

专题: 原子分子和材料物性数据

## 卷积神经网络辅助无机晶体弹性性质预测\*

刘宇杰<sup>1) #</sup> 王振宇<sup>1) #</sup> 雷航<sup>1) #</sup> 张国宇<sup>1)</sup> 戚家伟<sup>2) †</sup>高志斌<sup>1) ‡</sup> 孙军<sup>3)</sup> 宋海峰<sup>2)</sup> 丁向东<sup>3)</sup>

1) (西安交通大学材料科学与工程学院, 金属多孔材料全国重点实验室, 西安 710049)

2) (北京应用物理与计算数学研究所, 计算物理全国重点实验室, 北京 100088)

3) (西安交通大学材料科学与工程学院, 金属材料强度全国重点实验室, 西安 710049)

(2025年1月25日收到; 2025年3月18日收到修改稿)

无机晶体材料因具有优异的物理和化学特性, 在多个领域展现出广泛的应用潜力. 弹性性质 (如体积模量和剪切模量) 对预测材料的电导率、热导率及力学性能具有重要作用, 然而, 传统实验测量方法存在成本高、周期长等问题. 随着计算方法的进步, 理论模拟逐渐成为独立于实验的研究方法. 近年来, 基于图神经网络的机器学习方法在无机晶体材料的弹性性质预测中取得了显著成果, 尤其是晶体图卷积神经网络 (CGCNN) 在材料数据的预测和扩展方面表现出色. 本研究利用从 Matbench v0.1 数据集中收集的 10987 个材料的体积模量和剪切模量数据, 训练了两个 CGCNN 模型, 基于预训练的模型成功实现了对 80664 个无机晶体结构弹性模量的预测. 为保证数据质量, 筛选了材料电子带隙在 0.1—3.0 eV 之间, 并去除了含有放射性元素的化合物. 预测数据来源于两个主要数据集: 一是从 Materials Project 数据库中筛选出的 54359 个晶体结构, 构成 MPED 弹性数据集; 二是 Merchant 等 (2023 *Nature* 624 80) 通过深度学习和图神经网络方法发现的 26305 种晶体结构, 构成 NED 弹性数据集. 最终, 本研究预测了 80664 种无机晶体的体积模量和剪切模量, 丰富了现有的材料弹性数据资源, 并为材料设计提供了更多的数据支持. 本文数据集可在 <https://doi.org/10.57760/sciencedb.j00213.00104> 中访问获取.

关键词: 无机晶体材料, 弹性模量, 机器学习, 力学性质预测

PACS: 07.05.Tp, 61.72.-y, 62.20.D-

DOI: 10.7498/aps.74.20250127

CSTR: 32037.14.aps.74.20250127

## 1 引言

无机晶体材料因具有独特的物理和化学特性, 在电子<sup>[1,2]</sup>、光学<sup>[3,4]</sup>、热学<sup>[5]</sup>和力学<sup>[6,7]</sup>等多个领域展现出巨大的应用潜力, 其弹性性质 (如体积模量 ( $B$ ) 和剪切模量 ( $G$ ) 等) 对预测电导率<sup>[8]</sup>、热导率<sup>[9,10]</sup>和力学性能<sup>[11]</sup>等物理性质有关键作用. 例如, 自

20 世纪 50 年代 Pugh<sup>[12]</sup> 的研究以来, 体积模量与剪切模量之比已成为理解和预测材料延展性的重要指标<sup>[13-17]</sup>. 因此, 无机晶体材料弹性数据对预测材料性能至关重要, 且为功能材料 (如热电、压电、铁电材料) 的性能优化提供了重要的基础数据.

然而, 传统实验测量方法在获取无机晶体材料弹性模量方面存在成本高、周期长等局限. 随着算力提升和计算方法的改进, 通过理论模拟的方法可

\* 国家自然科学基金 (批准号: 12104356, 52250191)、计算物理全国重点实验室基金和国家重点研发计划 (批准号: 2023YFB4604100) 资助的课题.

# 同等贡献作者.

† 通信作者. E-mail: [xian\\_jiawei@iapcm.ac.cn](mailto:xian_jiawei@iapcm.ac.cn)‡ 通信作者. E-mail: [zhibin.gao@xjtu.edu.cn](mailto:zhibin.gao@xjtu.edu.cn)

以缓解实验测量方式的成本问题. 比如, 基于密度泛函理论 (DFT)<sup>[18]</sup> 的应力-应变关系分析法<sup>[19]</sup>, 这种方法通过求解材料的电子结构来精确预测其力学性能. 另一种常见的模拟方法是分子动力学 (MD) 模拟<sup>[20]</sup>, 即通过经典模型/机器学习势函数对材料的能量、原子位移和原子受力进行模拟, 从而估算出其弹性模量等力学性质. 除此之外, 介观尺度的蒙特卡罗模拟 (MC)<sup>[21]</sup> 和宏观尺度的有限元分析 (FEA)<sup>[22]</sup> 也可以作为有效的模拟工具, 前者适用于研究系统的统计性质, 后者则能模拟大规模的复杂结构与部件. 这些模拟方法获得的数据与实验测量结果比较接近, 但在处理大规模数据时, 通常会面临较高的计算成本和周期. 因此, 如何在保证预测精度的同时提升计算效率, 已成为当前亟待解决的核心问题. 人工智能 (AI) 和机器学习技术的快速发展, 使得获取更多无机晶体材料的弹性性质数据成为可能<sup>[23,24]</sup>, 并为传统模拟方法提供了更高效的替代方案.

数据驱动的方式已经成为一种扩大材料空间的有效方式<sup>[25]</sup>, 但获取高质量的大量数据仍然是一个挑战. 近年来的一些研究发现, 许多机器学习模型在预测晶体结构弹性性质中表现出较好的效果, 其中, Xie 和 Grossman<sup>[26]</sup> 提出的晶体图卷积神经网络 (crystal graph convolutional neural networks, CGCNN) 备受瞩目. 这种模型能够有效地将晶体信息转换为图信息, 进一步处理图结构数据, 捕捉节点间复杂映射关系, 提升特征学习能力, 预测材料的性质. 受 CGCNN 启发, 多个图神经网络模型应运而生, 如轨道图卷积神经网络 (orbital graph convolutional neural network, OGCNN)<sup>[27]</sup>、原子线图神经网络 (atomic line graph neural network, ALIGNN)<sup>[28]</sup>、图注意力图神经网络 (graph attention network graph neural network, GATGNN)<sup>[29]</sup>、连接优化晶体图网络 (connection optimized crystal graph network, coGN) 及其扩展版本 (connection optimized next-generation crystal graph network, coNGN)<sup>[30]</sup>. 如何在众多图卷积模型框架中选择在训练集之外依然保持较好泛化性的模型呢? 最近的一项研究为我们提供了启示, Omee 等<sup>[31]</sup> 评估了 8 种图神经网络 (graph neural network, GNN) 模型在 5 个分布外 (out-of-distribution, OOD) 测试集上的性能, 特别针对弹性数据集的结果表明, CGCNN 模型在 LOCO (leave-one-cluster-out) 和

SparseXsingle (single-point targets with the lowest structure density) 测试中取得了最小的平均绝对误差 (MAE), 在不同数据集上的表现稳定且准确, 这说明 CGCNN 模型在数据集之外具有很好的泛化能力和异常值发现与探索能力.

虽然目前在弹性模量数据库方面取得了一些进展, 2015 年开始, de Jong<sup>[32]</sup> 等结合实验数据与第一性原理计算, 设计了一种高通量的第一性原理计算方法, 系统性地研究了数千种无机晶体材料的弹性常数, 并构建了详细的弹性属性数据库. 此外, 国外主流数据库如 Materials Project<sup>[33]</sup> 包含了超过 10000 种材料的弹性常数数据; AFLOW<sup>[34]</sup> 提供了约 6000 种无机材料的弹性数据; OQMD (open quantum materials database)<sup>[35]</sup> 则涵盖了大约 4000 种材料的相关信息. 然而, 国内数据库相对较少, 虽然 Atomly<sup>[36]</sup> 数据库包含了大量的材料数据, 但无机晶体的弹性数据仍然占其中较少一部分. 因此, 建立一个富含无机晶体材料弹性性质的数据集显得尤为必要.

本文收集 10987 个体积模量和剪切模量的晶体结构以及对应的性质, 训练了两个弹性模量 CGCNN 模型, 基于预训练的 CGCNN 模型预测了收集的无机晶体结构的弹性模量, 扩大了整个材料空间的弹性模量数据集. 其流程大致如下: 利用 Matbench v0.1 中从 Materials Project 收集的包含 10987 个材料条目的体积模量和剪切模量数据集<sup>[37]</sup>, 训练了两个从晶体学信息文件 (CIF) 映射到体积模量和剪切模量的 CGCNN 模型. 由于带隙过大会导致较差的电导率, 且放射性元素对人体有害, 我们对收集的未含有模量信息的数据做了进一步筛选, 筛选标准为带隙在 0.1—3.0 eV 之间, 同时剔除包含放射性元素的单质及化合物. 预测的晶体结构数据主要来源于以下两部分: 1) 从 Materials Project<sup>[33]</sup> 数据库获取的晶体结构, 共筛选出 54359 个材料, 这些晶体结构组成的数据集记为 Materials Project elastic dataset (MPED) 数据集; 2) Merchant 等<sup>[38]</sup> 通过深度学习和图神经网络 (GNN) 发现的晶体结构, 共筛选出 26305 种结构, 这些晶体结构组成的数据集记为 nature elastic dataset (NED) 数据集. 最终, 预测了 80664 种无机晶体结构的体积模量和剪切模量, 一定程度上丰富了现有的弹性数据资源, 还能为功能材料性能设计与优化提供更多的数据支持.

