生成模型直接生成候选结构

Fig. 1. Two paradigms of data-driven materials design: forward design and inverse design. The black region in the figure represents the explored part of the materials chemical space, the region with colored dots indicates the unexplored space, and the region with yellow dots denotes the potential high-performance materials. The forward design process (top) starts from known material structures and gradually screens candidate materials via first-principles calculations and machine learning predictions; the inverse design process (bottom) starts from target properties and directly generates candidate structures through generative models.

噪声、误差乃至不完备性也可能对模型的训练与材料的生成带来不可忽视的影响。

综上所述, 目标性质导向的材料生成已成为材料逆向设计研究中的重要方向. 本文将围绕几类经典的材料生成模型 (CDVAE^[30], MatGAN^[31], MatterGen^[25]) 来系统梳理近年来该领域的主要研究进展, 重点关注生成模型中目标性质导向机制的设计方法与实际应用表现, 深入分析其在样本数据质量、目标性质引导效果、生成结构多样性等方面的表现, 并讨论未来该领域的发展潜力与挑战。

2 用于材料领域的生成模型

生成模型最初主要应用于图像和文本等领域^[12-16], 近年来逐渐被引入材料科学以实现结构的逆向生成. 然而, 与图像生成相比, 材料结构生成需满足晶体对称性、化学合理性与物理稳定性等严格约束, 这对模型设计提出了更高要求. 为此, 研究者对传统生成模型进行了适配改造, 使其具备生成合理晶体结构的能力. 目前常用于材料生成的生成模型包括融合了变分自编码器 (variational autoencoder, VAE)^[12] 与扩散过程的 CDVAE^[30]、基于对抗机制的 MatGAN^[31] 以及近年来发展的基于显式扩散建模的 MatterGen^[25]. 本节将简要介绍这些模型在材料结构生成任务中的基本原理与特点。

2.1 CDVAE

生成模型是一类能够学习样本中的数据分布并由此生成新样本数据的建模方法, 涵盖了从传统的概率图模型到现代的深度神经网络. 在深度学习框架下, 主流的生成模型通常可分为两大类: 一类是基于显式概率建模的方法, 如自回归模型 (autoregressive model)^[16]; 另一类是基于隐变量建模的方法, 如 VAE^[12]、生成对抗网络 (generative adversarial network, GAN)^[13]、归一化流 (normalizing flow)^[15]. 其中, VAE 由于训练稳定、隐空间可解析等优点, 被广泛应用于材料结构生成任务, 能够在隐空间中学习材料结构的概率分布并实现新结构的生成。

然而, 传统的 VAE 在晶体材料生成中的直接应用存在明显不足, 如生成结构往往缺乏周期性、稳定性差、物理约束弱, 难以满足实际材料设计的需求. 为了克服这些限制, Xie 等^[30] 提出了 CDVAE

(crystal diffusion variational autoencoder), 一种面向晶体材料结构生成的新型深度生成模型, 其核心目标是在严格的物理约束条件下, 实现稳定周期性材料结构的生成. 该模型首次将 VAE 与扩散过程^[14] 相结合, 引入物理归纳偏置, 并通过端到端训练方式学习晶体材料的潜在结构分布, 是材料逆向设计任务中的重要突破。

与早期的图像生成^[32] 或分子结构生成模型^[33,34] 相比, 晶体材料生成任务面临多项独特挑战. 首先, 晶体结构具有三维周期性, 生成过程中不仅要确保原子间的物理稳定性, 还需满足空间群对称性、平移/旋转/周期等不变性; 其次, 材料数据远比图像和分子数据稀缺, 特别是在真实晶体结构的三维坐标与晶格信息方面. 因此, 直接借用传统生成模型难以满足高精度材料设计的实际需求。

针对上述问题, CDVAE 采用了多层次的建模策略以确保生成结构的物理合理性和稳定性: 模型的编码器和解码器均采用适配晶体周期性的图神经网络 (periodic graph neural network, PGNN), 以保证对结构周期边界的准确处理, 并在隐空间采样过程中使用了 VAE 框架中的重参数化技术^[35]. 在隐空间采样后, 模型通过一个“性质聚合预测器”先对材料的组成、晶格参数与原子数进行预测, 随后再以这些结构信息作为先验, 对原子类型和三维坐标进行初始化. 在解码阶段, CDVAE 引入得分匹配网络 (score-based decoder)^[36], 通过扩散建模方式引导结构从高噪声状态逐步演化为低能稳定态. 该过程使用朗之万动力学^[37] 逐步修正原子位置并更新原子类型, 确保生成结构的局部稳定性与整体物理合理性 (见图 2(a)).

值得注意的是, CDVAE 的去噪目标不仅仅是从随机结构还原已有结构, 更是通过模拟从隐空间中采样初始结构, 再通过物理感知的得分网络将其转化为稳定材料, 从而实现真正意义上的从噪声中生长出晶体结构. 这一过程本质上引入了近似的弹性力场约束, 保证了生成材料的稳定性。

表 1 展示了不同的生成模型在不同数据集上的测试结果, 包括生成结构的结构有效性 (Validity)^[38,39]、结构覆盖率 (COV)^[40,41] 和性质分布偏差 (property statistics). 在多个公开数据集 (Perov-5^[42], Carbon-24^[43] 和 MP-20^[9]) 上的实验评估表明, CDVAE 在生成结构的物理合理性、新颖性等方面均优于现有的主流生成模型 (如 FTCP^[44],

Cond-DFC-VAE^[38], G-SchNet^[45] 等). 特别是在复杂三维无机材料体系 (如 MP-20) 中, CDVAE 展

现出更高的生成结构多样性与物理合理性, 具备显著的性质引导潜力.

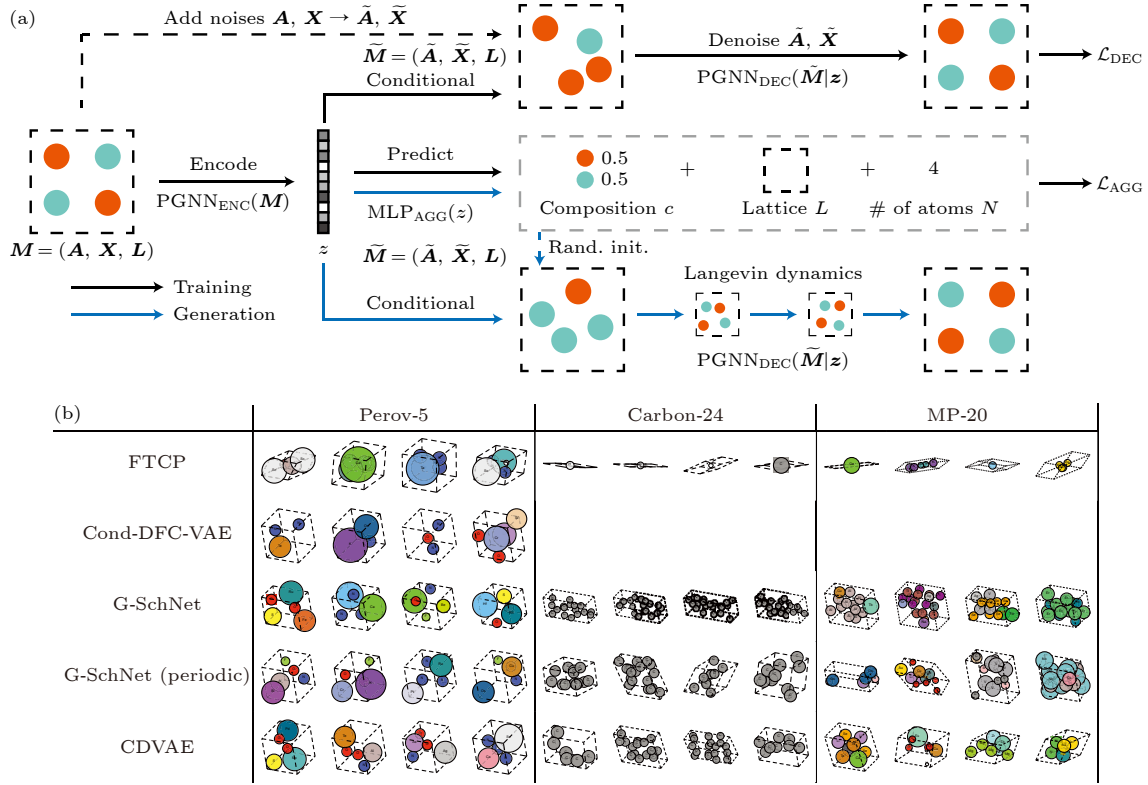


图 2 (a) CDVAE 模型的整体框架流程^[30]. 模型的输入为晶体的图结构表示, 通过周期性图神经网络 (PGNN) 进行编码, 结合性质聚合预测器预测结构聚合信息, 并由得分网络^[36] (图中使用 PGNN 进行解码的部分) 结合朗之万动力学 (Langevin dynamics)^[37] 生成稳定的三维晶体结构; (b) 不同的生成模型在不同的数据集上经过随机采样生成的晶体结构图

Fig. 2. (a) Framework of the CDVAE^[30]. The model takes the crystal graph representation as input, encodes it using a periodic graph neural network (PGNN) and predicts aggregated structural attributes through a property aggregation predictor. A score-based network^[36] (implemented using PGNN in the decoder) combined with Langevin dynamics^[37] is then used to generate stable three-dimensional crystal structures. (b) Crystal structures randomly generated by different generative models on different datasets.

表 1 不同的生成模型在不同的数据集上生成的结构在多个评估指标下的性能对比^[30]
Table 1. Comparison of generative model performance across datasets and evaluation metrics^[30].

Method	Data	Validity/%		COV/%		Property statistics		
		Struc.	Comp.	R.	P.	ρ	E	# elem.
FTCP	Perov-5	0.24	54.24	0.10	0.00	10.27	156.0	0.6297
	Carbon-24	0.08	—	0.00	0.00	5.206	19.05	—
	MP-20	1.55	48.37	4.72	0.09	23.71	160.9	0.7363
Cond-DFC-VAE	Perov-5	73.60	82.95	73.92	10.13	2.268	4.111	0.8373
G-SchNet	Perov-5	99.92	98.79	0.18	0.23	1.625	4.746	0.03684
	Carbon-24	99.94	—	0.00	0.00	0.9427	1.320	—
	MP-20	99.65	75.96	38.33	99.57	3.034	42.09	0.6411
P-G-SchNet	Perov-5	79.63	99.13	0.37	0.25	0.2755	1.388	0.4552
	Carbon-24	48.39	—	0.00	0.00	1.533	134.7	—
	MP-20	77.51	76.40	41.93	99.74	4.04	2.448	0.6234
CDVAE	Perov-5	100.0	98.59	99.45	98.46	0.1258	0.0264	0.0628
	Carbon-24	100.0	—	99.80	83.08	0.1407	0.2850	—
	MP-20	100.0	86.70	99.15	99.49	0.6875	0.2778	1.432

2.2 MatGAN

MatGAN(materials generative adversarial network)^[31]作为一种基于生成对抗网络(generative adversarial network, GAN)^[13]的生成模型,可用于探索无机材料组成空间并生成相应的材料结构.其主要是通过学习已有材料成分分布,自动生成化学上新颖、结构上合理的潜在材料组合,从而为材料筛选和发现提供候选空间扩展的有效路径.

