# Digital Discovery



**PAPER** 

View Article Online
View Journal | View Issue



Cite this: Digital Discovery, 2024, 3, 769

Received 15th January 2024 Accepted 27th February 2024

DOI: 10.1039/d4dd00024b

rsc.li/digitaldiscovery

# Learning conditional policies for crystal design using offline reinforcement learning†

Prashant Govindarajan, (1)\*\* \* Santiago Miret, (1)\*\* Jarrid Rector-Brooks, Carana Mariano Phielipp, Danarthanan Rajendran and Sarath Chandar

Navigating through the exponentially large chemical space to search for desirable materials is an extremely challenging task in material discovery. Recent developments in generative and geometric deep learning have shown promising results in molecule and material discovery but often lack evaluation with high-accuracy computational methods. This work aims to design novel and stable crystalline materials conditioned on a desired band gap. To achieve conditional generation, we: (1) formulate crystal design as a sequential decision-making problem, create relevant trajectories based on high-quality materials data, and use conservative Q-learning to learn a conditional policy from these trajectories. To do so, we formulate a reward function that incorporates constraints for energetic and electronic properties obtained directly from density functional theory (DFT) calculations; (2) evaluate the generated materials from the policy using DFT calculations for both energy and band gap; (3) compare our results to relevant baselines, including behavioral cloning and unconditioned policy learning. Our experiments show that conditioned policies achieve targeted crystal design and demonstrate the capability to perform crystal discovery evaluated with accurate and computationally expensive DFT calculations.

# 1 Introduction

The widespread enthusiasm in exploiting artificial intelligence (AI) for scientific discovery1 has resulted in various methodologies to integrate existing scientific knowledge and large databases to design and test new hypotheses more quickly. Recently, AI has shown favorable results in expediting the discovery of new chemical structural entities (e.g., small molecules, materials, and polymers).2-5 While several studies have focused on small molecule design for applications in drug discovery, there has also been an upsurge in attention for AI-based material discovery. 6-9 Among solid-state materials, crystalline substances are abundant in nature and are extensively used in industry for designing batteries, semiconductors, and photovoltaic systems. The set of known and experimentally observed crystalline materials is an infinitesimally tiny fraction (around 200 000) of the exponentially large chemical space spanning over 100 elements in the periodic table and 230 space groups in 3 dimensions.10,11 Determining a way to navigate through this large space to select chemical candidates with desired properties would be immensely beneficial for a plethora of applications like designing energy-efficient semiconductors and combatting climate change.

Besides the complex nature of the chemical space, designing stable crystalline materials using computational chemistry is a long-standing challenge primarily due to the time-consuming density functional theory (DFT) calculations to estimate energetic and electronic properties of materials. Previous works have utilized generative adversarial networks (GANs),12 diffusion models, 13,14 and reinforcement learning (RL), 15,16 in addition to advanced crystal representation schemes for generating crystals. 17,18 However, we identify two major gaps in the existing literature for AI-based material discovery. Firstly, most methods do not incorporate quantum mechanics-based first-principles calculations in the learning model, and instead use ML approximators. Studies that incorporate DFT computations in their ML pipeline for material design usually focus on smaller and very specific chemical systems (with limited number of elements or constraints on the space group) that might not generalize well to diverse chemical systems. 15,19 Secondly, state-of-the-art generative AI methods, such as diffusion models, predict the identities and positions of all atoms simultaneously, which is orthogonal to sequence based RL methods that also have more established exploration methods applicable to vast search spaces.

In this work, we develop a model that learns to sequentially construct crystal skeleton graphs by optimizing for both lower formation energy and desired band gap value (energy gap between the valence and conduction bands in solids), as computed by DFT. In our case, the crystal lattice parameters and

<sup>&</sup>quot;Mila-Quebec AI Institute, Polytechnique, Montréal, Canada. E-mail: prashant. govindarajan@mila.quebec

bIntel Labs, USA

<sup>&#</sup>x27;Mila-Quebec AI Institute, Université de Montréal, Canada

<sup>†</sup> Electronic supplementary information (ESI) available. See DOI: https://doi.org/10.1039/d4dd00024b

positions of atomic sites are known beforehand (crystal skeleton), and the task is to learn a conditional policy that can sequentially fill atoms to generate a stable and valid crystal with a desired band gap energy. To alleviate the issue of time-consuming DFT calculations when integrated into the scientific discovery loop, we apply offline reinforcement learning using the conservative Qlearning (CQL) approach,20 which is known to mitigate overestimation and out-of-distribution issues when agents are trained with static datasets in an offline manner. We construct a state transition dataset from high-quality nonmetallic crystal structures present in the Materials Project database. The reward function is carefully formulated to penalize high energies and large deviations from the desired band gap. Further, we leverage an expressive graph neural network (GNN) for crystal representation that ensures invariance to periodicity, translation, and rotation. Through our work, we aim to accelerate the process of high-throughput virtual screening (HTVS) for materials,21 where usually elements are combinatorially substituted in a known crystal structure and optimized using DFT calculations. Overall, our contributions are three-fold, as follows:

- (1) DFT evaluation of crystals designed with reinforcement learning: our targeted formulation of the reward function for offline RL is crafted from formation energy (per atom) and band gap values computed using first-principles DFT calculations, widely used in computational chemistry. The reward function penalizes high energy and large deviations from the desired band gap, resulting in a policy conditioned on a target band gap value.
- (2) Conservative offline reinforcement learning approach: using CQL as our offline RL framework, we show that conservatism, combined with the right amount of importance for the energy and band gap terms in the reward function, can result in an intuitive approach for generating crystals with a favorable shift in the distribution of properties of interest. Considering our task has a very sparse reward scheme, allows no exploration, and has a high dimensional action space and limited data, we highlight the important challenges that could be addressed in the future.
- (3) Open-source crystal structure design trajectory data: to ensure consistency in our reward calculation, we evaluate  $\sim$ 20k crystal structures using the Quantum Espresso<sup>22</sup> package for DFT calculation and subsequently construct offline RL trajectories based on the data. We release the dataset of trajectories and calculations as part of the paper to enable research to further improve our work. We use an open-source DFT calculator that is highly reproducible and consistent for all the structures evaluated. Prior work used different types of proprietary DFT software, which is difficult for the research community to reproduce.