## 2 方法

### 2.1 数据获取

Matbench\_v1.0 测试集<sup>[37]</sup>中包含 13 个不同材料属性, 通过 python 的 matminer 包可以实现数据的下载. 本文使用了该数据集中的两个收集自 Materials Project 数据库的数据集, 它们专门处理弹性属性: “matbench\_log\_gvrh”(用于从结构预测 DFT log<sub>10</sub> vrh 平均剪切模量) 和 “matbench\_log\_kvrrh”(用于从结构预测 DFT log<sub>10</sub> vrh 平均体积模量), 这两个数据集包含相同的材料条目, 都是 10987 个, 旨在通过 Voigt-Reuss-Hill (VRH) 平均方法预测剪切模量 ( $G$ ) 和体积模量 ( $B$ ). 由于这些数据集标准化和全面性的特点, 它们成为训练机器学习模型、预测关键弹性特性的理想选择.

获取高通量数据集可能具有一定挑战性. 然而, 最近机器学习的进展极大推动了稳定材料的发现. Merchant 等<sup>[38]</sup>利用深度学习和图神经网络 (GNN) 来拓展材料发现的范围, 尤其是无机晶体的研究. 他们的工作通过在凸包中添加 381000 个新条目, 扩大了已知材料的范围, 比之前的数据集增加了十倍. 我们通过 GitHub 访问了他们的数据集: [https://github.com/google-deepmind/materials\\_discovery](https://github.com/google-deepmind/materials_discovery). 该库的 “by\_composition” 文件夹中包含了 377221 个有效的 CIF 文件, 这些文件与 CGCNN 兼容. 此外, 还提供了一个汇总的 CSV 文件, 包含带隙、晶体对称性和分解能等数据. 这些材料包括 2—6 种元素, 原子序数范围为 2—106. 我们在其中筛选出带隙在 0.1—3.0 eV 之间的稳定结构, 有 30199 种, 排除掉对人体有害的放射性元素后, 这个数字降低至 26305. 另一方面, 来自 Materials Project 的其他数据进一步补充了这些数据集, 并支持通过开源 API 有针对性地检索材料属性, 具体来说, 通过筛选带隙在 0.1—3.0 eV 之间且无放射性元素的结构, 得到了 54359 种不同的结构. 总体而言, 共获得了 80664 种稳定的结构, 这些资源为高通量计算和分析提供了坚实的基础.

### 2.2 晶体图卷积神经网络 (crystal graph convolutional neural networks)

如图 1 所示, CGCNN 通过将晶体结构映射为图表示, 其中节点代表原子 (使用 92 维 one-hot 向量编码原子属性), 边表示原子间的化学键 (通过高

斯展开处理原子间距). 对于多体结构信息, 模型直接编码了键长信息, 而键角和二面角虽未显式表示, 但通过多层图卷积的消息传递机制进行隐式学习——每一层卷积层利用非线性函数将当前节点的特征、相邻节点的特征及连接边的特征进行融合, 更新节点表示, 从而逐步捕获更复杂的局部结构信息; 长程相互作用则通过设定合适的截断半径和多层卷积的迭代传递来获取. 对于局部畸变导致的对称性破缺 (如 Jahn-Teller 效应), CGCNN 通过精确记录每个原子的局部配位环境, 结合图卷积层的非线性变换能力, 可以有效捕捉这些结构畸变. 最终, 模型通过全局池化操作将所有原子特征整合为晶体表示, 既保持了局部结构信息的完整性, 又能有效表达全局的结构特征, 这种层级化的特征提取机制使得 CGCNN 能够在保持模型简洁性的同时, 准确描述晶体材料的本质特征.

本文设计了 3 个卷积层. 每个卷积层首先收集邻近原子、中心原子和键的信息, 并将这些特征拼接在一起. 随后, 特征通过一个全连接层, 并应用 Sigmoid 门控机制进行调节, 最后使用 softplus 激活函数进行非线性变换. 接下来, 通过一个池化层将原子级信息聚合到晶体级别, 并接入一个转换层, 将卷积特征转换为全连接层特征. 最终, 模型接入两个全连接隐藏层以进一步提取特征. 由于弹性性能与晶体结构密切相关, CGCNN 模型能够有效捕捉晶体结构的关键特征, 因此可直接用于从晶体结构预测弹性性能. 值得一提的是, 还有一些工作通过更精确的方式预测模量<sup>[39]</sup>. 为了提高模型的收敛性, 本文采用自适应矩估计 (adaptive moment estimation, Adam) 优化方法, 初始学习率设置为 0.001. Adam 结合了动量和均方根反向传播 (root mean square propagation, RMSProp) 的优点, 通过计算梯度的一阶矩 (均值) 和二阶矩 (未中心化的方差) 来动态调整学习率, 对初始学习率不敏感. 最终, 模型在测试集上的平均绝对误差 (MAE) 为 0.0981 log<sub>10</sub>(GPa) (剪切模量) 和 0.0790 log<sub>10</sub>(GPa) (体积模量), 验证了其预测精度. 为了进一步优化模型, 在 Adam 优化器的基础上, 手动调整了训练迭代步数、卷积层数量和隐藏层数量, 并选择验证集上表现最佳的模型. 通过计算训练集、验证集和测试集的 MAE 和  $R^2$  分数, 确定了最终模型. 最后, 关于该部分的更多信息可以参考补充材料 (online).

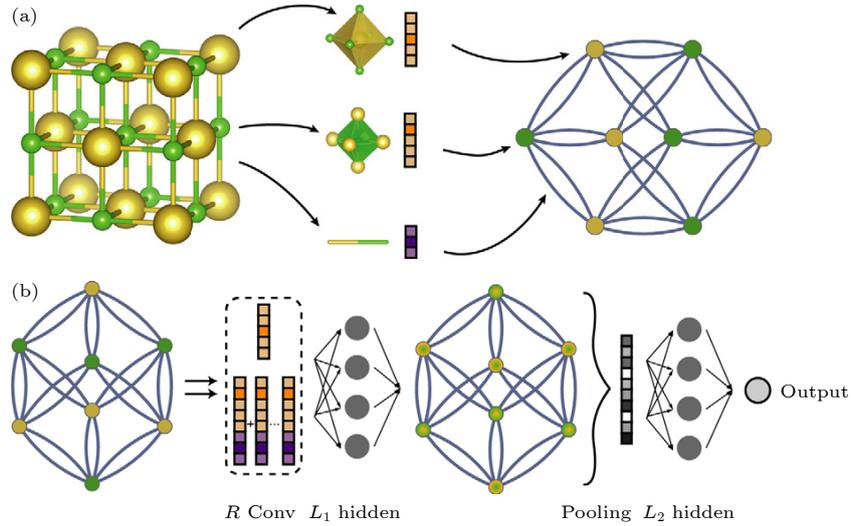


图1 晶体图卷积神经网络 (a) 晶体图的构造: 晶体结构被转换为图形, 其中节点代表单位原胞中的原子, 边缘则表示原子之间的连接. 每个节点和边都用对应于晶体中原子和化学键的向量进行表征; (b) 晶体图顶部的卷积神经网络结构: 在每个节点上构建  $R$  个卷积层和  $L_1$  个隐藏层, 从而形成一个新的图, 其中每个节点表示该原子的局部环境. 经过池化操作后, 将代表整个晶体的向量连接到  $L_2$  隐藏层, 然后连接到输出层, 以进行预测<sup>[26]</sup>

Fig. 1. Illustration of the crystal graph convolutional neural networks: (a) Construction of the crystal graph. Crystals are converted to graphs with nodes representing atoms in the unit cell and edges representing atom connections. Nodes and edges are characterized by vectors corresponding to the atoms and bonds in the crystal, respectively. (b) Structure of the convolutional neural network on top of the crystal graph.  $R$  convolutional layers and  $L_1$  hidden layers are built on top of each node, resulting in a new graph with each node representing the local environment of each atom. After pooling, a vector representing the entire crystal is connected to  $L_2$  hidden layers, followed by the output layer to provide the prediction<sup>[26]</sup>.

### 2.3 弹性性能

Jia 等<sup>[40]</sup> 提出了一种方法, 通过体积模量 ( $B$ ) 和剪切模量 ( $G$ ) 精确估算材料的其他弹性性质 (如泊松比、声速等), 这种方法比实验测量效率更高、周期更短. 因此, 当获得材料的体积模量 ( $B$ ) 和剪切模量 ( $G$ ) 时, 就可以估算出其他弹性性质. 而 CGCNN 模型能够较为准确地预测材料的体积模量 ( $B$ ) 和剪切模量 ( $G$ ), 所以, 基于 MPED 数据集和 NED 数据集的基本物理量 (如密度等), 结合上述方法, 可以估算出弹性特性和声速等一系列物理量. 具体计算方法如下<sup>[40,41]</sup>:

$$v_l = \sqrt{[B + (4/3)G]/\rho}, \quad (1)$$

$$v_t = \sqrt{G/\rho}, \quad (2)$$

$$v_s = \left[ \frac{1}{3} \left( \frac{1}{v_l^3} + \frac{2}{v_t^3} \right) \right]^{-1/3}, \quad (3)$$

其中  $v_l$ ,  $v_t$  和  $v_s$  分别为纵向声速、横向声速和平均声速,  $\rho$  为材料密度.