材料组成空间具有高度稀疏性和高维离散性^[46],传统的生成模型难以精准地学习材料组分共性规律.除此之外,合理的材料成分也需满足电中性^[47]、电负性匹配^[48]等隐含的物理规则,若不加以约束容易生成大量无效结构.为了应对这些挑战,MatGAN构建了包含生成器(generator, G)与判别器(discriminator, D)两部分的对抗学习系统(见图3(a)).模型以隐变量作为输入,生成器通过多层反卷积神经网络构造材料成分稀疏矩阵表示,再由判别器判定生成样本与真实样本的分布一致性.为提升训练稳定性,MatGAN引入 Wasserstein GAN(WGAN)^[49]架构,以 Wasserstein 距离优化对抗损失,从而缓解传统 GAN 训练中梯度不稳定的问题.

在生成结构新颖性和覆盖率方面,MatGAN 展现出显著优势.如图3(b)所示,通过 t-分布随机邻域嵌入(t-distributed stochastic neighbor embedding, t-SNE)^[50]降维可视化,MatGAN 生成的结构(蓝点)

广泛分布在训练样本(绿点)与留出样本(红点)之外,覆盖了更广泛的化学空间区域,说明其不仅具备复现训练分布的能力,还可以主动探索未知组分组合.在3个主流材料数据库(OQMD^[8], Materials Project^[9]和 ICSD^[7])上的实验表明,MatGAN 生成的200万个样本中,有超过98%是未出现在训练集中的新材料组合(见表2);其中对 ICSD 留出样本的复现率达到60.13%,而新样本数高达198万,这充分验证了 MatGAN 强大的材料生成能力和对于未知材料的探索能力.

表2 MatGAN 在 OQMD^[8], Materials Project^[9]与 ICSD^[7]三个数据库上的结构复现与新颖性统计^[31]

Table 2. Statistical results of structural recovery and novelty of MatGAN on the OQMD^[8], Materials Project^[9], and ICSD^[7] databases^[31].

	GAN-OQMD	GAN-MP	GAN-ICSD
Training sample #	251368	57530	25323
Leave out sample #	27929	6392	2813
Generated sample #	2000000	2000000	2000000
Recovery of training samples/%	60.26	47.36	59.54
Recovery of leave out/%	60.43	48.82	60.13
New samples	1831648	1969633	1983231

2.3 MatterGen

尽管基于隐空间建模的 CDVAE 在晶体结构生成方面取得了重要进展,但其生成过程仍存在一

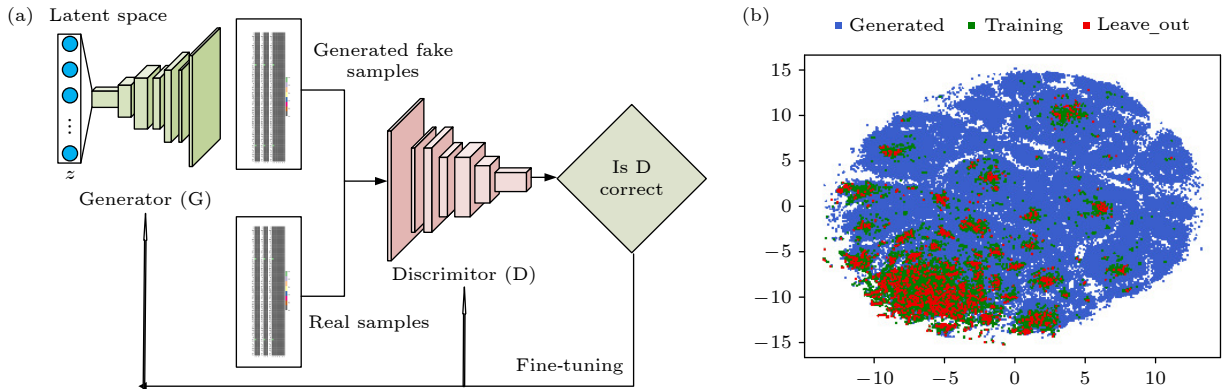


图3 (a) MatGAN的生成对抗网络架构示意图^[31],模型通过生成器从隐空间中采样隐向量,生成新的材料组成样本,判别器判断样本是否来源于真实数据分布,通过对抗训练优化生成效果;(b)使用 t-SNE^[50]方法对材料成分分布可视化降维图,图中蓝色点为 MatGAN 生成的样本,绿色点为训练集样本,红色点为保留未参与训练的 leave-out 样本

Fig. 3. (a) Schematic diagram of the generative adversarial network architecture of MatGAN^[31]. The model samples latent vectors from the latent space via the generator to produce new material compositions, while the discriminator determines whether the samples come from the real data distribution, thereby optimizing the generator through adversarial training. (b) t-SNE^[50] visualization of the material composition distribution after dimensionality reduction. Blue points represent samples generated by MatGAN, green points denote training set samples, and red points indicate leave-out samples not used during training.

定的局限性. 一方面, 隐空间的表达能力受限于编码器训练质量, 容易导致生成样本分布集中、结构多样性不足; 另一方面, 由于生成过程依赖隐变量, 模型在探索广阔的化学空间时可能面临覆盖率受限的问题. 为进一步突破传统基于高通量筛选^[51]、随机结构搜索^[52]等方法在新材料发现中的局限, Zeni 等^[25]提出了 MatterGen, 一种基于扩散模型^[14]的无机晶体结构生成框架, 该框架可在学习大量稳定材料分布规律的基础上高效探索广阔的化学空

间, 极大地提升了新型材料发现的可能性.

MatterGen 针对无机晶体材料的复杂周期性与对称性特征, 设计了分别针对原子类型、原子三维坐标与晶格参数的独立扩散过程 (见图 4(a)). 首先在正向过程 (forward process) 中, 稳定的材料结构 (原子类型 A_0 , 原子三维坐标 X_0 , 晶格参数 L_0) 逐步被加入噪声, 最终破坏为随机材料; 在反向过程 (reverse process) 中, 得分网络 (score network)^[53]从噪声中恢复出稳定结构. 为了精确建模材料的空

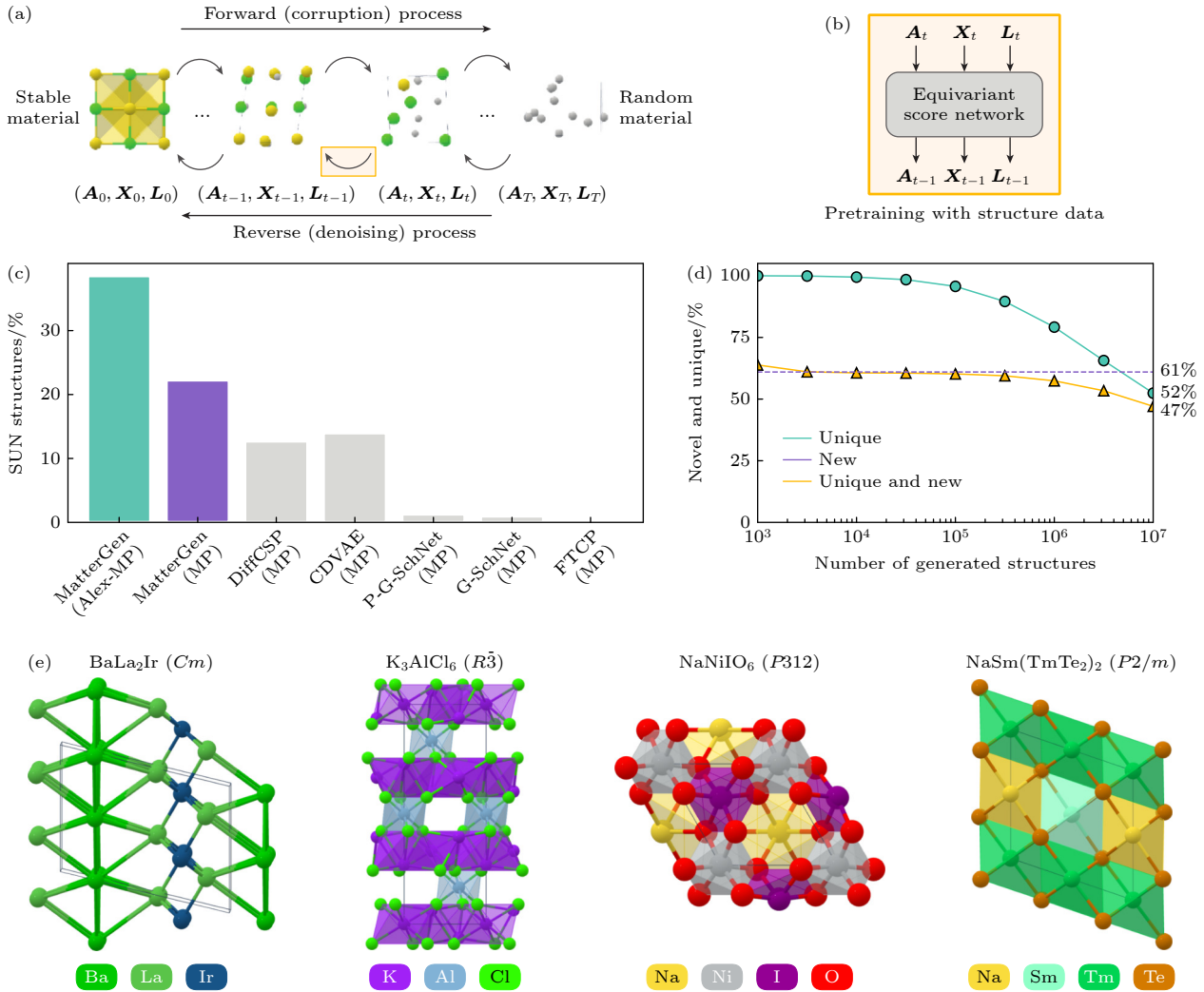


图 4 (a) MatterGen 的扩散建模流程示意图^[25], 正向过程 (forward process) 通过逐步加噪将稳定结构材料 (A_0, X_0, L_0) 转化为完全随机结构 (A_T, X_T, L_T), 反向过程 (reverse process) 则利用得分网络 (score network)^[53]引导去噪, 逐步重建出稳定晶体结构; (b) 等变性得分网络 (equivariant score network) 结构示意图; (c) 不同生成模型在 Materials Project^[9]数据集上新颖稳定生成材料的比例对比图; (d) MatterGen 生成材料在不同采样规模下的新颖性与唯一性统计图; (e) MatterGen 生成的部分代表性新型晶体结构图

Fig. 4. (a) Schematic diagram of the diffusion modeling process of MatterGen^[25]. The forward process gradually corrupts a stable material structure (A_0, X_0, L_0) into a completely random structure (A_T, X_T, L_T) by adding noise step-by-step, while the reverse process leverages a score network^[53] to guide denoising and progressively reconstruct a stable crystal structure. (b) Schematic diagram of the equivariant score network structure. (c) Comparison of the proportion of novel and stable structures generated by different models on the Materials Project^[9] dataset. (d) Statistical analysis of the novelty and uniqueness of materials generated by MatterGen at different sampling scales. (e) Representative examples of novel crystal structures generated by MatterGen.