### 2 Related works

#### 2.1 Automated materials design

Prior work has explored the application of various types of methods to crystal structure design, including evolutionary algorithms, simulated annealing, particle swarm optimization, and high-throughput screening.<sup>23–25</sup> Machine learning based methods have been more recently applied, primarily to molecular design problems, but also to periodic crystal structures.<sup>17,26</sup> Moreover, there have been notable works using machine learning

based methods to approximate the evaluation of material properties and behaviors.8,27 This includes approximating DFT outputs directly for different systems, such as ground-state crystal structures for a variety of applications, such as catalysts.28,29 The recent progress in graph neural networks and generative models has led to their successful application in materials design.29,30 GANs have been well explored for crystal structure design. 12,19,31 However, these approaches restrict the complexity of the problem to a fixed crystal system or a smaller chemical space.11 proposed a physics-guided GAN model using convolutional layers to learn the generative distribution of stable crystals, and the evaluation of generated crystals was done using DFT. CDVAE13 introduced a diffusion-based framework with highly expressive graph representation learning techniques to generate stable and valid crystal structures in 3 dimensions.32 used their Distributional Graphormer to generate structures of carbon polymorphs with the desired band gap.15 focused on building an online RL framework with DFT integrated reward function for surface reconstructions. However, they use the tightbinding version of DFT (DFTB), whose accuracy is lower than full DFT calculations. Other relevant works include ref. 33-35 and 36

#### 2.2 Offline reinforcement learning

Offline RL<sup>37,38</sup> enables for learning an optimal policy directly from existing trajectories, making it possible to utilize knowledge from known crystal structures. The ability to learn from previously determined crystal structures reduces the need for costly DFT calculations during training, which are necessary for online RL methods. Many recently proposed offline RL methods focus on managing distribution shift between the offline data and the learned policy,<sup>39-41</sup> with Conservative Q-Learning (CQL)<sup>20</sup> proving to be a particularly robust approach. CQL has shown success in training large capacity models and performing better with suboptimal data, which makes it a particularly good fit for our crystal structure design case.

# 3 Background

#### 3.1 Crystals

Solid-state crystals are characterized by ordered and periodic arrangement of atoms in 3 dimensional space. They consist of unit cells, which are the smallest group of atoms that form the repeating pattern of the crystal. A crystal's composition and arrangement of atoms give rise to distinct electronic properties usually determined by experimental or simulation-based density functional theory (DFT) calculations. In 3 dimensions, we can mathematically express the unit cell  $\boldsymbol{U}$  as follows.

$$U = \{w_1 l_1 + w_2 l_2 + w_3 l_3 | 0 \le w_i < 1\}, \tag{1}$$

where  $l_1, l_2, l_3 \in \mathbb{R}^3$  are primitive translation vectors that define the periodic translation symmetry of the crystal. Discrete linear transformations can be performed to obtain unit cells at different locations with  $\nabla = c_1 l_1 + c_2 l_2 + c_3 l_3$ , where  $c_1$ ,  $c_2$ , and  $c_3$  are integers, thus generating the entire 3-dimensional lattice. Therefore, a 3-dimensional lattice  $\Lambda$  is defined as all integral combinations of the lattice basis vectors.

$$\Lambda = \{c_1 \boldsymbol{l}_1 + c_2 \boldsymbol{l}_2 + c_3 \boldsymbol{l}_3 | c_i \in \mathbb{Z}\}.$$
 (2)

For a crystal with N atoms, where the atom positions are given by  $X = \{x_0, ..., x_{N-1}\}\$ , the corresponding position of atom uin a unit cell translated by  $c_1 \mathbf{l}_1 + c_2 \mathbf{l}_2 + c_3 \mathbf{l}_3$  is given by

$$\mathbf{x'}_{u} = \mathbf{x}_{u} + c_{1}\mathbf{l}_{1} + c_{2}\mathbf{l}_{2} + c_{3}\mathbf{l}_{3}.$$
 (3)

Further, there are 230 space groups in the 3-dimensional space, each of which describes a specific crystal symmetry. Every crystal in the database is associated with one space group number (1-230) depending on the arrangement of atoms in the crystal lattice. The order is based on the increasing complexity of symmetry elements and their combinations. For instance, space group number 1 is the simplest and least symmetric crystal system (triclinic), and 230 has the highest degree of symmetry (cubic).

#### 3.2 Crystal representation

A natural way to represent crystals is using graphs, with atoms as nodes and edges that connect neighboring or bonded atoms. However, using simple graphs is often not expressive enough to incorporate the inherent periodicity in crystals. In this work, we adopt multigraphs, following<sup>42</sup> to represent crystal structures. In multigraphs, two nodes can be connected by more than one type of edge. In the context of crystals, consider a graph  $\mathcal{G} = (V, E)$  with nodes (atoms)  $V = \{v_0, ..., v_{N-1}\}$  and edges (neighboring atoms)  $E = \{e_{uv,(c_1,c_2,c_3)} | 0 \le u \le N-1, 0 \le v \le N-1 \}$  $1, c_1, c_2, c_3 \in \mathbb{Z}, u, v \in V$ }. Here,  $e_{uv,(c_1,c_2,c_3)}$  is a directed edge from atom u to atom v in a unit cell translated by  $c_1 l_1 + c_2 l_2 + c_3 l_3$ . If  $c_1$  $=c_2=c_3=0$ , it corresponds to an edge between u and v in the same unit cell. Likewise, if  $c_1 = 1$ ,  $c_2 = c_3 = 0$  it corresponds to an edge between atom u in the original unit cell and atom v in a unit cell translated by  $l_1$ . This way, multigraphs carry information about the entire 3 dimensional structure of crystals.

#### 3.3 Offline reinforcement learning

While online RL methods demand frequent agent-environment interactions, offline RL exploits existing data,38 which is useful when receiving rewards or feedback from the environment is computationally expensive or physically implausible. As previously mentioned, our reward formulation depends on the energies and band gaps of crystals computed by DFT. Given that the time it takes to perform DFT simulation ranges between 6 seconds to more than 20 minutes for each input, depending on its size and type, it is highly infeasible to train an online reinforcement learning algorithm for this problem. Additionally, the high dimensional action space and the extremely complex reward landscape with narrow modes demand large amounts of exploration while learning in an online manner. Offline RL aims to learn from a static dataset  $\mathcal{D}$  consisting of state transitions, i.e.,  $(s_t, a_t, s_{t+1}, r_{t+1})$  obtained from a behavioral policy  $\pi_{\beta}(a|s)$  to learn an offline policy  $\pi_0(a|s)$ . However, directly adopting popular RL (e.g., deep O-learning) approaches in a data-driven manner causes two major issues - (1) the learned policy becomes out-of-distribution from the behavioral policy, and (2)

values of some states are over-estimated. Both these issues go hand-in-hand. Addressing these issues, 20 proposed conservative Q-learning (CQL), which regularizes Q-values by concurrently optimizing for the Bellman error to learn a conservative and lower-bound Q function. The optimization objective of the DQN<sup>43</sup> version (discrete action space) of CQL is given below

$$\min_{\theta} \omega \mathbb{E}_{s \sim \mathcal{D}} \left[ \log \sum_{\mathbf{a}'} \exp(Q_{\theta}(s, \mathbf{a}')) - \mathbb{E}_{s, \mathbf{a} \sim \mathcal{D}}[Q_{\theta}(s, \mathbf{a})] \right] + \frac{1}{2} \mathbb{E}_{s, \mathbf{a}, \mathbf{s}', r \sim \mathcal{D}} \left[ Q_{\theta}(s, \mathbf{a}) - \left( r + \gamma \max_{\mathbf{a}'} Q_{\theta'}(s', \mathbf{a}') \right) \right]^{2}.$$
(4)