另外, 已经证明泊松比 ( $\nu$ ) 可以由下式求得<sup>[42,43]</sup>:

$$\nu = \frac{x^2 - 2}{2x^2 - 2}, \quad (4)$$

其中  $x$  为纵向声速与横向声速之比, 即  $x = v_l/v_t$ .

先前的研究表明德拜温度  $\theta_D$  与平均声速  $v_s$  成正比, 所以德拜温度可以根据弹性模量计算出来<sup>[44]</sup>:

$$\theta_D = \frac{\hbar v_s}{k_B} \left( \frac{3N}{4\pi V} \right)^{1/3}, \quad (5)$$

其中,  $\hbar$  是约化普朗克常数,  $v_s$  为平均声速,  $k_B$  是玻尔兹曼常数,  $N$  表示原胞内的原子数,  $V$  是原胞体积.

### 2.4 机器学习性能评估指标

在机器学习中, 平均绝对误差 MAE 和决定系数  $R^2$  是评估回归模型表现的常用指标. MAE 衡量预测误差的平均幅度, 定义如下:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|, \quad (6)$$

其中,  $y_i$  为实际值,  $\hat{y}_i$  为预测值. MAE 值越小, 表示模型的预测精度越高. 此外, 判定系数  $R^2$  (coefficient of determination) 的具体计算公式为

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}, \quad (7)$$

其中,  $\bar{y}$  表示实际值的平均值.  $R^2$  值越接近 1, 表明模型拟合效果越好.

### 3 结果

#### 3.1 CGCNN 模型评估

在之前的研究中, Wang 等<sup>[45]</sup> 通过将机器学习设计策略用于高强度铝锂合金的开发中发现径向基函数神经网络 (radial basis function neural networks, RBF) 的预测能力优于反向传播神经网络 (back propagation neural networks, BP). 为了说明 CGCNN 模型的优势, 本文对比了 CGCNN 和其他机器学习模型的性能, 例如随机森林 (random forest)、极端梯度提升 (extreme gradient boosting, XGBoost)、支持向量回归 (support vector regression, SVR)、梯度提升 (gradient noosting) 和决策

树 (decision tree). 为了和晶体结构保持相关性, 选取原胞的平均原子序数、平均原子质量、平均电负性、空间群号、密度和每原子体积为特征构建模型. 图 2 展示了六种模型在预测剪切模量 ( $G$ ) 和体模量 ( $B$ ) 时的性能对比. 其中, 图 2(a) 与图 2(b) 分别展示了剪切模量模型的平均绝对误差 (MAE) 和决定系数 ( $R^2$ ), 图 2(c) 和图 2(d) 分别展示了预测体模量模型的平均绝对误差 (MAE) 和决定系数 ( $R^2$ ). 由于模型在验证集和测试集的性能可以评价其在训练集之外的泛化能力, 因为本文用验证集和测试集的 MAE 和  $R^2$  均值来评价模型性能的好坏, 图中 MAE 按照均值从小到大,  $R^2$  按照均值从大到小排列. 结果表明, CGCNN 模型在验证集和测试集上均表现出较低的 MAE 和较高的  $R^2$ , 显示出其在预测剪切模量和体模量方面具有更高的精度和可靠性.

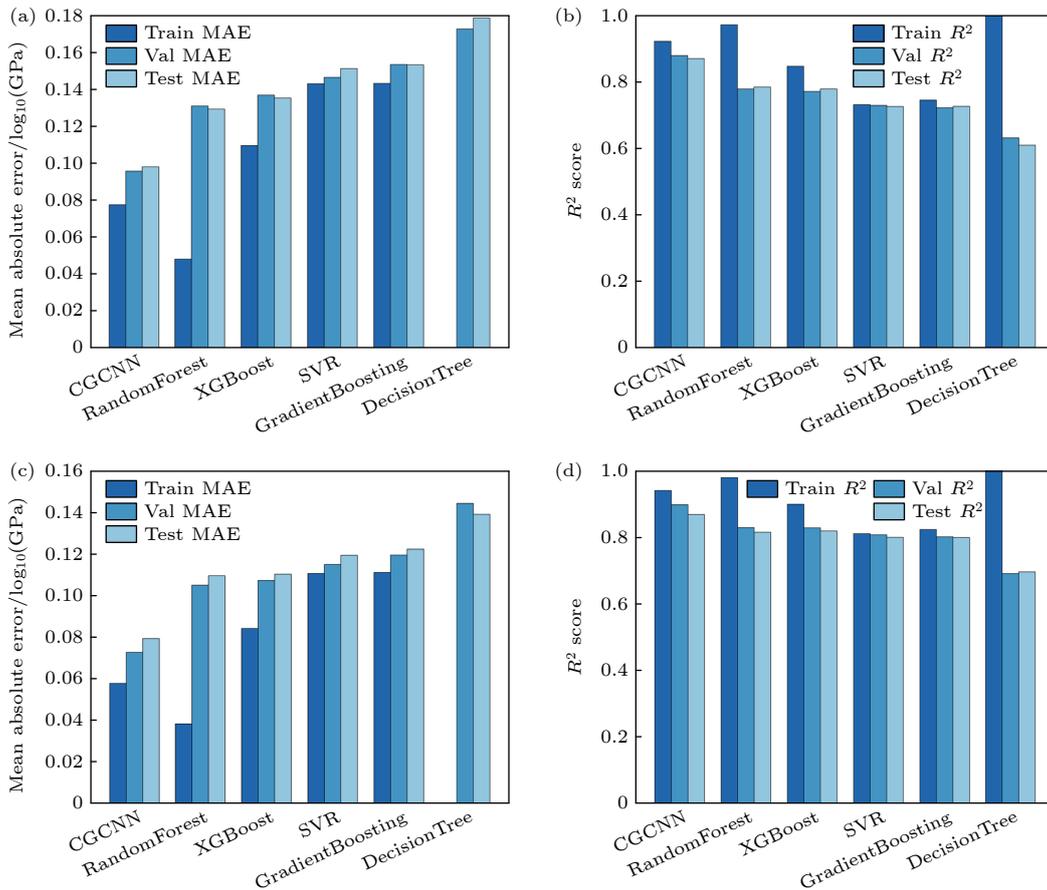


图 2 剪切模量 ((a), (b)) 和体模量 ((c), (d)) 的训练集 (Train)、验证集 (Val) 和测试集 (Test) 在晶体图卷积神经网络 (CGCNN)、随机森林 (random forest)、极限梯度提升 (XGBoost)、支持向量回归 (SVR)、梯度提升 (gradient boosting) 和决策树 (decision tree) 的平均绝对误差 (MAE) 和决定系数 ( $R^2$ )

Fig. 2. Mean absolute error (MAE) and coefficient of determination ( $R^2$ ) for the training set (Train), validation set (Val), and test set (Test) of shear modulus ((a), (b)) and bulk modulus ((c), (d)) in crystal graph convolutional neural network (CGCNN), random forest (RF), extreme gradient boosting (XGBoost), support vector regression (SVR), gradient boosting (GB), and decision tree (DT).

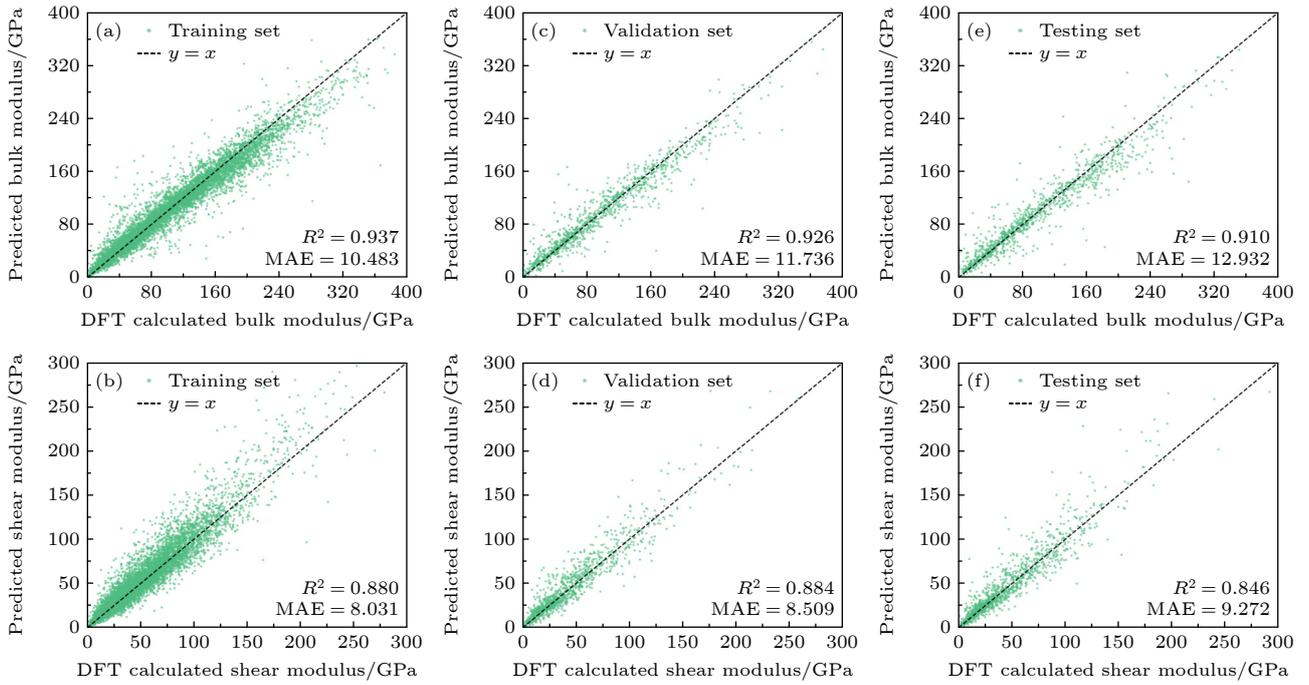


图3 利用CGCNN模型所预测的体积模量与剪切模量结果与DFT计算值对比分析 (a)和(b), (c)和(d)以及(e)和(f)分别为在训练集、验证集和测试集的结果

Fig. 3. Comparison between the volume modulus and shear modulus predicted by CGCNN model and the calculated values of DFT. (a) and (b), (c) and (d), (e) and (f) are the results in the train set, validation set, and test set, respectively.