间对称性与周期性, 如图 4(b) 所示, MatterGen 在得分网络中引入了等变性设计 (equivariant score network), 使模型在生成过程中自然保持晶体空间对称性, 大幅减少了无物理意义的结构。

在训练过程中, MatterGen 利用来自 Materials Project^[9] 和 Alexandria^[54] 数据库的约 60 万条稳定晶体数据进行预训练, 仅基于结构信息学习稳定材料的分布规律. MatterGen 在无条件生成任务中展现出卓越的性能. 如图 4(c) 所示, MatterGen 在生成新颖且稳定的材料比例上显著优于其他生成模型 (如 DiffCSP^[55], CDVAE^[30], FTCP^[44] 等), 其中基于 Alexandria+MP 联合训练的版本能够生成超过 30% 的新颖稳定结构. 即使是在大规模生成 (上百万结构) 条件下, MatterGen 仍能保持超过 50% 的新颖或唯一结构比例 (见图 4(d)), 验证了其强大的化学空间拓展能力. 除此之外, MatterGen 生成的晶体结构涵盖了不同的元素组合与晶体对称性体系 (见图 4(e)), 进一步证明了其生成结果在结构合理性与多样性方面的优越表现。

3 实现目标性质导向的模型与方法

尽管近年来生成模型在材料结构生成中取得了显著进展, 但“如何将目标性质有效引入生成过程”仍是当前逆向设计面临的关键挑战之一^[56]. 传统生成模型主要聚焦于生成结构的合理性与多样性, 然而在实际材料设计任务中, 研究者们更关心的是能否根据设定的性质需求 (如带隙、形成能、磁性等) 直接生成具有对应目标性质的全新材料. 这一问题涉及生成模型能否精确建模结构与性质之间复杂的映射关系。

由于材料结构与性质之间往往存在高度非线性、多尺度耦合与对原子尺度细节敏感依赖的问题^[57,58], 单纯依赖原始生成模型往往难以获得良好的性质一致性或精准的控制效果. 针对这一难题, 研究者们提出了多种不同的实现策略, 探索如何将目标性质嵌入模型的编码、解码或优化过程中, 以引导生成结构更符合预设的性质需求, 其中包括在编码器和解码器中引入目标性质向量进行建模 (如 Con-CDVAE)^[22]、融合结构掩码与性质掩码实现多重约束下的精细引导 (如 SCIGEN)^[23]、在扩散模型中通过适配器微调与引导机制实现性质控制 (如 MatterGen)^[25] 以及在生成阶段对隐变量进

行梯度优化以靠近目标性质的方法 (如隐变量优化 CDVAE)^[26]. 本节将逐一介绍这些策略的具体实现方式与作用效果, 并分析它们在引入目标性质导向机制方面的优劣与适用范围。

3.1 Con-CDVAE: 引入目标性质向量的编码-解码结构

在材料逆向设计中, 如何实现从目标性质出发生成对应结构一直是结构-性质耦合任务的核心问题^[59-61]. 传统生成模型如 CDVAE^[30], 虽然能在周期性边界条件下生成具有合理三维结构的无机晶体, 但其结构生成过程本质上是从无条件隐空间分布中采样, 缺乏对目标性质的控制机制. 因此, 在生成新结构的同时, 模型往往无法确保生成材料具备所需的物理性质, 尤其是在如带隙、形成能等对原子局域结构高度敏感的目标上, 结构与性质之间存在显著偏离^[62,63].

为解决上述问题, Ye 等^[22] 提出了 Con-CDVAE (conditional crystal diffusion variational auto-encoder) 模型, 通过在 CDVAE 框架中引入目标性质向量作为条件信息, 将目标导向的生成任务形式化为目标性质条件下的结构分布建模问题, 从而实现在保持晶体物理合理性与周期对称性的基础上, 生成具有预设性质的晶体结构。

如图 5(a) 所示, Con-CDVAE 在编码器和解码器中分别引入性质调控机制, 通过训练阶段构建目标性质与隐空间之间的对齐过程. 性质调控的核心机制来自一个独立的性质嵌入模块 (prop_{emb}), 可将连续性质 (如形成能、带隙) 和离散性质 (如晶体系统) 分别编码为固定维度的向量来进行嵌入. 其中, 连续性质通过高斯基展开再通过多层感知机 (multi-layer perceptron, MLP) 映射为隐向量, 离散性质则通过分类嵌入和 MLP 映射处理. 目标性质嵌入向量随后与编码器产生的隐变量 z 拼接, 构成含有目标性质信息的隐变量 z_{cond} 来传递给解码器进行解码. 这一设计确保了性质信息从编码阶段开始影响隐空间的分布建模过程, 从而使生成结构在隐空间中具备响应目标性质的能力。

为了增强生成结构与目标性质之间的一致性, Con-CDVAE 进一步引入了一个性质预测器模块 (predictor), 在训练阶段通过反向预测的方式对隐变量 z 解码所产生的性质值进行回归, 并与输入的真实性质进行对比, 构建额外的监督信号. 这一

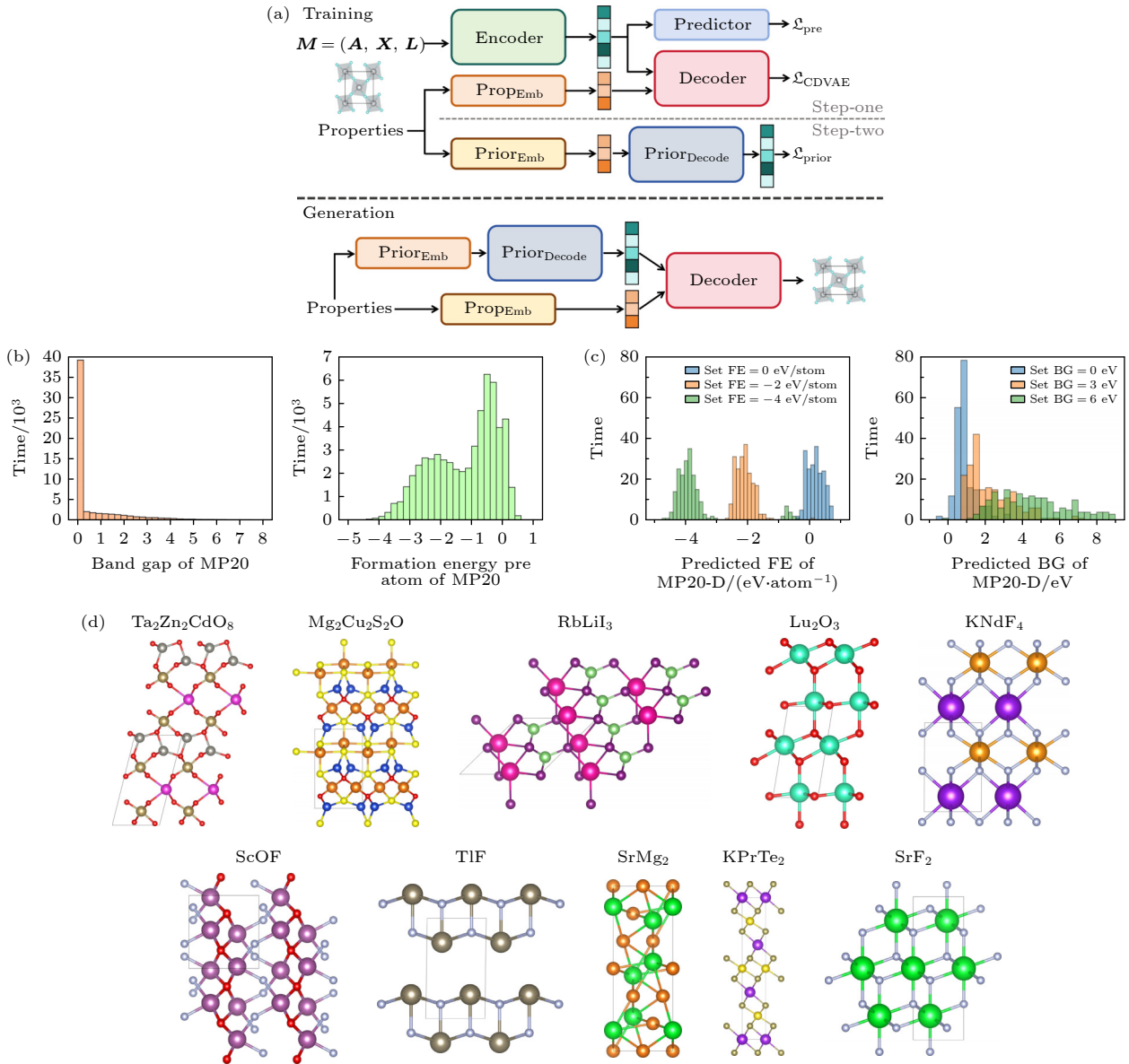


图 5 (a) Con-CDVAE 的训练与生成流程图^[22]. 训练阶段包括结构与目标性质的联合编码 (encoder)、性质预测器 (predictor) 和条件解码 (decoder), 同时训练先验 (prior) 模块^[14,36] 用于实现从目标性质生成隐变量 z . 生成阶段中, 输入目标性质后通过先验模块生成隐变量 z , 并与通过 $prop_{emb}$ 模块得到的性质嵌入向量进行拼接后得到含有目标性质信息的隐变量 z_{cond} , 将 z_{cond} 送入解码器后生成晶体结构. (b) MP-20^[9] 训练数据集中材料的带隙分布与形成能分布图. (c) Con-CDVAE 在设定的目标性质下的生成材料的性质分布图. (d) Con-CDVAE 生成的部分代表性新型晶体结构图

Fig. 5. (a) Training and generation workflow of Con-CDVAE^[22]. The training phase includes joint encoding of crystal structures and target properties via the encoder, a property predictor, and conditional decoding through the decoder. Simultaneously, a prior module^[14,36] is trained to enable the generation of latent variable z from target property inputs. In the generation phase, a target property is first provided, from which a latent variable z is generated using the prior module. This latent variable is then concatenated with a property embedding vector obtained from the $prop_{emb}$ module to derive the latent variable z_{cond} conditioned on the target property, which is passed to the decoder to generate crystal structures. (b) Distribution of band gap and formation energy in the MP-20 training dataset^[9]. (c) Property distributions of materials generated by Con-CDVAE under different target property settings. (d) Representative novel crystal structures generated by Con-CDVAE.

机制强化了隐空间中的结构-性质耦合程度, 使得具有相似目标性质的材料在隐空间中趋于聚集, 从而有助于提升模型在目标值附近的生成准确率.