Here,  $Q_{\theta}$  is the Q-network parametrized by  $\theta$ , and  $Q_{\theta'}$  is the target network.  $\omega$  controls the amount of conservatism, *i.e.*, higher the value of  $\omega$ , the more the preference for a conservative policy that better fits the data. When the action space is discrete, the learned discrete and deterministic offline policy is therefore

$$\pi_{o}(\boldsymbol{a}|\boldsymbol{s}) = \underset{\boldsymbol{a}}{\operatorname{argmax}} Q_{\theta}(\boldsymbol{s}, \boldsymbol{a}).$$
 (5)

#### 3.4 Density functional theory

DFT is a quantum mechanics-based simulation model that is used to compute the electronic structure of multi-atom systems, thereby estimating several properties including total energy, formation energy, and band gap. This is achieved by iteratively solving the Kohn-Sham equations.44 For evaluating crystal structures, we make use of the open-source Quantum Espresso software suite22 to perform self-consistent field (SCF calculations) using the Perdew-Burke-Ernzerhof (PBE) exchange-correlation functional. However, the PBE functional is known for its systematic underestimation of band gap energies,45 and is less accurate than functionals like HSE06 (ref. 46) or other self-energy approximations like GW.<sup>47</sup> Nevertheless, we used PBE because of its lower computational costs and superiority over DFTB. The output produced by the DFT simulation consists of two important properties that we are interested in - total energy (in Rydberg) and band gap (in eV units). Using total energy, we can also compute the formation energy (in eV per atom units). We also faced multiple new crystals failing to complete DFT simulation due to unknown properties (e.g., spin, magnetization) as part of our evaluation. Details about failure handling are provided in Section A.1.6.1.

#### 4 **Methods**

#### 4.1 **RL** formulation

The RL formulation of our problem follows a MDP defined as  $M = \langle \mathcal{S}, \mathcal{A}, \mathcal{T}, R, \gamma \rangle$ , where  $\mathcal{S}$  denotes the state space,  $\mathcal{A}$  denotes the action space,  $T(s'|s,s): S \times S \times A \rightarrow [0,1]$  is the environment transition probability function,  $R(\mathbf{s}, \mathbf{a}) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward function, and  $\gamma \in [0, 1]$  is a discount factor denoting the preference for long term rewards over short term rewards. In our setup, the state space consists of empty, partially or fully filled multigraphs  $(\mathcal{G}(V,E))$  of crystal structures. The action space A consists of atomic elements from which the agent has to choose to assign an atom at a given atomic site in a unit cell. Starting with initial state  $s_0$ , which is the graph  $\mathcal{G}_0$  of an empty

crystal skeleton, the sequential construction of a crystal of N atoms can be represented as a trajectory, as shown in Fig. 1a.

**4.1.1 Reward function.** For this property-driven crystal design problem, our reward function is expected to penalize high positive formation energies ( $E_{\rm form}$ ) and large deviations from a desired property of interest (e.g., band gap), whose value is denoted by  $\hat{p}$ . In the context of training an offline RL agent with batches of transitions, we aim to minimize the deviation between the ground truth property p of the crystal and  $\hat{p}$  (desired property). This bi-objective optimization can be addressed by using a linear combination of terms that individually optimize for lower energy and desired property. In other words, for a crystal with N atoms, the terminal reward, which is also equal to the return in this case, can be expressed in terms of its formation energy  $E_{\rm form}$  and ground truth property p as follows.

$$r_N(E_{\text{form}},\hat{p},p) = \alpha_1 g_E(E_{\text{form}}) + \alpha_2 g_p(p,\hat{p}). \tag{6}$$

Here,  $g_E(E_{\rm form})$  enforces lower formation energy,  $g_p(p,\hat{p})$  enforces p and  $\hat{p}$  to be close (e.g. distance function), and  $\alpha_1$  and  $\alpha_2$  are design parameters that control the importance of each of the terms. We use the exponentials of the negative formation energy of the crystal and the distance between the true and desired properties, yielding a terminal reward as follows:

$$r_N = \alpha_1 \exp\left(-\frac{E_{\text{form}}}{\beta_1}\right) + \alpha_2 \exp\left[-\frac{(p-\hat{p})^2}{\beta_2}\right].$$
 (7)

This introduces more design parameters  $\beta_1$ , and  $\beta_2$ , which essentially influence the sharpness of the mode of the exponential function; lower value of  $\{\beta_i\}_{i=1,2}$  results in a higher level of sharpness.<sup>2</sup>

#### 4.2 Q-Network and state representation

Our conditional Q-network  $Q_{\theta}(\mathbf{s}, \mathbf{a}; \hat{\mathbf{p}})$  consists of two components: (1) a graph neural network that extracts meaningful state representation of the input multigraph; (2) linear layers for computing Q-values from this representation. To represent and process multigraphs in an expressive manner, we adopt the MEGNet model,49 a universal graph machine learning framework for molecules and materials. MEGNet provides an effective way of iterative information exchange among node, edge and state features, which is particularly useful for chemical entities. For a crystal graph  $\mathcal{G}(V, E, \mathbf{v}; \hat{p})$  conditioned on the desired property  $\hat{p}$ , V and E are sets of nodes and edges, and  $\nu$  corresponds to the global state-level feature. For the N atoms in a unit cell, the categorical feature of the nodes  $H = \{h_u\}_{u=0}^{N-1}$  denote the one-hot encoding of the atom type in each of the nodes. It includes an additional dimension to indicate whether the node is currently filled or unfilled with an atom. Edges connect neighboring atoms based on the CrystalNN scheme proposed by of for determining the presence and type (i.e.,  $(c_1, c_2, c_3)$  triplet) of edges. The set of edge features  $\mathcal{T} = \{t_{uv,(c_1,c_2,c_3)}\}$ represents the Gaussian distance between the position of atom u in the reference unit cell and atom  $\nu$  in a unit cell shifted by  $c_1 l_1$  +  $c_2 l_2 + c_3 l_3$ .

$$t_{uv,(c_1,c_2,c_3)} = \exp\left[-\frac{d_{uv,(c_1,c_2,c_3)}^2}{\rho}\right],$$
 (8)

$$d_{uv,(c_1,c_2,c_3)} = \sqrt{(\mathbf{x}_v + c_1 \mathbf{l}_1 + c_2 \mathbf{l}_2 + c_3 \mathbf{l}_3 - \mathbf{x}_u)^2},$$
 (9)

where  $x_u, x_v \in \mathbb{R}^3$  are the positions (Cartesian coordinates) of atoms u and v in the reference unit cell. The state-level feature y is expressed as follows:

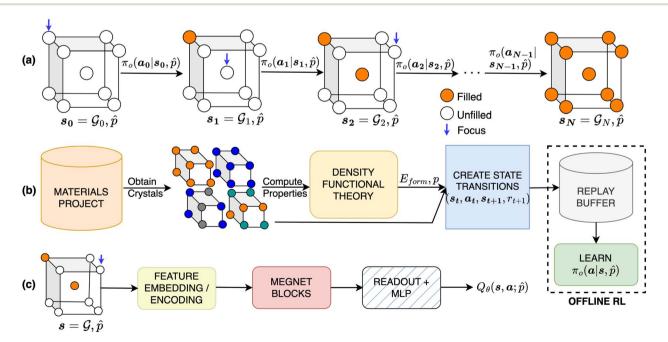


Fig. 1 (a) Our design approach centers on filling in the composition of predefined crystal using an RL policy. (b) To successfully train an RL policy, we obtain data from Materials Project,<sup>48</sup> recompute relevant property values using open-source DFT (Quantum Espresso<sup>22</sup>) and create trajectories for offline RL. (c) We train a graph neural network based policy based on MEGNet<sup>49</sup> to achieve property-conditioned crystal generation.

$$y = [z||f|, z = [a, b, c, \phi_1, \phi_2, \phi_3, S, \hat{p}],$$
 (10)

where, a, b, c are the lengths of the edges of the lattice  $(a = ||l_1||,$  $b = ||\boldsymbol{l}_2||, c = ||\boldsymbol{l}_3||, \phi_1, \phi_2, \phi_3$  are the angles of the lattice, S is the space group number of the crystal,  $\hat{p}$  is the desired property that the policy is conditioned on, and f is a categorical feature, which we refer to as focus - it instructs the policy which unfilled node to focus on for atom type prediction in the following step. The categorical features H and f are passed through embedding layers to obtain embedded feature maps  $\tilde{H}, \tilde{f}$ . Numerical features T and y are passed through multilayer perceptrons (MLPs)

$$\tilde{y} = \text{MLP}([z||\tilde{f}]). \tag{11}$$

A graph  $\tilde{\mathcal{G}}$  with embedded and encoded features is then passed through K MEGNet layers, followed by a readout layer (Appendix A.1.4) to obtain a graph-level representation, which is then passed through an MLP to obtain conditioned Q-values for all actions in A.

$$\tilde{\mathcal{G}}^{(k+1)} = \text{MEGNET}(\tilde{\mathcal{G}}^{(k)}) \forall k = 0, ..., K-1$$
 (12)

$$\psi(\tilde{\mathcal{G}}^{(K)}) = \text{READOUT}(\tilde{\mathcal{G}}^{(K)})$$
 (13)

$$Q_{\theta}(s = \mathcal{G}; \hat{p}) = \text{MLP}\Big(\psi\Big(\tilde{\mathcal{G}}^{(K)}\Big)\Big)$$
 (14)

#### 4.3 Dataset

For this study, we used a subset of the Materials Project database, referred to as MP-20, that was previously used by.13 MP-20 consists of ~45k metallic and nonmetallic crystals with different structures and compositions, covering 88 elements in the periodic table. All of them have at most 20 atoms. For our experiments, we excluded metallic crystals with zero band gap. # Metals constituted more than 60% of the data, leading to class imbalance challenges while conditioning the model with a nonzero band gap. Next, we used Quantum Espresso to determine the formation energies and band gaps of all nonmetallic crystals in the training and validation set. In the end, our training set included 8832 crystals, and our validation set included 2486 crystals.

#### 4.4 State transitions for offline RL

As shown in Fig. 1, we generated a static dataset for training the offline policy using episodic trajectories consisting of  $(s_t, a_t, s_{t+1}, s_{t+$  $r_{t+1}$ ) transitions from MP-20 crystals. We applied a deterministic policy  $\pi_{\beta}(a|s)$ , where the actions correspond to the original element identities of the atom at a specific position of interest in an empty or partially constructed crystal skeleton graph. Each trajectory of an episode starts with the initial state  $s_0$ , which is a graph  $\mathcal{G}_0$  of a crystal skeleton, where all atom identities are hidden. Through the focus feature f, we are explicitly providing the order of traversal through the nodes of the graph, thereby simplifying the problem further. To mitigate the effects of bias

due to this order dependency, we obtain up to 5 trajectories for each crystal by varying the order of nodes with breadth-first traversals of the graph from different source nodes. This way, we obtained  $\sim$ 520k transitions to train our offline RL policies.

#### 5 **Experiments**

In this study, we focus on training conditional CQL (c-CQL) models to design stable (i.e., low formation energy) crystals that have a desired band gap  $(\hat{p})$  of 1.12 eV, 2 eV, 3 eV, and 4 eV, which fall within the semiconductor range. To determine the amount of conservatism required for better performance, we varied  $\omega$  using weights of 1 and 5, with the latter being more conservative than the former. After an initial hyperparameter sweep, we choose the coefficients as follows:  $\alpha_1 = 1$ ,  $\alpha_2 = 10$ ,  $\beta_1$ = 5,  $\beta_2$  = 3. Our baselines are (1) random policy, (2) behavioral cloning (BC)§, and (3) Unconditional CQL (u-CQL) Policy (where  $\hat{p}$  is removed in the state feature vector and the reward is only in terms of  $E_{\text{form}}$ ). For evaluating the model, we start with an empty crystal skeleton graph  $\mathcal{G}_0$  as the initial state  $s_0$ , and perform a rollout using the learned conditional offline policy  $\pi_{o}(\boldsymbol{a}|\boldsymbol{s},\hat{p})$  to sequentially fill atoms in the crystal. We then perform a presimulation assessment of the generated crystals using the following metrics - (1) compositional validity: a generated crystal is valid if it has an overall neutral charge, as computed by SMACT,51 (2) accuracy, which is the fraction of correctly predicted atoms, (3) similarity, which measures the similarity of the predicted atoms with the ground truth, i.e., two atoms are similar if they belong to the same class of elements, and (4) novelty, which measures the fraction of valid crystals whose composition is not present already in the Materials Project database. Our results are shown in Table 1.

Next, we performed DFT simulation for all the valid crystals to estimate the total energy and band gap. The post-simulation metrics are (1) Average Formation Energy per atom  $\bar{E}_{form}$  of the policy-generated crystals (2) Earth Mover Distance (EMD) between the generated and true band gap distributions ( $\Gamma_{\text{true}}^p$ ), (3) Earth Mover Distance between the generated and true formation energy distributions ( $\Gamma_{\text{true}}^{E}$ ), (4) % of crystals that have the band gap value in the desired range  $(\nu)$ , which in our case is from  $\hat{p} - 0.25$  eV to  $\hat{p} +$ 0.25 eV, and (5) Out-of-distribution (OOD) design ( $\kappa$ ) – % of generated crystals that have band gaps in the desired range but whose corresponding ground truth crystals have band gaps outside the desired range. The results are shown in Fig. 2. Our initial set of experiments incorporated total energy  $(E_{tot})$  in the reward formulation instead of formation energy, where we also tuned the design parameters. The results of the same are detailed in Appendix A.3.