进一步地, 对本文模型进行了评估, 将 Matbench v0.1<sup>[37]</sup> 中的 10987 个数据分为训练集、验证集和测试集训练模型, 图3为训练的CGCNN模型的弹性模量的训练集、验证集和测试集的结果. 从图3(a)和图3(b)可以看到, 在训练集上模型表现较好, DFT计算值和模型的输出结果接近, 说明模型预测结果与DFT计算值之间具有较高的线性相关性. 剪切模量和体积模量的 $R^2$ 分别为0.936和0.880, 并且都具有较低的MAE, 均不超过11, 表明模型在训练集上的表现良好. 和图3(a)与图3(b)相似, 图3(c)和图3(d)以及图3(e)和图3(f)分别为验证集和测试集的结果, MAE和 $R^2$ 略逊色于训练集结果且较为合理. 这是因为模型在训练集上更能准确拟合数据, 并且测试集的结果表明模型具备一定的泛化能力, 但在某些样本上存在预测误差. 通过训练集、验证集和测试集的对比分析, 可以看出模型在训练集上的拟合效果较好, 并且在测试集上的表现仍具有较高的预测精度和可靠性.

### 3.2 预测数据集的信息统计

3.1节已验证CGCNN在预测剪切模量和体积模量方面的出色性能. 基于此, 本文将该模型应

用于更大规模的数据集, 包括来自MPED数据集的54359种材料以及NED数据集筛选的26305种材料<sup>[38]</sup>. 对所预测的数据集进行了详细的统计分析, 主要目的在于: 通过晶系分布(含70+元素、7大晶系)、原子构型及成分特征验证数据代表性, 发现低对称性晶系(单斜/三斜)占比高、氧化物主导等特征符合材料学规律; 揭示一些结构特征, 低对称性晶系(单斜/三斜)占比高、氧化物主导等符合材料学规律, 兼具少量复杂结构. 需要说明的是, 尽管该数据集在材料组分与结构表征方面优势显著, 但关键力学参数(如剪切模量 $G$ 、体积模量 $B$ )的缺失限制了其深度应用. 为此, 本研究通过建立物性预测模型, 系统补充了缺失参数, 有效拓展了数据集的应用维度. 此外, 本文公开了元素频率表、晶系分布图等统计细节, 研究者可通过这些可视化数据快速定位目标样本(如特定元素或晶系材料), 显著降低数据筛选成本.

具体来说, 将MPED和NED两个数据集中的晶系、原胞内的原子数和元素出现次数分布用统计图刻画出来(如图4和图5). 其中, 图4是来源于MPED数据集的晶系、原胞原子数和元素统计结果, 图4(a)展示了数据集中7种晶系的分布情况. 单斜晶系(Monoclinic)占比最高, 为29.6%,

对应 16101 个结构; 三斜晶系 (Triclinic) 次之, 占比 26.4%, 对应 14461 个结构; 正交晶系 (Orthorhombic) 占比 19.4%, 包含 10858 个结构; 四方晶系 (Tetragonal) 和三方晶系 (Trigonal) 分别占 7.5% (4100 个) 和 7.5% (4077 个); 立方晶系 (Cubic) 和六方晶系 (Hexagonal) 分别占 6.9% (3721 个) 和 2.5% (1361 个). 图 4(b) 是原胞中原子数量的分布直方图, 总体来看, 原胞中原子数分布较为广泛, 其中原子数量较少 (小于 150) 的结构占据绝大多数. 随着原胞中原子数量的增加, 出现频率显著下降. 特别是当原胞中原子数超过 250 后, 频率显著降低, 但仍有少量复杂的晶体结构原胞原子数接近达到 444 个. 图 4(c) 展示了数据集中 77 种元素的出现频率分布. 横轴包含数据集中所有出现的元素, 按照出现频率从高到低排列, 纵轴为对应的元素出现次数. 其中, 氧 (O) 元素出现频率最高, 显

著高于其他元素, 表明氧化物在数据集中占据主导地位. 其他常见元素包括锂 (Li)、硫 (S)、镁 (Mg)、钠 (Na)、铁 (Fe) 等, 均在材料中多次出现. 稀有气体元素如氙 (Xe)、氪 (Kr) 和铑 (Rh) 等出现频率最低, 表明这些元素仅在极少数材料中存在. 分布呈现出明显的长尾效应, 大多数元素的频率集中在较低范围, 只有少数元素出现次数非常高. 另一方面, 图 5 来源于 NED 数据集. 图 5(a) 数据显示三斜和单斜晶系占据主导地位, 而六方晶系则极为稀少. 由图 5(b) 得知大部分材料的原胞原子数较低, 原子数在 3—40 之间的结构占主要比例. 图 5(c) 中, 氧化物主导了数据集的材料构成, 少数元素 (如氧、硒) 占据显著比例, 而稀土元素则频率较低. 这两组图表直观地展示了材料在晶体结构、原子数量和化学成分等方面的分布特征, 为进一步研究材料的性质提供了重要的统计依据.

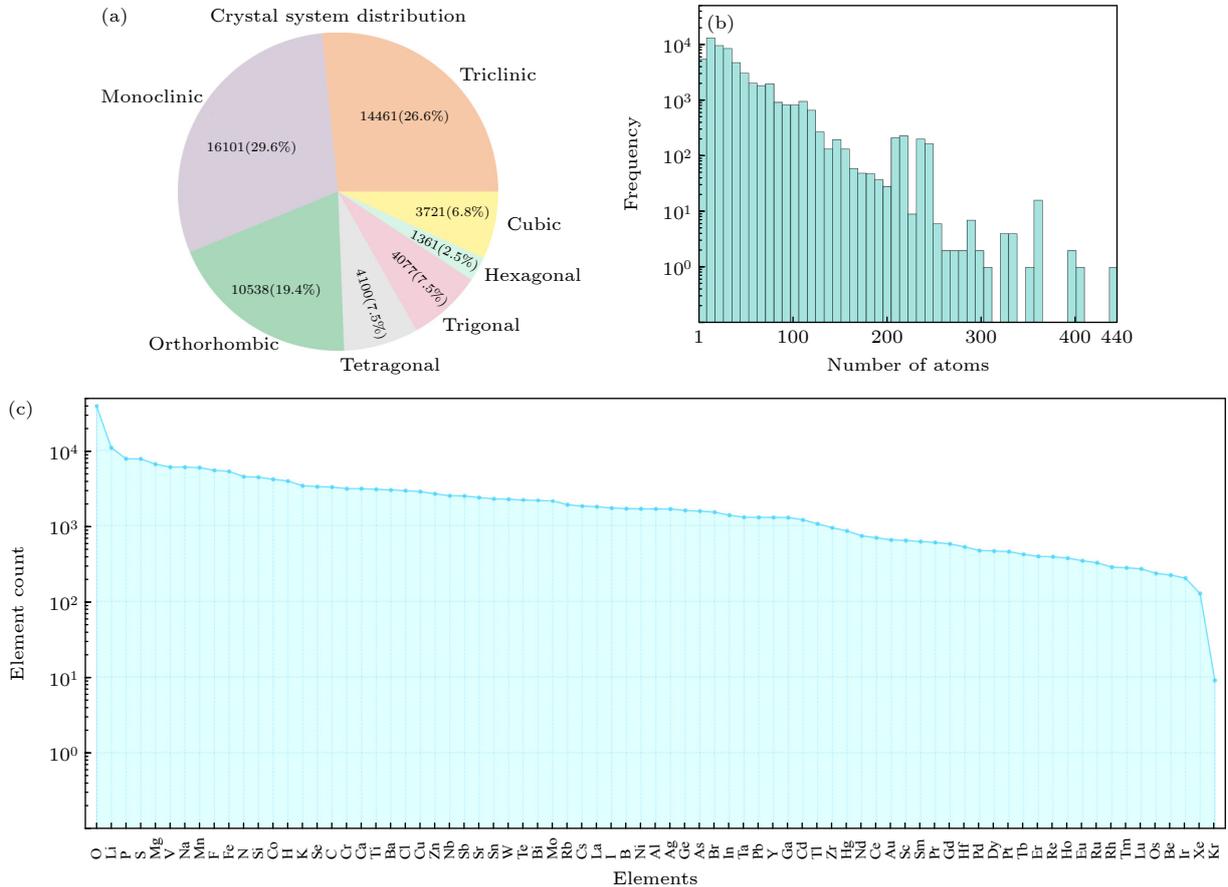


图 4 来自 MPED 数据集的预测数据集的统计分析 (a) 7 种晶系分布, 单斜晶系最常见 (16101 个结构), 其次是三斜晶系 (14461 个结构), 最少的是六方晶系 (1361 个结构); (b) 数据集原胞中的原子数范围 (1—444 个原子) 的分布; (c) 元素分布, 显示了 77 种不同元素的出现频率. 该数据集包括过渡金属、主族元素和稀土元素, 其中氧的出现频率最高

Fig. 4. Statistical analysis of predictive datasets from MPED: (a) The distribution of 7 crystal systems, with monoclinic being the most common (16101 structures), followed by triclinic (14461 structures), while hexagonal is the least one (1361 structures); (b) distribution of range of number of atoms in the primitive cell (1—444 atoms) across the dataset; (c) elemental distribution that illustrates the frequency of 77 distinct elements. The dataset encompasses transition metals, main group elements, and rare earth elements, with oxygen showing the highest frequency.