在模型训练方面, Con-CDVAE 采用了两阶段

的训练机制. 第 1 阶段固定输入晶体结构与目标性质, 训练 Con-CDVAE 主模型 (encoder, decoder, $prop_{emb}$, predictor) 以学习目标性质条件下的隐变量生成与解码机制; 第 2 阶段则在冻结上述模块的

基础上, 单独训练一个基于 DDPM(denoising diffusion probabilistic model)^[14] 的 prior 模块^[14,36], 学习如何仅从目标性质向量出发生成合理的隐变量 z . 该模块通过在训练阶段引入逐步加噪与反向去噪过程, 模拟隐空间上的扩散采样机制, 完成训练后, 通过指定目标性质向量就可以利用 prior 从随机噪声中反向生成隐向量, 再由 decoder 生成对应晶体结构, 实现真正意义上的从性质生成结构的逆向生成能力.

结果显示, Con-CDVAE 在 MP-20^[9] 数据集上实现了目标性质导向的生成, 在目标形成能与目标带隙的控制精度方面均表现出一定的响应性和集中度. 如图 5(b), (c) 所示, 在设置目标形成能为 0, -2 和 -4 eV/atom 的条件下, 模型生成结构的形成能分布呈现出明显偏移趋势, 且在训练集中含目标性质训练数据较少的情况下 (如形成能为 -4 eV/atom) 也能做到不错的引导效果, 表明模型成功实现了一定的性质引导能力. 同时, 如图 5(d) 所示, 模型生成的代表性结构也展示出良好的空间群对称性与三维几何合理性, 进一步验证了 Con-CDVAE 在结构质量与性质对齐方面的综合性能.

尽管 Con-CDVAE 已在目标性质建模方面迈出了重要一步, 但该模型仍存在一些待优化之处. 例如, 虽然 Con-CDVAE 在形成能这种连续性质的生成引导方面展现出良好效果, 能够有效控制生成材料的形成能分布靠近预设目标, 但对于如带隙这类对局域原子排布极为敏感的性质, 其引导效果仍不够显著 (见图 5(c)), 生成结构的目标一致性相对较弱. 此外, 模型也比较依赖于训练集中的性质数据分布, 对于含带隙这种复杂性数据较少的区域的引导效果可能不会特别明显. 因此, 未来研究可考虑引入多尺度性质图建模^[64] 或基于 Transformer 的条件控制结构^[65] 等一系列方法来增强模型对复杂或离散目标性质的生成引导能力.

3.2 SCIGEN: 融合结构约束与引导的 DiffCSP 扩展模型

现有的生成模型如 CDVAE^[30], UniMat^[66] 与 DiffCSP^[55] 虽已展现出在材料空间探索的巨大潜力, 但大多仍依赖于数据库分布进行无条件生成, 难以主动控制生成结构的几何构型或磁性拓扑. 因此, 为了应对当下材料生成模型难以实现特定几何约束

下材料设计的挑战, Okabe 等^[23] 提出了 SCIGEN (structural constraint integrated generative model), 一种在 DiffCSP^[55] 扩散生成模型基础上进一步改进的生成模型.

SCIGEN 的核心目标是融合结构合理性与目标性质一致性, 实现可控几何结构的晶体生成. 相比于仅依赖空间群对称性建模的传统方法^[65,67], SCIGEN 引入了结构掩码机制, 通过明确指定几何构型作为约束, 在多步扩散去噪过程中引导结构演化来生成目标几何构型的结构. 这种方式不仅提升了结构生成的物理稳定性与可控性, 更使模型具备了面向特定材料 (如量子磁性材料) 进行定向设计的能力.

如图 6(a) 所示, SCIGEN 通过设定好的几何规则 (如 Kagome, Triangular 或 Lieb 晶格) 作为先验输入, 通过结构中部分已知的晶胞信息、原子数、键长与对称性构建结构掩码 (structure mask). 在多步扩散去噪过程中, 模型以掩码作为控制信号, 对约束区域保持不变, 仅对未约束区域进行生成, 从而实现对结构拓扑的精细控制. 图 6(b) 进一步展示了 SCIGEN 的掩码机制: 在原子类型、原子坐标与晶胞参数 3 个结构特征维度分别构建掩码矩阵, 以确保在不同生成阶段对不同类型的几何约束加以保留或采样. 该机制首次在晶体结构生成任务中实现了对几何构型的直接控制, 是 SCIGEN 模型的核心创新之一.

为了验证其结构控制能力与性质引导效果, SCIGEN 在 Lieb 晶格这一代表性的二维几何结构上进行测试. 图 6(c) 为理想的 Lieb 晶格拓扑, 如图 6(d) 所示, SCIGEN 在该几何约束下生成了三元材料 Ce_3NaTb_4 , 其原子布局保留了典型的 Lieb 结构, 并具备合理的周期性与键长, 其能带结构如图 6(e) 所示, 显示出接近费米能级的平带特征, 预示其可能具有非平庸的磁性或拓扑态行为^[68].

虽然 SCIGEN 在几何结构控制方面展现出显著优势, 尤其在生成 Lieb, Kagome 等具有潜在物理性能的晶格结构方面取得突破, 但其仍面临若干挑战. 首先, SCIGEN 在引入结构约束后, 需生成大量候选材料 (如近 800 万) 并经历多阶段筛选与 DFT 计算, 最终仅有少数材料满足目标结构与性质, 这导致生成效率较低. 其次, 模型性能高度依赖于结构初始化策略与充足的训练数据支持, 在面对罕见结构模式或低资源区域时, 泛化能力仍受限.

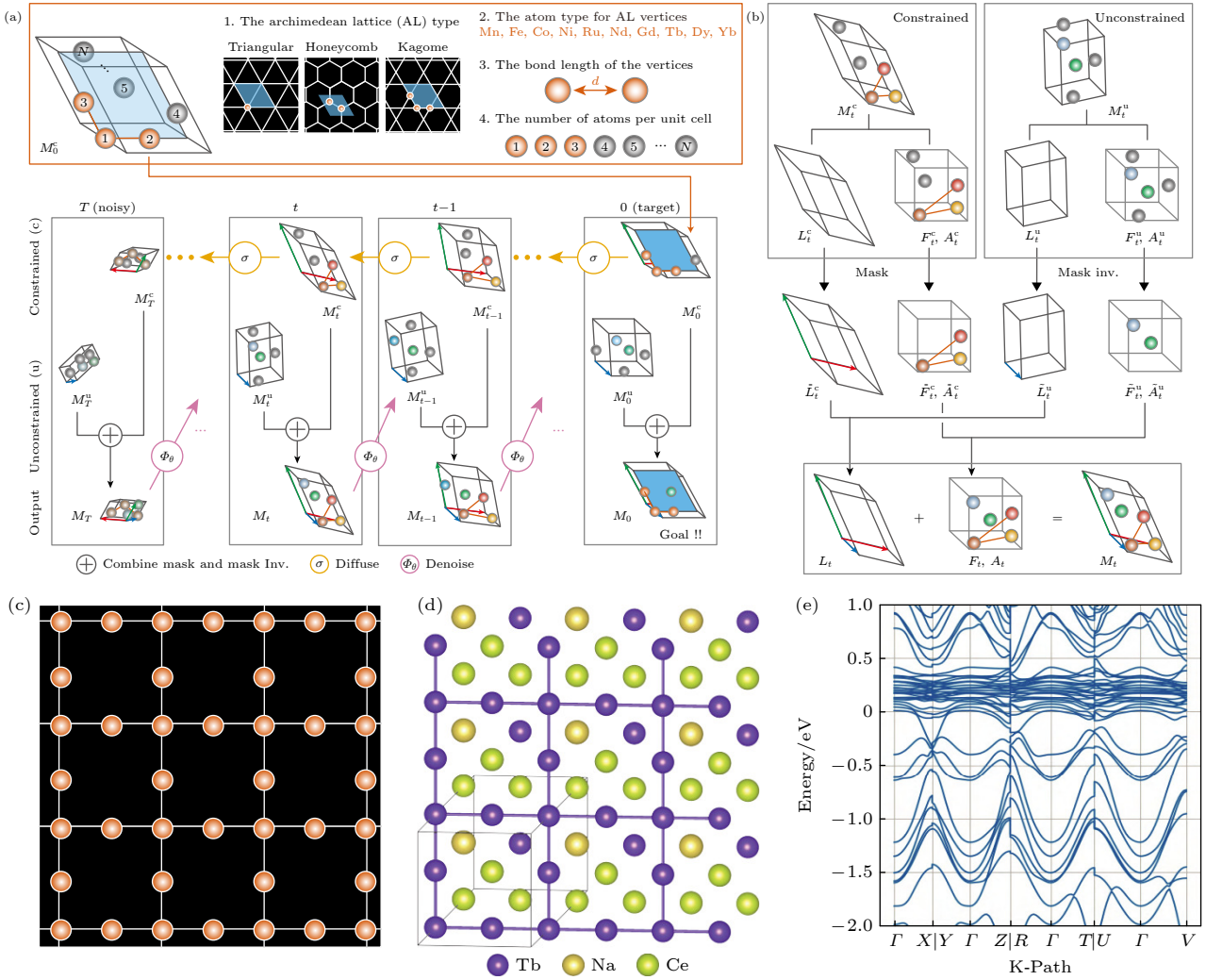


图 6 (a) SCIGEN 模型的材料结构生成过程示意图^[23], 该模型以几何规则为引导, 通过多步扩散与去噪过程, 结合掩码机制, 将具有几何约束的结构与待生成结构融合, 引导生成结构逐步逼近目标结构; (b) SCIGEN 的结构生成中几何约束与非约束结构的合并与掩码机制过程图; (c) Lieb 晶格的二维示意图; (d) SCIGEN 在 Lieb 几何约束下生成的代表性结构 Ce_3NaTb_4 ; (e) Ce_3NaTb_4 的能带结构图

Fig. 6. (a) Schematic illustration of the material structure generation process in the SCIGEN model^[23]. The model is guided by geometric rules and utilizes a multi-step diffusion and denoising process, combined with a masking mechanism, to integrate geometrically constrained structures with unconstrained ones, thereby gradually guiding the generation toward the target structure. (b) Diagram showing the merging of constrained and unconstrained structures and the masking mechanism used in SCIGEN during structure generation. (c) A 2D schematic of the Lieb lattice. (d) A representative structure, Ce_3NaTb_4 , generated by SCIGEN under Lieb lattice constraints. (e) Band structure of the generated Ce_3NaTb_4 material.

因此, 未来可通过结合稀疏正则化技术等一系列方法来减少冗余生成与筛选成本, 从而进一步对模型进行改进.

3.3 微调扩散模型 MatterGen: 适配器调控实现目标性质控制

MatterGen^[25] 通过引入适配器模块^[69], 可以对已训练好的模型进行微调, 从而进一步实现了对目标性质的精确调控. 该模块被插入到 MatterGen 的等变得分网络 (equivariant score network) 中,

通过接受性质向量 \mathbf{c} 的输入^[70] 从而动态调节得分网络在每一步去噪过程中的输出特征. 这个机制可以使模型在预训练基础上, 通过少量带标签数据微调模型部分参数, 实现从“生成合理结构”到“生成满足特定目标性质的结构”的功能跃迁.