#### 5.1 Analysis of pre-simulation metrics

For all band gap targets, as seen in Table 1, the more conservative models (i.e.,  $\omega = 5$ ) generally perform better in terms accuracy,

<sup>‡</sup> Metallic crystals, being conductors, have a zero band gap because of the overlapping conduction and valence bands.

<sup>§</sup> Trained with supervised classification loss.

<sup>¶</sup> Classes - transition metals, post-transition metals, group 1 metals, group 2 metals, nonmetals, lanthanides, actinides, halogens, and noble elements (Appendix A.1.1.2).

Table 1 Pre-simulation metrics for all four band gap targets compared with random, BC and unconditioned policy baselines. More conservative models generally perform well in accuracy, similarity, and validity, but generate less novel crystals. BC outperforms CQL models in accuracy and validity

	Accuracy (%	%)	Similarity (%)		Validity (%)		Novelty (%)	
CQL weight	$\omega = 1$	$\omega = 5$	$\omega=1$	$\omega = 5$	$\omega=1$	$\omega = 5$	$\omega=1$	$\omega = 5$
Random	0.0115		0.1254		NaN		NaN	
BC	49.37		69.04		81.84		51.34	
u-CQL	46.30	48.69	67.73	68.64	79.17	80.78	51.15	48.26
c-CQL ( $\hat{p} = 1.12 \text{ eV}$ )	43.31	47.72	65.66	69.24	78.75	80.66	63.51	50.95
$c$ -CQL ( $\hat{p} = 2$ eV)	42.73	47.97	65.35	68.87	80.23	79.99	65.57	51.18
c-CQL ( $\hat{p} = 3 \text{ eV}$ )	43.59	48.16	65.95	69.67	79.15	81.15	65.08	50.11
$\text{c-CQL }(\hat{p} = 4 \text{ eV})$	43.40	46.63	65.92	68.29	79.87	78.29	65.55	51.03

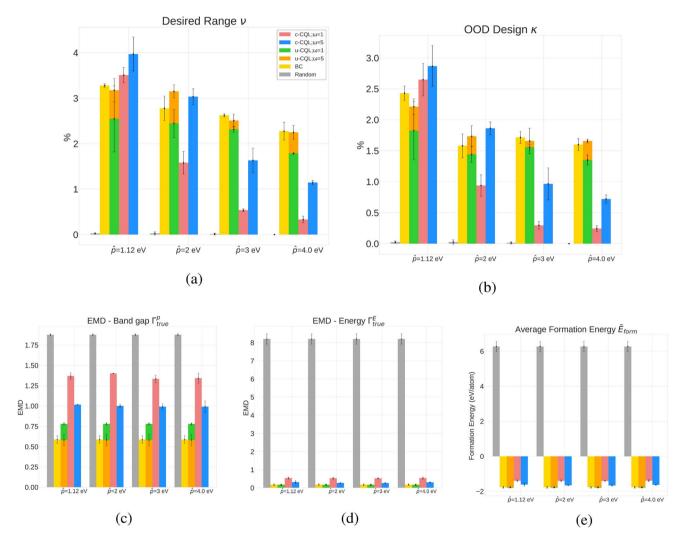


Fig. 2 Results for conditioned CQL policies on four band gap design targets (1.12 eV, 2 eV, 3 eV and 4 eV) with formation energy in the reward function (eqn (7)). Conditioned and more conservative policies perform well in the  $\kappa$  and  $\nu$  metric when the target is lower, while unconditioned policies, including behavioral cloning, perform better at reproducing the original distribution. Random policies fail to reproduce the original distribution and achieve desired properties. (a) % Desired range for the four band gaps targets for various policies. Conditioned policies outperform random policy and compete with unconditional policies in designing crystal in the desired property range ( $\hat{p} - 0.25$ ,  $\hat{p} + 0.25$ ). (b) % of generated crystals with property in the desired range with corresponding ground truth crystals outside the desired range. Conditioned policies outperform baselines for lower band gap targets (1.12 eV and 2 eV). (c) Band gap EMD (generated  $\nu$ s. true) for various policies showing that unconditioned policies reproduce the original dataset better. Lower value indicates more resemblance to the true distribution. (d) Formation energy EMD (generated  $\nu$ s. true) for various policies showing that unconditioned policies reproduce the original dataset better. Lower value indicates more resemblance to the true distribution. (e) Average formation energy for various policies yielding valid crystals with energy below 0. The average formation energy of randomly generated crystals is high and positive. Lower is better.

similarity, and validity. The metrics were also influenced by the magnitude of the reward function - the higher the magnitude, lower the accuracy, and in most cases, the lower the validity of generated structures (pre-simulation results with total energy in Tables 2 and 3). This is interesting because when the magnitude of the reward is lower or  $\omega$  is higher, the conservative term in the CQL objective in eqn (4) becomes dominant, resulting in the net maximization of Q-values of state-action pairs present in the dataset. Higher  $\omega$  also results in lower novelty scores. Behavioral cloning (BC), trained with no reward signal, performed the best in accuracy and validity (Table 1), which can be attributed to BC's better prediction capacities attributed to supervised learning. However, this might not be helpful from the perspective of property-driven crystal design where the CQL-based policies outperform BC in  $\nu$  and  $\kappa$  in some cases, as described next in Section 5.2 outlining relevant case studies.

#### 5.2 Band gap design case studies: targeting 1.12 eV, 2 eV, 3 eV & 4 eV

The results in Fig. 2, which include a well-performing policy for all the design cases, show some clear trends: (1) for lower band gap targets (i.e. 1.12 eV, 2 eV), conditioned policies (with  $\omega = 5$ ) generate more materials in the desired property range when the corresponding true materials have properties outside the desired range (Fig. 2a). Examples of such materials are shown in Fig. 3. (2) Greater conservatism leads to more materials in the desired range as shown by the fact that  $\omega = 5$  outperforms  $\omega = 1$  in all design cases. (3) Unconditioned policies manage to recreate the original distributions better than conditioned distributions. This is shown by better performance in Fig. 2c and d, holding for both band gap and formation energy. (4) Random policies are not effective in generating valid and desired crystal structures. It is likely that the