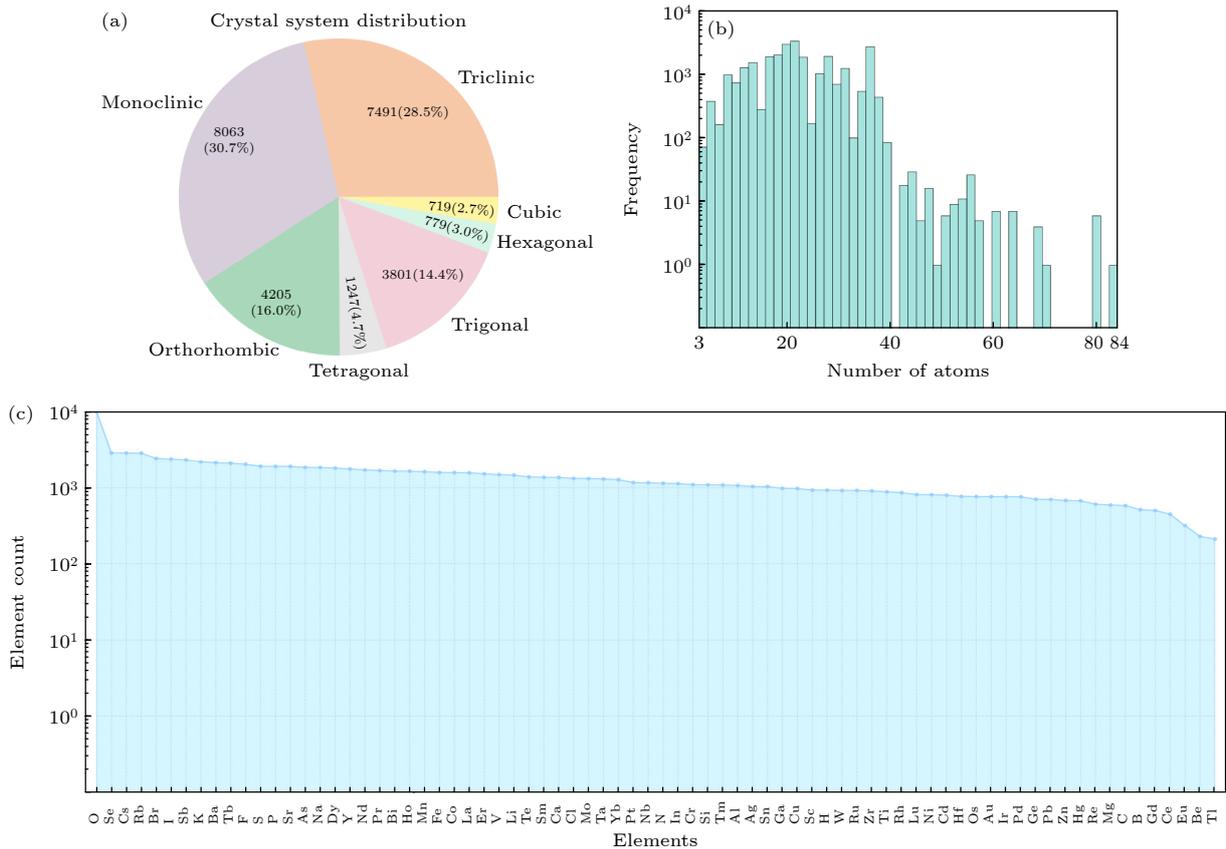


图 5 来自 NED 数据集所预测数据集的统计分析 (a) 7 种晶系分布, 单斜晶系最常见 (8063 个结构), 其次是三斜晶系 (7491 个结构), 最少的是六方晶系 (779 个结构); (b) 数据集原胞内原子数范围 (3—84 个原子) 的分布; (c) 元素分布, 显示了 76 种不同元素的出现频率. 该数据集包括过渡金属、主族元素和稀土元素, 其中氧的出现频率最高

Fig. 5. Statistical analysis of predictive datasets from NED: (a) The distribution of 7 crystal systems, with monoclinic being the most common (8063 structures), followed by triclinic (7491 structures), while hexagonal is the least one (779 structures); (b) distribution of the range of the number of atoms in the primitive cell (3–84 atoms) across the dataset; (c) elemental distribution illustrating the frequency of 76 distinct elements. The dataset encompasses transition metals, main group elements, and rare earth elements, with oxygen showing the highest frequency.

### 3.3 弹性性质的预测

进一步地, 利用 CGCNN 模型分别预测了 MPED 数据集与来自 NED 数据集中材料的剪切模量和体积模量, 获取了大量弹性性质相关的数据 (结果如表 A1 和表 A2). MPED 数据集和 NED 数据集的剪切模量和体模量统计分布以及相互关系分别如图 6 与图 7 所示, 通过散点图结合边缘直方图的方式直观呈现了数据特征. 通过直观的可视化手段, 将剪切模量与体积模量在两个不同数据集分布特征和关联性清晰呈现, 为进一步分析材料性能之间的关系提供了重要参考. 散点图的横轴为剪切模量, 纵轴为体积模量, 以不同颜色区分不同晶体结构. 具体来说, 图 6(a) 与图 7(a) 是所有材料的剪切模量与体模量分布, 图 6(b)—(h) 和图 7(b)—(h) 分别是三斜晶系、单斜晶系、正交晶系、三方晶系、四方晶系、六方晶系、立方晶系, 对称性由低到

高. 由散点图可以看出, 材料的剪切强度和体积模量存在紧密关联, 当材料剪切强度提高时, 其抵抗压缩的能力也会同步增强. 此外, 图 6 和图 7 还绘制了两条  $B/G$  比值的线, Pugh 比率 ( $B/G$ ) 在之前的工作中认为与晶体化合物的延展性相关, 并且进一步与泊松比相关<sup>[12]</sup>. 条形图展示了各晶系材料的剪切模量和体模量的统计分布. 分布显示, 大多数材料的剪切模量和体积模量集中在 10—100 GPa 的区域. 同时从各晶系的数据分布可以看出, 高对称性晶系 (如立方晶系和六方晶系) 的数据点更集中分布在图的右上方区域, 表现出更高的剪切模量和体积模量. 这一结果为进一步研究材料性能的关联性提供了重要参考. 最后, 若想查看更多数据, 请访问数据集 <https://doi.org/10.57760/sciencedb.j00213.00104>.