图 7(a) 展示了 MatterGen 在微调阶段的结构框架: 模型以微调数据集中带有标签的目标性质 \mathbf{c} 作为输入, 让适配器模块影响等变得分网络的去噪过程, 从而实现生成样本在目标性质维度上的对齐. 通过对 MatterGen 施加不同的目标约束, 可使

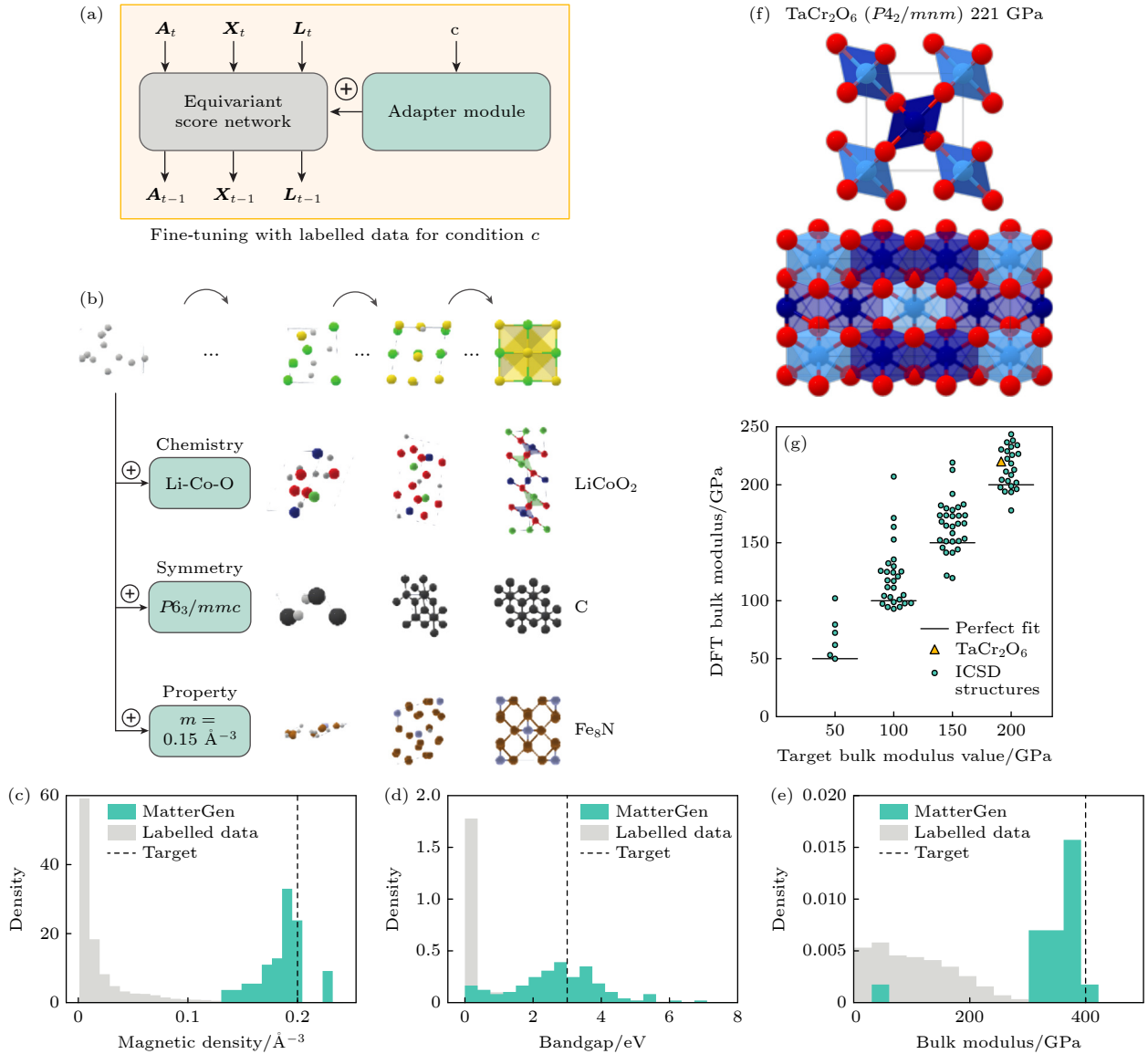


图 7 (a) MatterGen 微调阶段的模型结构示意图^[25], 目标性质 c 被送入适配器模块 (adapter module), 并与等变性得分网络 (equivariant score network) 的输出结合, 引导结构从随机状态逐步恢复至符合目标性质的晶体结构; (b) MatterGen 在不同目标性质条件下生成满足不同目标性质的晶体结构; (c)–(e) MatterGen 在目标性质分别为磁密度 (c)、带隙 (d)、体积模量 (e) 的条件下生成材料与对应微调数据集的材料性质分布图; (f) 实验上成功合成的 MatterGen 以体积模量 200 GPa 为目标性质生成的 TaCr_2O_6 结构图; (g) MatterGen 在目标性质为不同体积模量条件下生成材料的性质的 DFT 验证结果图

Fig. 7. (a) Schematic illustration of the MatterGen model during the fine-tuning stage^[25]. The target property condition c is fed into the adapter module and combined with the output of the equivariant score network to guide the generation process from random structures toward crystal structures satisfying the target property. (b) MatterGen-generated crystal structures under different target property constraints, demonstrating the model's ability to satisfy various design conditions. (c)–(e) Property distribution comparisons between MatterGen-generated materials and the fine-tuning dataset under target conditions of magnetic density (c), bandgap (d), and bulk modulus (e). (f) Structure of TaCr_2O_6 , successfully synthesized experimentally, generated by MatterGen with a target bulk modulus of 200 GPa. (g) DFT-validated property results of structures generated by MatterGen under different target bulk modulus values, illustrating the model's predictive accuracy and target controllability.

MatterGen 根据不同的约束来生成对应目标性质的材料 (见图 7(b)).

在具体生成效果上, MatterGen 展现出了对目标性质较好的集中性. 如图 7(c)–(e) 所示, 在分别设定目标性质为磁密度 0.2 \AA^{-3} 、带隙 3 eV、体积

模量 400 GPa 的条件下, MatterGen 生成的材料在目标性质附近的分布密度明显高于微调训练集中带标签样本的密度, 说明引入的适配器模块即使是在目标性质标签稀缺的条件下仍然具备良好的泛化性能和目标性质引导能力. 图 7(f) 则展示了

实验上成功合成的由 MatterGen 通过设定目标高体积模量作为目标性质来生成的 TaCr_2O_6 结构. MatterGen 通过设定不同的目标值来生成的材料在 DFT 验证下仍能保持与目标性质的高度一致性 (见图 7(g)), 这也充分体现了 MatterGen 在目标性质引导方面的强大性能.

然而, 尽管 MatterGen 在目标性质引导的生成方面取得了一定的突破, 但也存在一定的局限性. 该模型的预训练依赖于大量高质量的晶体结构数据, 训练所使用的数据量通常达到几十万甚至上百万级别, 除此之外, 即使是在微调阶段, 适配器所需要的带标签的微调数据集的数据量仍然高达上万甚至几十万条. 因此, 对于实验中较难获取或获取成本较高的复杂性质 (如热导率、载流子迁移率等) 的数据, 模型可能难以做到较好的性质引导效果. 这种对数据规模的依赖无疑限制了模型的泛化能力. 如何将 MatterGen 与数据增广^[71] 或小样本学习^[72] 等方法结合来降低数据门槛可能是未来模型优化的重要方向.

3.4 隐变量优化 CDVAE: 通过隐变量搜索实现目标性质导向生成

生成模型与性质预测模型的协同融合正推动材料逆向设计从“被动生成”走向“目标性质导向生成”的新阶段^[73,74]. 尽管当前的一些专门用于晶体材料结构生成的生成模型 (如 CDVAE^[30]) 在晶体结构生成中已展现出稳定性与周期性建模能力, 但其生成过程通常受限于训练数据的分布, 难以精准控制生成结构的目标性质. 在面对目标性质导向的材料设计任务时, 这类模型往往缺乏主动调控的能力, 难以在高维稀疏的材料空间中精准定位目标性质所在的区域. 针对该问题, Song 等^[26] 提出了隐空间优化策略 (latent space optimization), 即通过优化生成模型隐变量 z 实现对生成结构性质的直接调控. 该方法的关键在于不修改模型结构, 而是通过隐变量搜索这种后引导的方式来实现目标性质到结构的逆向映射.

Song 等^[26] 开发了 MAGECS (materials generation with efficient global chemical space search), 一种集成了预训练生成模型 CDVAE、监督图神经网络 (graph neural network, GNN)^[75] 和鸟群算法 (bird swarm algorithm, BSA)^[76] 的逆向设计

框架 (见图 8(a)), 用于高效探索全局化学空间并生成具有目标性质的材料结构. 在 CO_2 还原反应 (CO_2RR) 合金电催化剂的设计中, 该框架首先利用预训练好的 CDVAE 模型在结构空间中从隐变量 z 中生成大量晶体表面结构, 并枚举所有可能的吸附位点; 然后借助预训练好的 GNN 模型快速预测这些结构在目标中间体 CO 上的吸附能 (ΔE_{CO}) 并将其作为性质优化指标, 利用 BSA 通过对性质指标的迭代优化与反馈在隐空间中持续优化隐变量 z 的分布; 最后引导生成器生成满足目标性质的结构.

图 8(b) 展示了 3 组不同来源 (MAGECS 生成、原始数据库、CDVAE 生成) 的结构在目标性质范围 $|\Delta E_{\text{CO}} + 0.67| \leq 0.2 \text{ eV}$ 内的分布情况. 在生成数量从 100 扩展到 5 万的过程中, MAGECS 生成结构的分布始终显著偏向目标性质区域, 其均值逐渐逼近目标线, 最终有 34.4% 的生成结构落入目标范围, 远高于数据库 (13.3%) 与 CDVAE 生成 (13.8%) 的水平. 这表明了 MAGECS 框架在保持生成结构多样性的同时具备更强的目标性质引导能力. 最后在实验上成功合成了 MAGECS 生成的 5 种具有目标性质的结构 (见图 8(c)). 为了验证生成结构性质预测的可靠性与一致性, 研究者们对实验上成功合成的两种具有高法拉第效率的结构进行了电催化 CO_2 还原反应, 如图 8(d) 所示, Pd_5Sn_2 与 CuAl 结构在不同气氛 (Ar/CO_2) 下的电流密度曲线展现出良好的电催化性能, 其中 $\text{Pd}_5\text{Sn}_2\text{-CO}_2$ 体系在 -0.7 V 时达到显著的催化电流, 这一实验证据显示了通过隐变量搜索获得的生成结构在真实物理性质上的可靠性.