	Exa	mple 1	Example 2		
$\hat{p}$	True Crystal	Generated Crystal	True Crystal	Generated Crystal	
1.12 eV	$p=2.157 \ Rb_2TlInBr_6$	$p = 1.251$ $TlGaTe_2$	p = 3.280 Cs <sub>2</sub> KYI <sub>6</sub>	$p=1.082$ $Cs_2RbInI_6$	
2 eV	$p = 1.049$ $TlGaTe_2$	$p=2.113 \ KAlTe_2$	$p = 5.370$ $Mg_2P_2O_7$	$p=1.943 \ MgSnP_2O_7$	
3 eV	$p = 3.959$ $Cs_2Si(HO_2)_2$	$p = 3.243$ $Rb_2Si(HO_2)_2$	$p=1.150$ $Cs_2CuBiBr_6$	$p=2.977$ $K_2NaBiBr_6$	
4 eV	$p=5.097$ $Na_3InF_6$	$p=4.221 \ SrCaLa_4Zr_2O_{12}$	$p = 0.801$ $TlIn_2GaTe_4$	$p=3.979 \ KNaCl_2$	

Fig. 3 Examples of cases where the crystal generated by our model has the band gap in the desired range, i.e.,  $(\hat{p} - 0.25, \hat{p} + 0.25)$ , while the ground truth crystal has the band gap outside the desired range. In most cases, it can be observed that some of the elements are common in the true and generated crystals. This indicates selective atomic substitutions for favorable band gap shifts

random policy generated a small subset of valid metal-like crystals given the close to zero average band gap shown in Appendix A.1.3. Random policy generates many unrealistic crystals with high and positive formation energies (Fig. 2e), and many of the DFT runs validating the crystals failed (Table 2). Results of pre-simulation metrics included in Table 1 also show poor performance of random policy. As shown in Fig. 2, the higher values of  $\hat{p}$  are more challenging because: (1) most samples in the dataset have a lower band gap value (Appendix A.2) making the number of samples with a higher band gap that get exposed to the model while training a very small fraction, (2) underestimation of band gaps by DFT, which causes an unfavorable shift from the expected band gap distribution.

#### 6 Limitations

The important limitations of this work are that the scope is limited to predicting only the atom types, given all other information about the skeleton of the crystal and the order of traversal, and the training data is small and limited to nonmetals. Considering computational challenges attributed to DFT calculations, we had to restrict our design parameter space to a very small set, but it would be interesting to see the results after an extensive analysis after training models with several values of  $\omega$ ,  $\alpha_1$ ,  $\alpha_2$ ,  $\beta_1$  and  $\beta_2$ . A recent study showed that the performance of offline RL algorithms is influenced by the bias of the dataset generated by the behavioral policy and the strength of the reward signals.52 These aspects should be analyzed in the context of the crystal design problem with DFT-based reward signals for choosing the most appropriate offline RL algorithm and design parameters. Due to the significant underestimation of band gaps by DFT, many of the generated crystals had an estimated band gap value of 0.0, which severely hindered our evaluation and analyses. This explains the very low fraction of generated crystals having a greater band gap. Further, resolving the imbalance in the data due to the large number of samples in the lower band gap regions could help in learning better policies for generating crystals with higher band gap.

#### 7 Conclusion

We show that it is possible to train reinforcement learning based policies that can design valid crystal compositions conditioned on a crystal structure skeleton and a target property, such as the band gap, evaluated on precise and expensive computational chemistry engines, such as DFT. We demonstrate that offline RL methods can be used to learn distributions of design trajectories for valid crystal structures and provide tuning based on desired properties. While our results suggest that one can train policies for materials design problems, there is still significant space for future work to improve the performance, robustness, and capabilities of the RL policies. First, our approach can be extended to include additional design variables, such as crystal lattice parameters and atomic positions, for greater design flexibility to design more performant materials. Second, the dataset we used for offline RL is still limited in size given the large cost of generating the dataset in a consistent manner and evaluating the reward function for structures generated by the policy. This leaves significant room for future work in creating large pretraining datasets and accelerating the evaluation of crystal structures through more optimized high-throughput DFT or machine learning based approximators. Third, much algorithmic work remains in designing better policies for material design that can further improve the performance of conditional design.

# Data availability

The source code and models for this work can be obtained from <a href="https://github.com/chandar-lab/crystal-design">https://github.com/chandar-lab/crystal-design</a>. The offline trajectory datasets can be obtained from <a href="https://zenodo.org/records/10626005">https://zenodo.org/records/10626005</a>.

# Conflicts of interest

The authors declare no conflicts of interests.

# A Appendices

#### A.1 Experimental details

- **A.1.1 Pre-simulation metrics.** Pre-simulation metrics were computed for crystals designed by the policies prior to performing simulation using DFT they are (1) accuracy, (2) similarity (3) compositional Validity (referred to as validity for simplicity), and (4) novelty. Further details on how to calculate them are provided below.
- A.1.1.1 Accuracy. Accuracy is measured as the percentage of predicted atoms that match the ground truth. Note that the accuracy in this case is computed globally across atoms predicted in all the crystals present in the validation dataset.

Accuracy(%) =

#predicted atoms that exactly match the ground truth
Total number of predicted atoms

We can also measure the fraction of crystals that were reconstructed to match the ground truth exactly. However, this was a very small percentage ( $\sim$ 2–7%) for all the models.

- A.1.1.2 Similarity. While accuracy measures the fraction of exact matches, our similarity metric considers a prediction as a match if the predicted atom and the ground truth atom belong to the same category. The categories are defined as follows.
  - (1) Group 1: Li, Na, K, Rb, Cs
  - (2) Group 2: Be, Mg, Ca, Sr, Ba
- (3) Transition metals: Sc, Ti, V, Cr, Mn, Fe, Co, Ni, Cu, Zn, Y, Zr, Nb, Mo, Tc, Ru, Rh, Pd, Ag, Cd, Hf, Ta, W, Re, Os, Ir, Pt, Au, Hg
  - (4) Nonmetals: H,B, C, N, O, Si, P, S, As, Se, Te
  - (5) Halogens: F, Br, Cl, I
  - (6) Noble: Xe, Ne, Kr, He
- (7) Lanthanides: La, Ce, Pr, Nd, Pm, Sm, Eu, Gd, Tb, Dy, Ho, Er, Tm, Yb, Lu
  - (8) Actinides: Ac, Th, Pa, U, Np, Pu

Similarity(%) =

#predicted atoms that have same category as ground truth
Total number of predicted atoms

A.1.1.3 Compositional validity. We follow<sup>13</sup> and compute the compositional validity of crystals using SMACT.<sup>51</sup>

$$Validity(\%) = \frac{\#valid\ crystals}{Total\ number\ of\ crystals}$$

A.1.1.4 Novelty. To assess the novelty aspect of our approach, we compute the fraction of valid generated crystals whose compositions are novel, i.e., when the compositions are not present in the Materials Project Database. We utilised the API of Materials Project (mpr.summary.search function) to retrieve crystals with matching compositions. Note that our novelty percentage is conditioned on the valid crystals, and we do not query invalid compositions. Hence, in Table 1, while other metrics are computed by dividing the total number of crystals in the validation set, novelty is computed by dividing the number of valid crystals generated by the model.