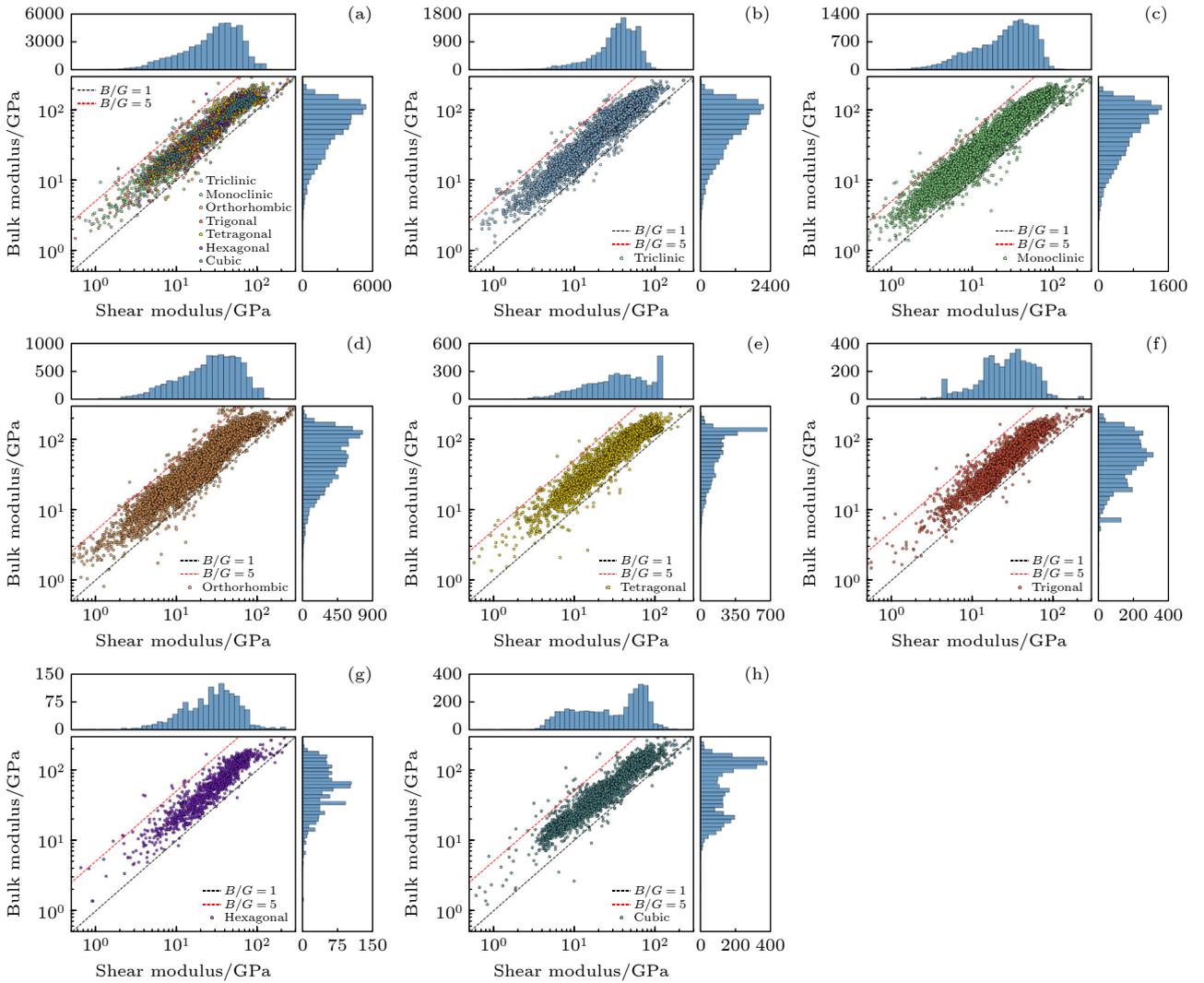


图6 MPED数据集中不同材料的剪切模量与体模量分布 (a) 所有材料的剪切模量与体模量分布, 不同颜色代表不同的晶系; (b) 三斜晶系; (c) 单斜晶系; (d) 正交晶系; (e) 三方晶系; (f) 四方晶系; (g) 六方晶系; (h) 立方晶系. 条形图展示了各晶系材料的剪切模量和体模量的统计分布

Fig. 6. Shear modulus and bulk modulus distributions of different materials in the MPED dataset: (a) Shear modulus vs. bulk modulus distributions for all materials, with different colors representing different crystal systems; (b) triclinic; (c) monoclinic; (d) orthorhombic; (e) trigonal; (f) tetragonal; (g) hexagonal; (h) cubic. The bar graphs show the statistical distribution of shear and bulk moduli for each crystal system material.

## 4 结论

本研究基于 CGCNN 模型, 对材料的弹性性能进行了系统的训练、预测和分析. 基于 CGCNN 训练了两个弹性模量模型, 同时结合来自 MPED 和 Merchant 等<sup>[38]</sup>发现的新材料, 针对新材料的剪切模量和体积模量进行了深入探索, 最终形成了一个包含 80664 种晶体弹性性质的数据集. 结果表明, CGCNN 能够精准捕捉晶体结构中局部化学环境的特征, 对剪切模量和体积模量的预测具有较高的准确性, 其中 MAE 值低于 13,  $R^2$  值接近 1, 充分验证了模型的可靠性和泛化能力.

通过对两个预测数据集的统计分析, 发现低对称性的晶系材料占比更高, 氧化物在化学组成中占据主导地位, 原胞原子数主要集中于较低范围, 稀土元素的出现频率显著低于常见元素. 这些统计结果不仅符合自然界中的材料分布特征, 也为进一步研究提供了重要依据. 弹性模量的可视化结果显示, 剪切模量与体积模量之间具有显著的正相关性, 体现了它们在物理性质上的耦合特征.

为了丰富材料的弹性性能, 基于剪切模量和体积模量计算声速、泊松比及德拜温度等物理参数, 为材料性能的多维度研究提供了基础支持. 对超过 8 万种稳定材料结构的弹性性能预测结果表明,

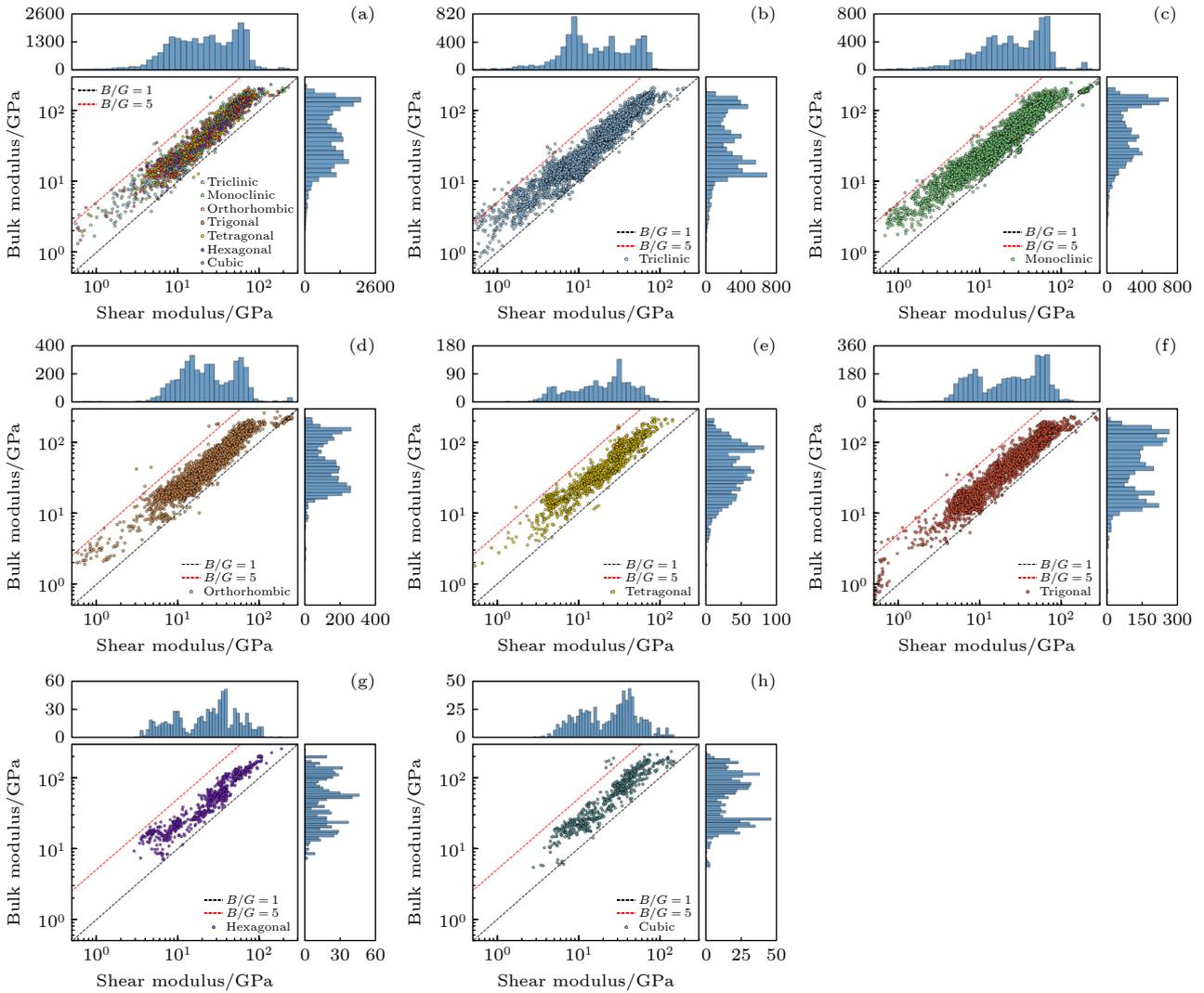


图 7 NED 数据集中各晶体结构材料的模量分布 (a) 整体剪切模量-体模量分布 (颜色区分晶系); (b) 三斜晶系; (c) 单斜晶系; (d) 正交晶系; (e) 三方晶系; (f) 四方晶系; (g) 六方晶系; (h) 立方晶系. 条形图统计了各晶系材料的剪切模量和体模量分布

Fig. 7. Distribution of moduli for various crystal structure materials in the NED dataset: (a) Overall shear modulus-bulk modulus distribution (color-coded by crystal system); (b) triclinic system; (c) monoclinic system; (d) orthorhombic system; (e) trigonal system; (f) tetragonal system; (g) hexagonal system; (h) cubic system. Bar charts illustrate the distribution of shear modulus and bulk modulus for materials in each crystal system.

CGCNN 在大规模数据集上的适用性和高效性, 为加速新材料发现和优化设计提供了强有力的工具. 因此, 本研究训练了两个弹性模量 CGCNN 模型, 同时证明了 CGCNN 在材料弹性性能预测中的强大能力, 结合大规模数据分析揭示了材料性质的分布规律与物理关联性, 为材料科学领域提供了新的研究思路和方法.

感谢西安交通大学高性能计算平台的支持.

### 数据可用性声明

支撑本研究成果的数据集可在科学数据银行 <https://doi.org/10.57760/sciencedb.j00213.00104> 中访问获取.

### 附录 MPED 与 NED 数据集中无机晶体材料的结构参数与物理性能摘要表

表 A1 和表 A2 分别是 MPED 和 NED 数据集无机晶体材料基础物理特性及预测值 (部分). 其中,  $N$ ,  $\rho$ ,  $V$  和  $M$  分别是原胞中的原子数、密度、原胞的体积和原子总质量;  $B$  和  $G$  分别为通过 CGCNN 网络得到的体积模量和剪切模量的预测值;  $v_l$ ,  $v_t$ ,  $v_s$ ,  $\nu$  和  $\theta_D$  分别为纵向声速、横向声速、平均声速、泊松比和德拜温度, 并可由 (1)–(5) 式求得. 完整数据可以在 <https://doi.org/10.57760/sciencedb.j00213.00104> 中下载.