总体来看, 隐变量优化 CDVAE 通过在隐空间中的有效指导搜索, 为实现目标性质导向的晶体结构生成提供了新思路, 在加速材料发现、提升目标结构生成的成功率方面展现出巨大潜力. 然而, 该方法本质上依赖在隐空间中进行基于目标性质的偏导向搜索, 可能导致生成结构的物理合理性和新颖性下降^[77]. 此外, 由于性质预测器在训练分布之外的泛化能力有限, 当优化结果偏离训练分布或数据发生漂移时, 模型对生成材料的性质预测可能失真, 从而降低了整体方法的可靠性与实用性. 因此, 如何平衡目标性质精度与生成结构多样性并增强预测器对新颖结构的鲁棒性是该方法未来改进的重要方向.

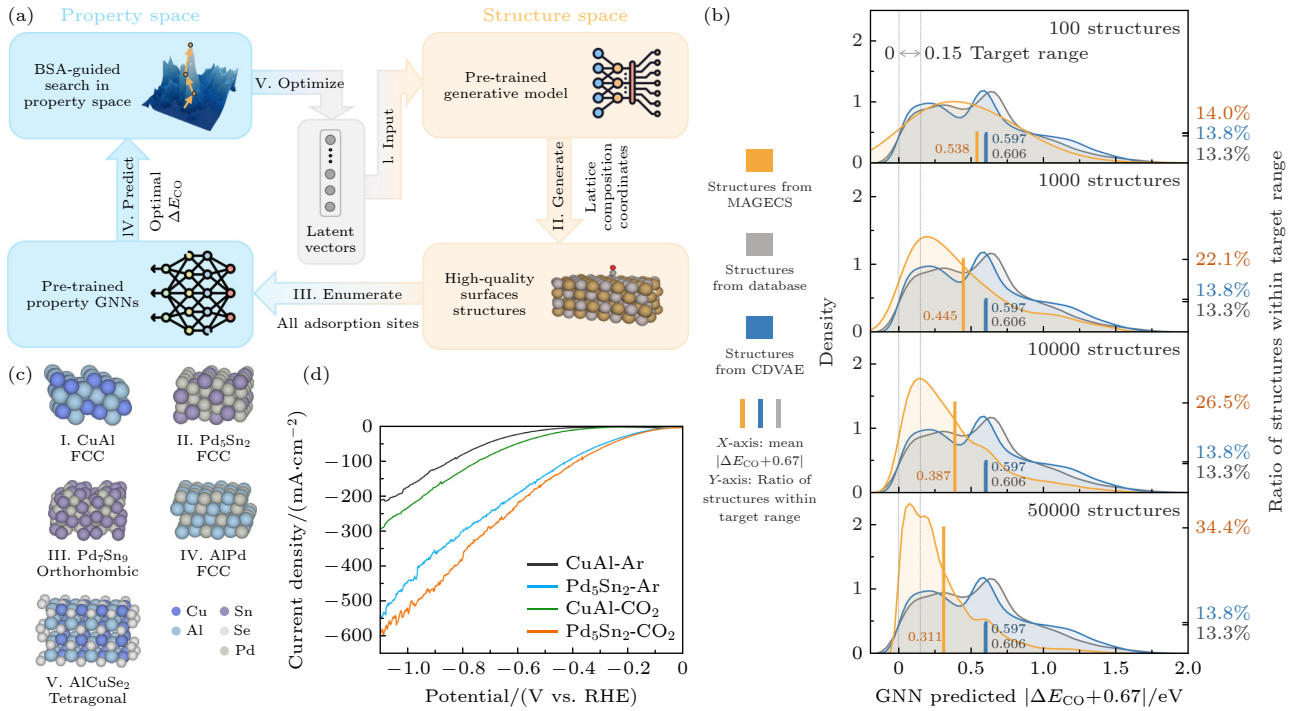


图 8 (a) MAGECS(materials generation with efficient global chemical space search) 框架结构图^[26], 该框架包括生成模型 CDVAE、监督图神经网络 GNN 和鸟群算法 BSA, 通过在隐空间中搜索最优隐变量以生成满足目标性质的材料结构; (b) MAGECS, CDVAE 生成结构和原始数据库结构在目标性质区间 $|\Delta E_{CO} + 0.67| \leq 0.2$ eV 的分布对比图; (c) 实验上成功合成的 MAGECS 生成的 5 种具有目标性质的结构; (d) CuAl 和 Pd₅Sn₂ 的电催化 CO₂ 还原反应性能曲线

Fig. 8. (a) Framework diagram of MAGECS (materials generation with efficient global chemical space search)^[26]. The framework integrates a generative model (CDVAE), a supervised graph neural network (GNN), and a bird swarm algorithm (BSA) to search for optimal latent vectors in the latent space and generate materials that meet target properties. (b) Comparison of the distribution of structures from MAGECS, CDVAE, and the original database within the target property range $|\Delta E_{CO} + 0.67| \leq 0.2$ eV. (c) Five target-property-aligned structures generated by MAGECS and successfully synthesized in experiments. (d) Electrocatalytic CO₂ reduction performance curves of CuAl and Pd₅Sn₂.

本节梳理了近年来 3 种主流的将目标性质机制引入生成模型的方法, 并以图示方式清晰呈现了它们在模型结构中的具体介入位置和作用路径 (见图 9)。目前主流的将目标性质机制引入生成模型的策略可分为 3 种类型: 第 1 类是以 Con-CDVAE^[22] 为代表的条件编码 (conditional-encoding) 方法, 通过在编码阶段将目标性质编码为隐变量 z' , 与结构编码得到的隐变量 z 进行拼接得到的带有目标性质信息的隐变量 z'' 输入到解码器中, 从而在生成初始阶段实现性质引导。第 2 类方法是以 SCIGEN^[23] 与 MatterGen^[25] 为代表的解码器微调 (fine-tuning decoder) 方法, 其核心思想是在解码器或去噪过程中引入目标性质引导模块, 如 adapter 层或微调的得分网络, 从而在生成过程中逐步实现向目标性质方向的演化。第 3 类方法是以 MAGECS^[26] 为代表的隐变量优化 (latent vector optimization) 方法, 通过在生成后的隐空间中对隐变量 z 进行目标性

质优化, 借助外部搜索算法或预测器在隐空间中迭代调整 z 之后得到符合目标性质分布的 z' , 输入到解码器中来实现不依赖模型结构修改的目标性质导向生成。

从性质引导效果来看, 3 类方法在不同性质上均实现了一定的引导效果。表 3 展示了 3 类方法在不同目标性质上的生成成功率。这里生成成功率 (success rate) 定义为生成材料中符合目标性质区间要求的材料所占总的生成材料的比例。表中的 5 个性质分别是形成能 (E_f)、带隙 (E_g)、磁密度 (M)、体积模量 (K) 和 CO 上的吸附能 (ΔE_{CO})。Con-CDVAE 这类方法在形成能这种连续性性质上的引导效果表现较好 (成功率可达到 40.1%), 但在带隙这类复杂性性质上的引导效果较弱 (仅 20.7%), 且这类方法比较依赖于训练数据分布。MatterGen 这类方法通过微调数据集相比于 Con-CDVAE 在带隙这个性质的引导效果上有了进一步的提升 (可达

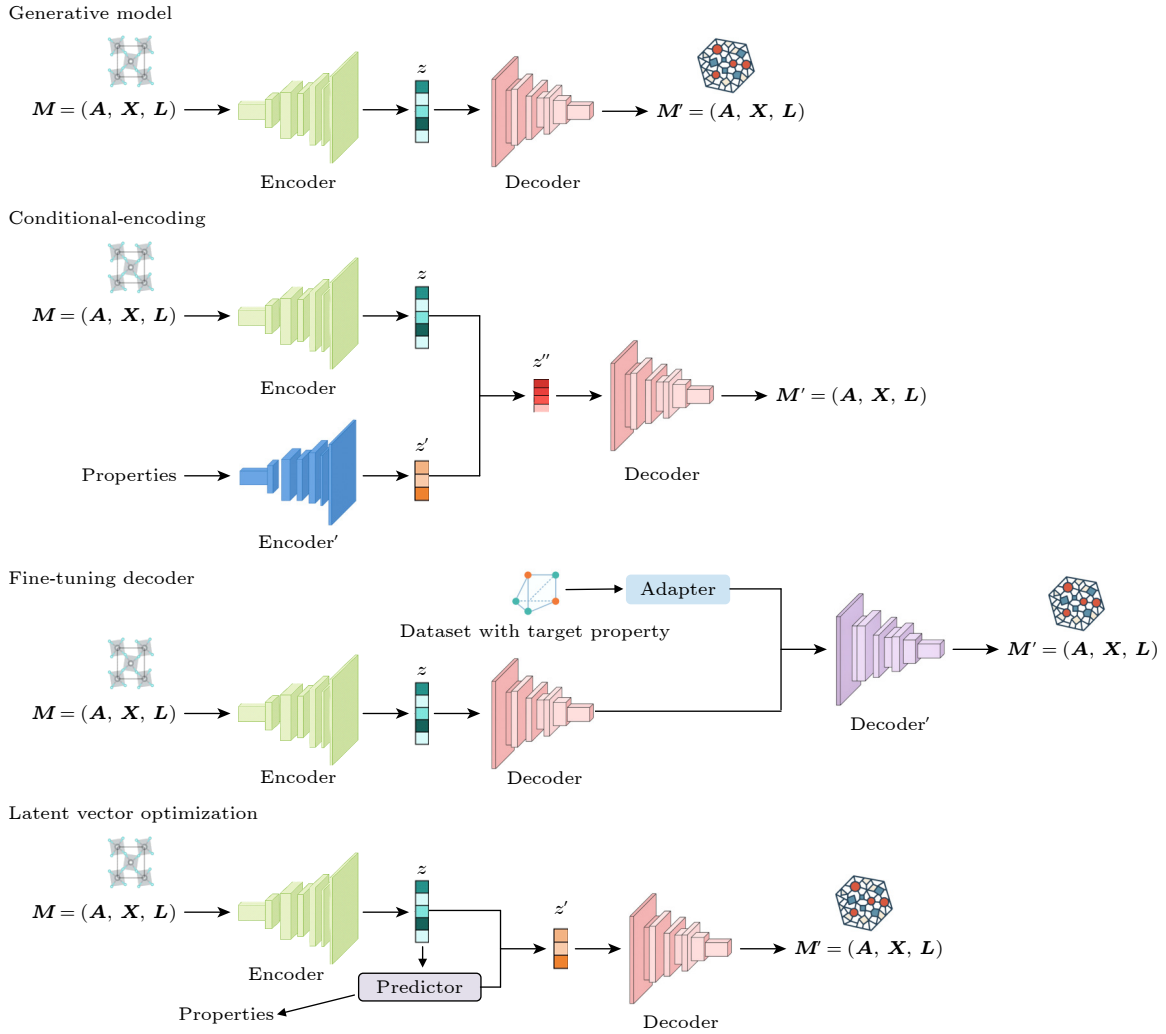


图 9 三种将目标性质导向机制引入生成模型的方法

Fig. 9. Three methods for incorporating property-guidance mechanisms into generative models.