$$Novelty(\%) = \frac{\#crystals \ with \ novel \ compositions}{\#valid \ crystals}$$

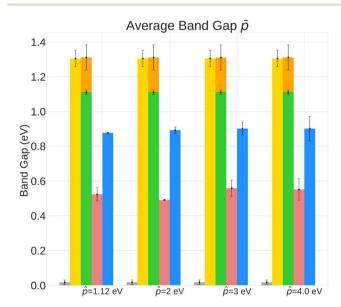
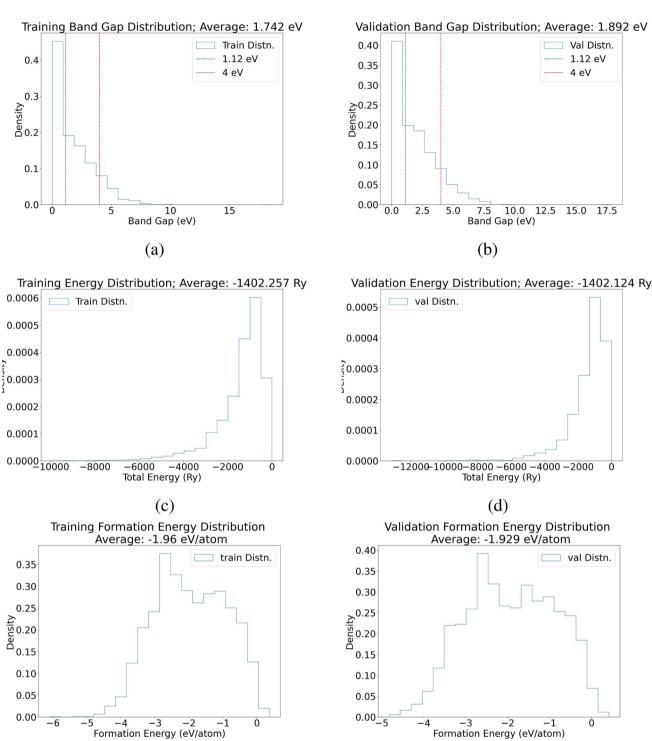


Fig. 4 Analysis of average band gap of generated crystals in the validation set. BC and unconditional policies have an average band gap closer to the ground truth average (1.892 eV). Random policy failes to generate crystals with higher band gap.

Table 2 % Generated valid crystals that successfully underwent DFT simulation, for random policy and each of the trained models. Most of the crystals generated by the random policy failed DFT simulation

	% DFT success	
CQL weight	$\omega=1$	$\omega = 5$
Random	15.18	
BC	68.25	
u-CQL	68.97	70.29
c-CQL ( $\hat{p} = 1.12 \text{ eV}$ )	58.59	66.82
c-CQL ( $\hat{p} = 2 \text{ eV}$ )	56.04	67.99
$c$ -CQL ( $\hat{p} = 3 \text{ eV}$ )	56 <b>.</b> 55	68.38
c-CQL $(\hat{p} = 4 \text{ eV})$	55.64	66.19



Distribution of band gap (eV), total energy (Ry), and formation energy (eV per atom) for training and validation datasets.

A.1.2 Post-simulation metrics. Post-simulation metrics were computed for crystals designed by the policies after performing simulation using DFT. As indicated in Appendix A.1.6.1, crystals that failed DFT simulation were not included while computing post-simulation metrics. Details on how to calculate them are provided below.

(e)

A.1.2.1 Average formation energy. Following Appendix A.1.7, the formation energy was calculated for all the generated and valid crystals that successfully underwent DFT simulation. The average formation energy if therefore,

(f)

$$\overline{E}_{\text{form}} = \sum_{i=1}^{N} \frac{E_{\text{form},i}}{N} (\text{eV per atom})$$

**Table 3** Pre-simulation metrics for band gap design case of 1.12 eV with  $(\alpha_1 - \alpha_2 - \beta)$  corresponding to the terms of the reward function in eqn (16) and **best by metric highlighted**. Unconditional policies perform better on pre-simulation metrics while conditioned policies produce target designs shown as in Fig. 7 and discussed in Section 5.2

	Accuracy (%)		Similarity (%)		Validity (%)	
CQL weight	$\omega = 1$	$\omega = 5$	$\omega = 1$	$\omega = 5$	$\omega = 1$	$\omega = 5$
Random	0.0115	w – 3	0.1254	w – 3	NaN	w – 3
BC	52.26		71.98		85.00	
uCQL	49.77	51.53	70.85	71.26	81.50	82.54
(0-5-1)	38.64	48.85	61.23	69.38	69.99	77.84
(0-5-3)	43.02	46.43	65.01	67.04	73.57	78.44
(0-10-1)	36.54	43.72	59.3	65.18	73.33	80.81
(0-10-3)	35.16	42.42	57.48	64.15	71.20	81.30
(1-5-1)	42.11	47.72	64.00	68.12	75.62	80.29
(1-5-3)	40.59	47.57	63.70	67.26	72.93	76.51
(1 - 10 - 1)	35.02	43.18	58.63	65.13	67.82	75.14
(1 - 10 - 3)	35.38	43.81	57.23	65.58	61.87	77.19

where *N* is the number of valid crystals whose formation energy values were computed successfully using DFT.

A.1.2.2 EMD (band gap). The Earth Mover's Distance (EMD) was computed to determine the distributional distance between the properties of generated crystals and the ground truth crystals in the validation dataset. For band gap, the  $\Gamma^p_{\rm true}$  was calculated as follows.

$$\Gamma_{\text{true}}^p = \text{EMD}(\{p_i\}_{i=1}^M, \{\tilde{p}_i\}_{i=1}^N)$$

where M is the total number of crystals in the validation set, and N is the number of valid generated crystals that successfully underwent DFT simulation.  $p_i$  is the property value of the ith crystal in the validation set, and  $\tilde{p}_j$  is the property value of the jth valid crystal generated by the model.

A.1.2.3 EMD (formation energy). Similar to  $\Gamma_{\text{true}}^p$ , EMD between the true and generated formation energy distributions,  $\Gamma_{\text{true}}^E$  were computed as follows.

$$\Gamma_{\text{true}}^E = \text{EMD}(\{E_{\text{form},i}\}_{i=1}^M, \{\tilde{E}_{\text{form},i}\}_{i=1}^N)$$

 $E_{{
m form},i}$  is the property formation energy (ev/atom) of the *i*th crystal in the validation set, and  $\tilde{E}_{{
m form},j}$  is the property value of the *j*th valid crystal generated by the model.

A.1.2.4 Desired range. The desired range metric ( $\nu$ ) is the fraction of generated crystals whose property (here, band gap) lies between  $\hat{p}-0.25$  and  $\hat{p}+0.25$ , where  $\hat{p}$  is the target property. For simplicity and easier analysis, the denominator of this fraction is the total number of crystals in the validation set. This way, the metric provides a way to quantitatively compare the corresponding percentages across different models.

$$\nu = \frac{\text{#generated crystals in the property range}(\hat{p} - 0.25, \hat{p} + 0.25)}{M}$$

Here, *M* is the number of crystals in the validation set.