表 A1 MPED 数据集无机晶体材料基础物理特性及预测值 (部分). 这里, ID-number 和 Formula 分别是材料编号和化学式

Table A1. Fundamental physical properties (partial) and predicted values of inorganic crystalline materials from MPED datasets. The CIF files of these materials were obtained from the Materials Project. Here, ID-number and Formula represent the material ID and chemical formula, respectively.

ID-number	Formula	$N$	$\rho$	$V$	$M$	$B$	$G$	$v_l$	$v_t$	$v_s$	$\nu$	$\theta_D$
mp-1000	BaTe	2	4.938	89.094	264.927	31.764	23.469	2180.121	3573.541	2407.744	0.204	160.498
mp-10009	GaTe	8	5.1549	254.251	789.292	24.095	16.757	1802.955	3001.379	1994.276	0.218	93.722
mp-1001012	Sc <sub>2</sub> ZnSe <sub>4</sub>	14	3.254	289.440	567.162	53.623	32.397	3155.374	5454.806	3502.473	0.249	157.640
mp-1001015	Y <sub>2</sub> ZnS <sub>4</sub>	14	3.675	335.691	742.961	60.652	25.843	2651.771	5087.157	2967.069	0.314	127.104
mp-1001016	Sc <sub>2</sub> ZnSe <sub>4</sub>	14	4.687	333.879	942.322	54.940	22.543	2193.172	4258.659	2455.876	0.320	105.395
mp-1001019	MgSc <sub>2</sub> Se <sub>4</sub>	14	4.086	349.578	860.114	52.875	22.985	2371.850	4521.352	2652.741	0.310	112.113
mp-1001021	Y <sub>2</sub> ZnSe <sub>4</sub>	14	4.811	385.950	1118.121	55.070	22.939	2183.662	4219.640	2444.462	0.317	99.958
mp-1001023	BeC <sub>2</sub>	6	1.879	58.402	66.067	132.395	102.494	7386.608	11967.830	8148.016	0.192	625.248
mp-1001024	Y <sub>2</sub> MgS <sub>4</sub>	14	3.173	345.765	660.753	56.994	26.037	2864.435	5375.943	3200.229	0.302	135.747
mp-1001034	MgIn <sub>2</sub> Se <sub>4</sub>	14	5.031	376.146	1139.562	39.515	21.476	2066.136	3680.578	2299.251	0.270	94.830
mp-1001069	Li <sub>48</sub> P <sub>16</sub> S <sub>61</sub>	125	1.743	2652.952	2784.713	19.812	7.267	2041.845	4114.028	2291.557	0.337	49.283
mp-1001079	LiC <sub>2</sub> N <sub>2</sub>	10	1.505	130.116	117.952	56.823	20.405	3681.742	7471.454	4133.696	0.340	242.869
mp-10013	SnS	2	3.596	69.620	150.775	17.613	5.617	1249.772	2642.016	1406.249	0.356	101.772
mp-1001594	C <sub>4</sub> O <sub>3</sub>	84	1.656	1155.735	1152.492	19.101	12.904	2791.530	4682.464	3090.023	0.224	87.663
mp-1001604	LuTiS <sub>2</sub>	4	7.377	99.825	443.480	49.490	20.396	1662.754	3224.127	1861.754	0.319	119.486
mp-1001611	LuTiSe <sub>2</sub>	4	8.001	111.508	537.270	43.737	22.793	1687.844	3043.848	1880.122	0.278	116.295
mp-1001780	LuCuS <sub>2</sub>	4	6.522	77.056	302.643	74.239	35.316	2327.021	4313.132	2597.493	0.295	181.731
mp-1001786	LiScS <sub>2</sub>	4	2.700	71.362	116.027	58.972	36.372	3670.409	6309.130	4072.100	0.244	292.285
mp-1001790	LiO <sub>3</sub>	4	2.130	42.828	54.939	46.463	28.415	3652.317	6292.720	4052.874	0.246	344.878
mp-1001831	LiB	4	2.099	28.090	35.504	111.075	134.490	8004.910	11762.661	8727.079	0.069	854.731

表 A2 NED 数据集无机晶体材料基础物理特性及预测值 (部分). 这里, Filename 表示文件名

Table A2. Basic physical properties and predicted values of inorganic crystalline materials (part) from NED datasets. Here, Filename represents the file name.

Filename	$N$	$\rho$	$V$	$M$	$G$	$B$	$v_l$	$v_t$	$v_s$	$\nu$	$\theta_D$
FIrS	3	7.798	51.805	243.280	28.413	54.027	3433.128	1908.824	2125.825	0.276	244.862
AuGeP	3	7.381	67.619	300.580	23.064	55.970	3427.627	1767.655	1979.213	0.319	208.603
GdHO	3	7.384	39.190	174.257	62.945	113.409	5169.778	2919.774	3247.588	0.266	410.537
LiPrPtSn	4	9.285	82.565	461.643	31.112	78.216	3590.578	1830.554	2051.127	0.324	222.617
ErLiPdSn	4	8.792	75.424	399.330	36.874	81.235	3851.257	2047.962	2288.361	0.303	255.968
BaBiHgNa	4	6.817	138.827	569.887	11.187	24.989	2419.500	1281.048	1431.855	0.305	130.688
BeGeHLA	4	5.801	63.421	221.566	49.688	90.981	5206.069	2926.621	3256.448	0.269	385.920
AlHKSb	4	3.004	104.402	188.848	14.352	23.461	3765.877	2185.915	2425.631	0.246	243.454
EuHgNaSb	4	7.135	115.739	497.304	15.654	30.762	2690.122	1481.228	1650.873	0.282	160.097
LiNiSmSn	4	7.617	72.963	334.704	36.441	70.798	3958.873	2187.199	2437.061	0.280	275.632
DyLiPdSn	4	8.557	76.573	394.571	35.786	81.074	3879.627	2045.067	2286.509	0.308	254.475
N <sub>2</sub> SSe <sub>2</sub>	5	2.175	166.436	217.998	2.459	2.521	1632.981	1063.352	1165.878	0.132	107.904
LiNaSe <sub>2</sub> Zn	5	3.916	107.396	253.260	17.754	31.924	3767.961	2129.286	2368.236	0.265	253.647
BrGeLa <sub>2</sub> Rh	5	6.436	137.585	533.260	27.302	50.532	3675.249	2059.620	2292.318	0.271	226.057
CsHgNaS <sub>2</sub>	5	4.774	146.289	420.615	9.852	18.449	2572.057	1436.510	1599.245	0.273	154.518
AlAs <sub>2</sub> CsMg	5	3.863	143.606	334.035	20.025	28.181	3769.434	2276.944	2517.185	0.213	244.713
Br <sub>2</sub> GeSmY	5	4.803	163.091	471.714	19.469	32.496	3488.675	2013.383	2235.315	0.250	208.287
As <sub>2</sub> Ca <sub>2</sub> Sr	5	3.392	155.481	317.619	28.093	37.392	4697.360	2877.774	3176.871	0.200	300.774
KLiMnTe <sub>2</sub>	5	3.860	153.218	356.177	12.890	26.438	3361.705	1827.331	2038.601	0.290	193.953
Al <sub>2</sub> C <sub>2</sub> Yb	5	6.426	64.862	251.024	88.642	125.838	6162.150	3713.927	4106.697	0.215	520.350