表 3 三类性质导向方法在不同性质上的引导效果

Table 3. Effectiveness of three property-targeted methods across various properties.

Method	Data	Training sample	Success rate/%				
			E_f	E_g	M	K	ΔE_{CO}
Con-CDVAE ^[22]	MP-20	71665	36.3	18.4	—	—	—
	MP-40	108039	40.1	20.7	—	—	—
	OQMD	616412	38.8	19.3	—	—	—
MatterGen ^[25]	Alex-MP-20	607683	—	31.6	48.7	55.4	—
MAGECS ^[26]	GASpy	13000	—	—	—	—	34.4

31.6%), 且在磁密度 (48.7%) 和体积模量 (55.4%) 上都有较好的引导效果. 但该方法训练和微调都需要大量数据, 训练和微调数据集数据量一般都是几十甚至上百万的数量级. MAGECS 这类方法则通过隐变量优化实现了 CO 吸附能上的引导 (34.4%), 但这种方法生成结构的新颖性和稳定性容易受预测器精度的影响.

尽管这 3 类方法从不同层面实现了目标性质与结构生成的耦合, 但从“性质可控”到“精确生成”的路径仍面临诸多挑战. 首先, 现有方法大多集中在单目标性质控制, 对于复杂性质耦合 (如带隙与材料结构稳定性兼顾) 仍缺乏有效机制, 未来亟需发展多性质协同建模与权重平衡策略. 其次, 数据稀缺问题仍然制约模型的性能发挥, 尤其是在二维

材料、磁性材料等新兴体系中, 高质量标签数据的获取成本较高, 因此, 小样本学习^[72]、迁移学习^[27]与数据增广^[71]等策略仍具有重要研究价值。

最后, 如何保证生成的结构在化学上合理并具有实际可合成性是生成模型在未来发展的重大挑战. 从实验角度来看, 可合成性涉及前驱体选择、反应条件与动力学路径, 远超单纯的热力学稳定性考量; 从理论角度来看, 形成能、电荷平衡或动力学稳定性指标虽能提供初步判断, 但与实际可合成性之间仍存在明显差距. 近期的相关工作表明, 除了将空间群、电中性等物理约束嵌入生成流程^[78]外, 还可以通过机器学习和大模型方法直接预测材料的可合成性、可能的前驱体及合成路径^[79], 从而在材料设计阶段兼顾实验可行性与理论稳定性. 未来可将物理约束与数据驱动方法紧密结合, 发展专门面向材料可合成性的预测模型与生成框架, 从而推动生成模型向生成具有实际合成可行性的材料发展.

4 总结与展望

随着生成模型在材料科学中的不断深化应用, 目标性质导向的材料结构生成方法正逐步成为推动材料逆向设计范式转变的核心路径^[73,74]. 相比于传统的基于高通量筛选^[4-6]或先验组合^[80]的结构发现方法, 将目标性质嵌入生成流程中, 不仅能够提升模型的生成有效性和实用性, 更为按需设计高性能材料提供了新的技术支持.

尽管目前已有如条件生成、扩散过程调控、隐变量优化等多种方法被提出用于实现性质导向的晶体结构生成, 但在实际应用中仍存在诸多挑战: 目标性质的引导效果仍受限于训练数据的分布, 面对极端性质区间或罕见结构构型时模型生成能力明显下降. 更为关键的是, 如何保证生成模型所产生的结构在化学上合理并具有可合成性, 仍是当前亟待解决的难题. 除此之外, 当前模型多以单一性质为目标, 难以同时满足多目标甚至约束性质的联合优化需求, 以及模型生成结果的物理合理性与稳定性仍依赖后验筛选与外部验证, 尚未实现从结构生成到性质保障的端到端闭环.

未来的发展方向可围绕数据、模型与系统集成等多个方面深入推进. 在数据方面, 可借助小样本学习等策略来提升模型在稀疏区域的泛化能力. 在模型结构方面, 可通过融合结构归纳偏置与多模态

条件建模机制来进一步增强生成模型对多目标协同设计的适应性. 同时结合大规模材料数据库与实验反馈, 构建主动学习的闭环系统, 推动生成模型由关注结构生成本身, 转向集成结构稳定性评估与性质预测功能, 实现从生成到筛选再到性质判断的全流程闭环. 随着理论方法的持续完善、数据资源的不断扩充以及计算平台的升级迭代, 基于生成模型的性质导向材料设计有望在能源、催化、电子等关键材料领域加速实现从理论预测到实验验证的转化, 成为推动新材料发现的重要驱动力.

参考文献

- [1] Jiang X, Xue D Z, Bai Y, Wang W Y, Liu J J, Yang M L, Su Y J 2025 *Rev. Mater. Res.* **1** 100010
- [2] Wu M F, Zhang S Y, Ren J 2025 *APL Mater.* **13** 020601
- [3] Merchant A, Batzner S, Schoenholz S S, Aykol M, Cheon G, Cubuk E D 2023 *Nature* **624** 80
- [4] Butler K T, Davies D W, Cartwright H, Isayev O, Walsh A 2018 *Nature* **559** 547
- [5] Zhong M, Tran K, Min Y M, Wang C H, Wang Z Y, Dinh C T, De Luna P, Yu Z Q, Rasouli A S, Brodersen P, Sun S, Voznyy O, Tan C S, Askerka M, Che F L, Liu M, Seifitokaldani A, Pang Y J, Lo S C, Ip A, Ulissi Z, Sargent E H 2020 *Nature* **581** 178
- [6] Rao Z Y, Tung P Y, Xie R W, Wei Y, Zhang H B, Ferrari A, Klaver T P C, Körmann F, Sukumar P T, Kwiatkowski da Silva A, Chen Y, Li Z M, Ponge D, Neugebauer J, Gutfleisch O, Bauer S, Raabe D 2022 *Science* **378** 78
- [7] Hellenbrandt M 2004 *Crystallogr. Rev.* **10** 17
- [8] Kirklin S, Saal J E, Meredig B, Thompson A, Doak J W, Aykol M, Rühl S, Wolverton C 2015 *npj Comput. Mater.* **1** 1
- [9] Jain A, Ong S P, Hautier G, Chen W, Richards W D, Dacek S, Cholia S, Gunter D, Skinner D, Ceder G, Persson K A 2013 *APL Mater.* **1** 011002
- [10] Hastrup S, Strange M, Pandey M, Deilmann T, Schmidt P S, Hinsche N F, Gjerding M N, Torelli D, Larsen P M, Riis-Jensen A C, Gath J, Jacobsen K W, Jørgen Mortensen J, Olsen T, Thygesen K S 2018 *2D Mater.* **5** 042002
- [11] Lu S H, Zhou Q H, Chen X Y, Song Z L, Wang J L 2022 *Natl. Sci. Rev.* **9** nwac111
- [12] Kingma D P, Welling M 2019 *Found. Trends Mach. Learn.* **12** 307
- [13] Goodfellow I J, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, Courville A, Bengio Y 2014 *arXiv: 1406.2661v1 [stat. ML]*
- [14] Ho J, Jain A, Abbeel P 2020 *arXiv: 2006.11239v2 [cs. LG]*
- [15] Stimper V, Liu D, Campbell A, Berenz V, Ryll L, Schölkopf B, Hernández-Lobato J M 2023 *J. Open Source Softw.* **8** 5361
- [16] Hoogeboom E, Gritsenko A A, Bastings J, Poole B, Berg R van den, Salimans T 2022 *arXiv: 2110.02037v2 [cs. LG]*
- [17] Hong T, Chen T K, Jin D L, Zhu Y, Gao H, Zhao K, Zhang T Y, Ren W, Cao G X 2025 *npj Quantum Mater.* **10** 12
- [18] Jin L Z, Du Z J, Shu L, Cen Y, Xu Y F, Mei Y F, Zhang H 2025 *Nat. Commun.* **16** 1210
- [19] Sanchez-Lengeling B, Aspuru-Guzik A 2018 *Science* **361** 360
- [20] Zhang K, Chen T, Abbas Y, Jan S U, Zhou Z H, Chu S Q,