参考文献

- [1] Koester S J, Schaub J D, Dehlinger G, Chu J O 2006 *IEEE J. Sel. Top. Quantum Electron.* **12** 1489
- [2] Seo D, Gregory J, Feldman L, Tolk N, Cohen P 2011 *Phys. Rev. B* **83** 195203
- [3] Parola S, Julián-López B, Carlos L D, Sanchez C 2016 *Adv. Funct. Mater.* **26** 6506
- [4] Sanchez C, Lebeau B, Chaput F, Boilot J P 2003 *Adv. Mater.* **15** 1969
- [5] Beekman M, Cahill D G 2017 *Cryst. Res. Technol.* **52** 1700114
- [6] Tan J C, Cheetham A K 2011 *Chem. Soc. Rev.* **40** 1059
- [7] Reddy C M, Krishna G R, Ghosh S 2010 *CrystEngComm* **12** 2296
- [8] Keyes R W 1968 *Solid State Physics.* **20** 37
- [9] Wang X, Shu G, Zhu G, Wang J S, Sun J, Ding X, Li B, Gao Z 2024 *Mater. Today Phys.* **48** 101549
- [10] Chen L, Tran H, Batra R, Kim C, Ramprasad R 2019 *Comput. Mater. Sci.* **170** 109155
- [11] Lake G J, Thomas A G 1967 *Proc. R. Soc. London, Ser. A* **300** 108
- [12] Pugh S 1954 *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science* **45** 823
- [13] Niu H, Chen X Q, Liu P, Xing W, Cheng X, Li D, Li Y 2012 *Sci. Rep.* **2** 718
- [14] Gschneidner Jr K, Russell A, Pecharsky A, Morris J, Zhang Z, Lograsso T, Hsu D, Chester Lo C, Ye Y, Slager A, Kesse D 2003 *Nat. Mater.* **2** 587
- [15] Greaves G N, Greer A L, Lakes R S, Rouxel T 2011 *Nat. Mater.* **10** 823
- [16] de Jong M, Olmsted D L, van de Walle A, Asta M 2012 *Phys. Rev. B* **86** 224101
- [17] de Jong M, van der Zwaag S, Sluiter M 2012 *Int. J. Mater. Res.* **103** 972
- [18] Payne M C, Teter M P, Allan D C, Arias T, Joannopoulos J D 1992 *Rev. Mod. Phys.* **64** 1045
- [19] Wang J, Yang X, Zeng Z, Zhang X, Zhao X, Wang Z 2017 *Comput. Mater. Sci.* **138** 135
- [20] Plimpton S 1995 *J. Comput. Phys.* **117** 1
- [21] Grigoras S, Gusev A, Santos S, Suter U 2002 *Polymer* **43** 489
- [22] Belytschko T, Black T 1999 *Int. J. Numer. Methods Eng.* **45** 601
- [23] Zhao Y, Yuan K, Liu Y, Louis S Y, Hu M, Hu J 2020 *J. Phys. Chem. C* **124** 17262
- [24] Chibani S, Coudert F X 2020 *APL Mater.* **8** 080701
- [25] Chew A K, Sender M, Kaplan Z, Chandrasekaran A, Chief Elk J, Browning A R, Kwak H S, Halls M D, Afzal M A F 2024 *J. Cheminform.* **16** 31
- [26] Xie T, Grossman J C 2018 *Phys. Rev. Lett.* **120** 145301
- [27] Karamad M, Magar R, Shi Y, Siahrostami S, Gates I D, Barati Farimani A 2020 *Phys. Rev. Mater.* **4** 093801
- [28] Choudhary K, DeCost B 2021 *npj Comput. Mater.* **7** 185
- [29] Louis S Y, Zhao Y, Nasiri A, Wang X, Song Y, Liu F, Hu J 2020 *Phys. Chem. Chem. Phys.* **22** 18141
- [30] Ruff R, Reiser P, Stühmer J, Friederich P 2024 *Digital Discovery* **3** 594
- [31] Omee S S, Fu N, Dong R, Hu M, Hu J 2024 *npj Comput. Mater.* **10** 144
- [32] de Jong M, Chen W, Angsten T, Jain A, Notestine R, Gamst A, Sluiter M, Ande C K, van der Zwaag S, Plata J J 2015 *Sci. Data* **2** 150009
- [33] Jain A, Ong S P, Hautier G, et al. 2013 *APL Mater.* **1** 011002
- [34] Curtarolo S, Setyawan W, Hart G L, et al. 2012 *Comput. Mater. Sci.* **58** 218
- [35] Saal J E, Kirklin S, Aykol M, Meredig B, Wolverton C 2013 *JOM* **65** 1501
- [36] Atomly 2025 <https://atomly.net> Accessed: 2025-01-13
- [37] Dunn A, Wang Q, Ganose A, Dopp D, Jain A 2020 *npj Comput. Mater.* **6** 138
- [38] Merchant A, Batzner S, Schoenholz S S, Aykol M, Cheon G, Cubuk E D 2023 *Nature* **624** 80
- [39] Shuai C, Liu W, Li H, Wang K, Zhang Y, Xie T, Chen L, Hou H, Zhao Y 2023 *Int. J. Plast.* **170** 103772
- [40] Jia T, Chen G, Zhang Y 2017 *Phys. Rev. B* **95** 155206
- [41] Qin G, Huang A, Liu Y, Wang H, Qin Z, Jiang X, Zhao J, Hu J, Hu M 2022 *Mater. Adv.* **3** 6826
- [42] Belomestnykh V N, Tesleva E P 2004 *Tech. phys.* **49** 1098
- [43] Belomestnykh V N 2004 *Tech. Phys. Lett.* **30** 91
- [44] Toher C, Oses C, Plata J J, Hicks D, Rose F, Levy O, de Jong M, Asta M, Fornari M, Buongiorno Nardelli M, Curtarolo S 2017 *Phys. Rev. Mater.* **1** 015401
- [45] Wang H, Duan Z, Guo Q, Zhang Y, Zhao Y 2023 *CMC-Comput. Mater. Continua* **77** 1393

SPECIAL TOPIC—Atomic, molecular and materials properties data

# Machine learning-driven elasticity prediction in advanced inorganic materials via convolutional neural networks\*

LIU Yujie<sup>1) #</sup> WANG Zhenyu<sup>1) #</sup> LEI Hang<sup>1) #</sup> ZHANG Guoyu<sup>1)</sup>XIAN Jiawei<sup>2) †</sup> GAO Zhibin<sup>1) ‡</sup> SUN Jun<sup>3)</sup>SONG Haifeng<sup>2)</sup> DING Xiangdong<sup>3)</sup>

1) (State Key Laboratory of Porous Metal Materials, School of Materials Science and Engineering, Xi'an Jiaotong University, Xi'an 710049, China)

2) (National Key Laboratory of Computational Physics, Institute of Applied Physics and Computational Mathematics, Beijing 100088, China)

3) (State Key Laboratory for Mechanical Behavior of Materials, School of Materials Science and Engineering, Xi'an Jiaotong University, Xi'an, 710049, China)

( Received 25 January 2025; revised manuscript received 18 March 2025 )

## Abstract

Inorganic crystal materials have shown extensive application potential in many fields due to their excellent physical and chemical properties. Elastic properties, such as shear modulus and bulk modulus, play an important role in predicting the electrical conductivity, thermal conductivity and mechanical properties of materials. However, the traditional experimental measurement method has some problems such as high cost and low efficiency. With the development of computational methods, theoretical simulation has gradually become an effective alternative to experiments. In recent years, graph neural network-based machine learning methods have achieved remarkable results in predicting the elastic properties of inorganic crystal materials, especially, crystal graph convolutional neural networks (CGCNNs), which perform well in the prediction and expansion of material data.

In this study, two CGCNN models are trained by using the shear modulus and bulk modulus data of 10987 materials collected in the Matbench v0.1 dataset. These models show high accuracy and good generalization ability in predicting shear modulus and bulk modulus. The mean absolute error (MAE) is less than 13 and the coefficient of determination ( $R^2$ ) is close to 1. Then, two datasets are screened for materials with a band gap between 0.1 and 3.0 eV and the compounds containing radioactive elements are excluded. The dataset consists of two parts: the first part is composed of 54359 crystal structures selected from the Materials Project database, which constitute the MPED dataset; the second part is the 26305 crystal structures discovered by Merchant et al. (2023 *Nature* **624** 80) through deep learning and graph neural network methods, which constitute the NED dataset. Finally, the shear modulus and bulk modulus of 80664 inorganic crystals are predicted in this study. This work enriches the existing material elastic data resources and provides more data support for material design. All the data presented in this paper are openly available at <https://doi.org/10.57760/sciencedb.j00213.00104>.

**Keywords:** inorganic crystal materials, elastic modulus, machine learning, data prediction

**PACS:** 07.05.Tp, 61.72.-y, 62.20.D-

**DOI:** 10.7498/aps.74.20250127

**CSTR:** 32037.14.aps.74.20250127

\* Project supported by the National Natural Science Foundation of China (Grant Nos. 12104356, 52250191), the Funding of National Key Laboratory of Computational Physics, China, and the National Key Research and Development Program of China (Grant No. 2023YFB4604100).

# These authors contributed equally.

† Corresponding author. E-mail: [xian\\_jiawei@iapcm.ac.cn](mailto:xian_jiawei@iapcm.ac.cn)

‡ Corresponding author. E-mail: [zhibin.gao@xjtu.edu.cn](mailto:zhibin.gao@xjtu.edu.cn)

## 卷积神经网络辅助无机晶体弹性性质预测

刘宇杰 王振宇 雷航 张国宇 戚家伟 高志斌 孙军 宋海峰 丁向东

## Machine learning-driven elasticity prediction in advanced inorganic materials via convolutional neural networks

LIU Yujie WANG Zhenyu LEI Hang ZHANG Guoyu XIAN Jiawei GAO Zhibin SUN Jun SONG Haifeng DING Xiangdong

引用信息 Citation: *Acta Physica Sinica*, 74, 120702 (2025) DOI: 10.7498/aps.74.20250127

CSTR: 32037.14.aps.74.20250127

在线阅读 View online: <https://doi.org/10.7498/aps.74.20250127>

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

#### 基于机器学习的无机磁性材料磁性基态分类与磁矩预测

Classification of magnetic ground states and prediction of magnetic moments of inorganic magnetic materials based on machine learning

物理学报. 2022, 71(6): 060202 <https://doi.org/10.7498/aps.71.20211625>

#### 基于机器学习和第一性原理计算的Janus材料预测

Prediction of magnetic Janus materials based on machine learning and first-principles calculations

物理学报. 2024, 73(23): 230201 <https://doi.org/10.7498/aps.73.20241278>

#### 基于机器学习的菱形穿孔石墨烯负泊松比效应预测与优化

Prediction and optimization of negative Poisson's ratio in rhombic perforated graphene based on machine learning

物理学报. 2025, 74(9): 096201 <https://doi.org/10.7498/aps.74.20241624>

#### 胶体聚合物弹性模量的微观理论: 键长的效应

Microscopic theory for elastic modulus of colloidal polymers: Effect of bond length

物理学报. 2021, 70(12): 126401 <https://doi.org/10.7498/aps.70.20210128>

#### 机器学习辅助的WC-Co硬质合金硬度预测

Hardness prediction of WC-Co cemented carbide based on machine learning model

物理学报. 2024, 73(12): 126201 <https://doi.org/10.7498/aps.73.20240284>

#### 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>