- Xie G C, Ullah S, Akram M Z, Zhang J, Xuan Y M, Gong J R 2021 *Matter* **4** 1054
- [21] Manzoor A, Zhang Y, Aidhy D S 2021 *Comput. Mater. Sci* **198** 110669
- [22] Ye C Y, Weng H M, Wu Q S 2024 *Comput. Mater. Today* **1** 100003
- [23] Okabe R, Cheng M, Chotrattanapituk A, Hung N T, Fu X, Han B, Wang Y, Xie W, Cava R J, Jaakkola T S, Cheng Y, Li M 2024 arXiv: 2407.04557v1 [cond-mat. mtrl-sci]
- [24] Ye C, Wang Y, Xie X, Zhu T, Liu J, He Y, Zhang L, Zhang J, Fang Z, Wang L, Liu Z, Weng H, Wu Q 2025 arXiv: 2505.00076v1 [cond-mat. mtrl-sci]
- [25] Zeni C, Pinsler R, Züchner D, Fowler A, Horton M, Fu X, Wang Z, Shysheya A, Crabbé J, Ueda S, Sordillo R, Sun L, Smith J, Nguyen B, Schulz H, Lewis S, Huang C W, Lu Z, Zhou Y, Yang H, Hao H, Li J, Yang C, Li W, Tomioka R, Xie T 2025 *Nature* **639** 624
- [26] Song Z L, Fan L F, Lu S H, Ling C Y, Zhou Q H, Wang J L 2025 *Nat. Commun.* **16** 1053
- [27] Chen X Y, Lu S H, Chen Q, Zhou Q H, Wang J L 2024 *Nat. Commun.* **15** 5391
- [28] Choudhary K 2024 arXiv: 2405.03680v2 [cond-mat. mtrl-sci]
- [29] Lu S H, Zhou Q H, Guo Y L, Wang J L 2022 *Chem* **8** 769
- [30] Xie T, Fu X, Ganea O E, Barzilay R, Jaakkola T 2021 arXiv: 2110.06197v3 [cs. LG]
- [31] Dan Y B, Zhao Y, Li X, Li S B, Hu M, Hu J J 2020 *npj Comput. Mater.* **6** 84
- [32] Chen S, Ge C, Zhang S, Sun P, Luo P 2025 arXiv: 2504.07963v1 [cs. CV]
- [33] Shi C, Xu M, Zhu Z, Zhang W, Zhang M, Tang J 2020 arXiv: 2001.09382v2 [cs. LG]
- [34] Xu M, Yu L, Song Y, Shi C, Ermon S, Tang J 2022 arXiv: 2203.02923v1 [cs. LG]
- [35] Kingma D P, Welling M 2022 arXiv: 1312.6114v11 [stat. ML]
- [36] Song Y, Ermon S 2019 *Proceedings of the 32nd Annual Conference on Neural Information Processing Systems Vancouver, Canada, December 8–14, 2019* p11895
- [37] Lemons D S, Gythiel A 1997 *Am. J. Phys.* **65** 1079
- [38] Court C J, Yildirim B, Jain A, Cole J M 2020 *J. Chem. Inf. Model.* **60** 4518
- [39] Davies D W, Butler K T, Jackson A J, Skelton J M, Morita K, Walsh A 2019 *J. Open Source Softw.* **4** 1361
- [40] Xu M, Luo S, Bengio Y, Peng J, Tang J 2021 arXiv: 2102.10240v3 [cs. LG]
- [41] Ganea O E, Pattanaik L, Coley C W, Barzilay R, Jensen K F, Green W H, Jaakkola T S 2021 arXiv: 2106.07802v1 [physics. chem-ph]
- [42] Castelli I E, Landis D D, Thygesen K S, Dahl S, Chorkendorff I, Jaramillo T F, Jacobsen K W 2012 *Energy Environ. Sci.* **5** 9034
- [43] Pickard C J <https://archive.materialscloud.org/record/2020.0026/v1> [2024-04-20]
- [44] Ren Z K, Tian S I P, Noh J, Oviedo F, Xing G Z, Li J L, Liang Q H, Zhu R M, Aberle A G, Sun S J, Wang X N, Liu Y, Li Q X, Jayavelu S, Hippalgaonkar K, Jung Y, Buonassisi T 2022 *Matter* **5** 314
- [45] Gebauer N W A, Gastegger M, Schütt K T 2020 arXiv: 1906.00957v3 [stat. ML]
- [46] Zhang H, Georgescu A B, Yerramilli S, Karpovich C, Apley D W, Olivetti E A, Rondinelli J M, Chen W 2025 arXiv: 2412.17283v2 [cond-mat. mtrl-sci]
- [47] Murphy L R, Meek T L, Allred A L, Allen L C 2000 *J. Phys. Chem. A* **104** 5867
- [48] Sessa F, Rahm M 2022 *J. Phys. Chem. A* **126** 5472
- [49] Arjovsky M, Chintala S, Bottou L 2017 arXiv: 1701.07875v3 [stat. ML]
- [50] Maaten L van der, Hinton G 2008 *J. Mach. Learn. Res.* **9** 2579
- [51] Green M L, Choi C L, Hattrick-Simpers J R, et al. 2017 *Appl. Phys. Rev.* **4** 011105
- [52] Pickard C J, Needs R J 2011 *J. Phys. Condens. Matter* **23** 053201
- [53] Song Y, Sohl-Dickstein J, Kingma D P, Kumar A, Ermon S, Poole B 2021 arXiv: 2011.13456v2 [cs. LG]
- [54] Schmidt J, Wang H C, Cerqueira T F T, Botti S, Marques M A L 2022 *Sci. Data* **9** 64
- [55] Jiao R, Huang W, Lin P, Han J, Chen P, Lu Y, Liu Y 2023 arXiv: 2309.04475v2 [cond-mat. mtrl-sci]
- [56] Sultanov A, Crivello J C, Rebafka T, Sokolovska N 2023 *J. Chem. Inf. Model.* **63** 6986
- [57] Xie T, Grossman J C 2018 *Phys. Rev. Lett.* **120** 145301
- [58] Ward L, Agrawal A, Choudhary A, Wolverton C 2016 *npj Comput. Mater.* **2** 1
- [59] Han X Q, Wang X D, Xu M Y, Feng Z, Yao B W, Guo P J, Gao Z F, Lu Z Y 2025 *Chin. Phys. Lett.* **42** 027403
- [60] Long T, Zhang Y, Zhang H 2024 arXiv: 2409.19124v1 [cond-mat. mtrl-sci]
- [61] Chen L, Zhang W, Nie Z, Li S, Pan F 2021 *J. Mater. Inform.* **1** 4
- [62] Chen Y, Wang X, Deng X, Liu Y, Chen X, Zhang Y, Wang L, Xiao H 2024 arXiv: 2408.07608v1 [cond-mat. mtrl-sci]
- [63] New A, Pekala M, Pogue E A, Le N Q, Domenico J, Piatko C D, Stiles C D 2023 arXiv: 2309.12323v1 [cond-mat. mtrl-sci]
- [64] Banik S, Dhabal D, Chan H, Manna S, Cherukara M, Molinero V, Sankaranarayanan S K R S 2023 *npj Comput. Mater.* **9** 1
- [65] Cao Z, Luo X, Lv J, Wang L 2024 arXiv: 2403.15734v2 [cond-mat. mtrl-sci]
- [66] Yang S, Cho K, Merchant A, Abbeel P, Schuurmans D, Mordatch I, Cubuk E D 2024 arXiv: 2311.09235v2 [cs. LG]
- [67] Jiao R, Huang W, Liu Y, Zhao D, Liu Y 2024 arXiv: 2402.03992v2 [cs. LG]
- [68] Tsai W F, Fang C, Yao H, Hu J 2015 *New J. Phys.* **17** 055016
- [69] Ho J, Salimans T 2022 arXiv: 2207.12598v1 [cs. LG]
- [70] Zhang L, Rao A, Agrawala M 2023 arXiv: 2302.05543v3 [cs. CV]
- [71] Shuaibi M, Kolluru A, Das A, Grover A, Sriram A, Ulissi Z, Zitnick C L 2021 arXiv: 2106.09575v1 [cs. LG]
- [72] Guo Z, Zhang C, Yu W, Herr J, Wiest O, Jiang M, Chawla N V 2021 arXiv: 2102.07916v1 [cs. LG]
- [73] Ryan K, Lengyel J, Shatruk M 2018 *J. Am. Chem. Soc.* **140** 10158
- [74] Chenebua E T, Nganbe M, Tchagang A B 2024 *npj Comput. Mater.* **10** 198
- [75] Fung V, Zhang J, Juarez E, Sumpter B G 2021 *npj Comput. Mater.* **7** 1
- [76] Varol Altay E, Alatas B 2020 *Artif. Intell. Rev.* **53** 1373
- [77] Butler K T, Choudhary K, Csanyi G, Ganose A M, Kalinin S V, Morgan D 2024 *npj Comput. Mater.* **10** 1
- [78] Wu Y L, Li X Y, Guo R, Xu R Q, Ju M G, Wang J L 2025 *Natl. Sci. Rev.* **12** nwaf081
- [79] Song Z L, Lu S H, Ju M G, Zhou Q H, Wang J L 2025 *Nat. Commun.* **16** 6530
- [80] Seko A, Hayashi H, Tanaka I 2018 *J. Chem. Phys.* **148** 241719

SPECIAL TOPIC—AI + Physical Science

Goal-property-guided material generation: Toward on-demand construction via inverse design of materials*

LIU Zhanghe¹⁾ CHEN Xinyu¹⁾ ZHOU Qionghua^{1)2)†} WANG Jinlan^{1)2)‡}

1) (*Key Laboratory of Quantum Materials and Devices of Ministry of Education,
School of Physics, Southeast University, Nanjing 211189, China*)

2) (*Suzhou Laboratory, Suzhou 215004, China*)

(Received 24 July 2025; revised manuscript received 10 September 2025)

Abstract

In recent years, the application of machine learning in materials science has significantly accelerated the discovery of new materials. In particular, when combined with traditional methods such as first-principles calculations, machine learning models have proven effective in screening potential high-performance materials from existing databases. However, these methods are largely limited by the known chemical spaces, making it difficult to achieve the active design of novel material structures. To overcome this limitation, generative models have become a promising tool for inverse material design, providing new avenues for exploring unknown structures and property spaces. Although existing generative models have achieved initial progress in crystal structure generation, achieving property-guided material generation remains a significant challenge. In this review paper, we first introduce the representative generative models recently applied to materials generation, including CDVAE, MatGAN, and MatterGen, and analyzes their basic abilities and limitations in structural generation. We then focus on strategies for incorporating target properties into generative models to generate the property-guided structure. Specifically, we discuss four representative methods: Con-CDVAE based on target property vectors, SCIGEN with integrated structural constraints and guidance mechanisms, a fine-tuned version of MatterGen leveraging adapter-based property control, and a CDVAE latent space optimization strategy guided by property objectives. Finally, we summarize the key challenges faced by property-guided generative models and provide an outlook on future research directions. This review aims to offer researchers a systematic reference and inspiration for advancing property-driven generative approaches in material design and provides researchers with a systematic reference and insight into the advancement of property-driven generative methods for materials design.

Keywords: machine learning, generative models, inverse design, property-guided

PACS: 07.05.Mh, 91.60.Ed, 81.05.Zx

DOI: [10.7498/aps.74.20250989](https://doi.org/10.7498/aps.74.20250989)

CSTR: [32037.14.aps.74.20250989](https://cstr.cn/32037.14.aps.74.20250989)

* Project supported by the National Key Research and Development Program of China (Grant No. 2021YFA1500703), the National Natural Science Foundation of China (Grant Nos. 22033002, T2321002, 22373013), and the Frontier Leading Technology Basic Research Major Project of Jiangsu Provincial Science and Technology Planning Special Fund, China (Grant Nos. BK20222007, BK20232012).

† Corresponding author. E-mail: qh.zhou@seu.edu.cn

‡ Corresponding author. E-mail: jlwang@seu.edu.cn

目标性质导向的材料生成：迈向按需构筑的材料逆向设计

刘章赫 陈新宇 周颢桦 王金兰

Goal–property–guided material generation: Toward on–demand construction via inverse design of materials

LIU Zhanghe CHEN Xinyu ZHOU Qionghua WANG Jinlan

引用信息 Citation: *Acta Physica Sinica*, 74, 240701 (2025) DOI: 10.7498/aps.74.20250989

CSTR: 32037.14.aps.74.20250989

在线阅读 View online: <https://doi.org/10.7498/aps.74.20250989>

当期内容 View table of contents: <http://wulixb.iphy.ac.cn>

您可能感兴趣的其他文章

Articles you may be interested in

生物大分子过渡态搜索算法及其中的机器学习

Transition state searching for complex biomolecules: Algorithms and machine learning

物理学报. 2023, 72(24): 248701 <https://doi.org/10.7498/aps.72.20231319>

机器学习辅助绝热量子算法设计

Machine learning assisted quantum adiabatic algorithm design

物理学报. 2021, 70(14): 140306 <https://doi.org/10.7498/aps.70.20210831>

实现散射场强整形的微散射体阵列逆向设计方法

Inverse design method of microscatterer array for realizing scattering field intensity shaping

物理学报. 2021, 70(1): 010202 <https://doi.org/10.7498/aps.70.20200825>

机器学习结合固溶强化模型预测高熵合金硬度

Machine learning combined with solid solution strengthening model for predicting hardness of high entropy alloys

物理学报. 2023, 72(18): 180701 <https://doi.org/10.7498/aps.72.20230646>

蛋白质计算中的机器学习

Machine learning for *in silico* protein research

物理学报. 2024, 73(6): 069301 <https://doi.org/10.7498/aps.73.20231618>

机器学习模型预测稀土化合物的热力学稳定性

Machine learning model predicted thermodynamic stability of rare earth compounds

物理学报. 2025, 74(13): 130201 <https://doi.org/10.7498/aps.74.20250362>