A.1.2.5 OOD design. Through the  $\kappa$  metric, we compared the number of crystals generated (from the validation set) whose property value lies in the desired range, i.e.,  $(\hat{p}-0.25,\hat{p}+0.25)$ , but the corresponding ground truth property is outside the desired

range (hence, OOD crystals). This indicates that the model has learned to place atoms such that the property shifts from a value outside the desired range to within the range. Similar to  $\nu$ , the denominator is M, the number of crystals in the validation set.

$$\kappa = \frac{\text{\#OOD crystals}}{M}$$

**A.1.3** Additional post-simulation results. As part of the post-simulation analysis, we also investigated the average band gap of the crystals designed by each model (Fig. 4).

**A.1.4 MEGNet.** In our work, we adopted the MEGNet<sup>29</sup> model to process crystal graphs and extract state representation, as part of the Q-network  $Q_{\theta}$ . The important hyperparameters of the model are listed below.

- Number of MEGNet blocks: 3
- Node embedding dimensions: 16
- Edge embedding dimensions: 1
- State embedding dimensions: 8
- READOUT Function: order-invariant set2set53

**A.1.5 Offline RL.** We adopt conservating Conservative Q-Learning (CQL)<sup>20</sup> as the offline RL approach. The important hyperparameters of our training process is listed below.

- Number of steps trained: 250 000
- Discount factor: 0.99
- Batch size: 1024
- Learning rate:  $3 \times 10^{-4}$
- Soft target network update rate:  $5 \times 10^{-3}$
- Optimizer: Adam

**A.1.6 DFT parameters (Quantum Esperesso).** For performing DFT calculations, we use the Quantum Espresso v7.1 (ref. 22) simulation suite. The details of the DFT parameters are given below. For simplicity, this configuration was used for all crystals, and the evaluation is consistent for the training and generated crystals. Note that we do not perform structure relaxation in any of the cases.

- Calculation: SCF
- Pseudopotentials: solid-state pseudopotentials (SSSP) version 1.3.0 obtained from https://www.materialscloud.org/ discover/sssp/table/efficiency
  - Tolerance: 10<sup>-6</sup>
  - Number of bands: 256
  - *k*-points: (3–3–3)
- Occupations: fixed (since our training set consists only of nonmetallic crystals)
  - Diagonalization: David
  - ecutrho: 245
  - ecutwfc: 30
  - mixing\_beta: 0.7
  - degauss: 0.001
  - Default charge: 0
  - Maximum iterations: 1000

A.1.6.1 Handling failures. It is important to note that DFT can be best leveraged once we know certain properties of the crystals – for example, charge, magnetization, and metallicity. Considering the difficulty in determining these properties for completely unknown crystals, we standardized the evaluation

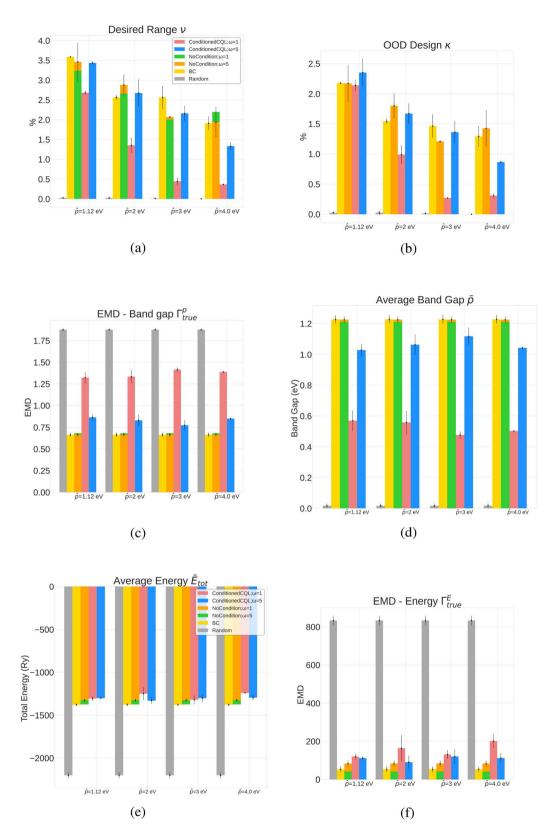


Fig. 6 Results for conditioned CQL policies on all band gap design targets. Conditioned and more conservative policies perform better in the  $\kappa$ metric in some cases, while unconditioned policies, including behavioral cloning, perform better at reproducing the original distribution. Random policies fail to reproduce the original distribution and achieve desired properties.

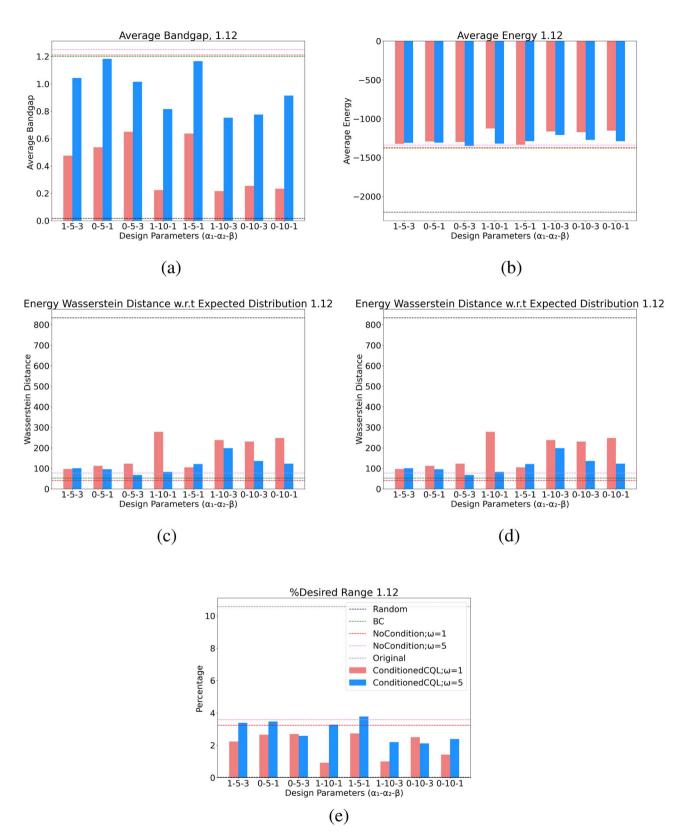


Fig. 7 Full design parameter values for all learned policies for the band gap design case of 1.12 eV. Nomenclature of the table is  $(\alpha_1 - \alpha_2 - \beta)$  corresponding to the terms of the reward function in eqn (16).

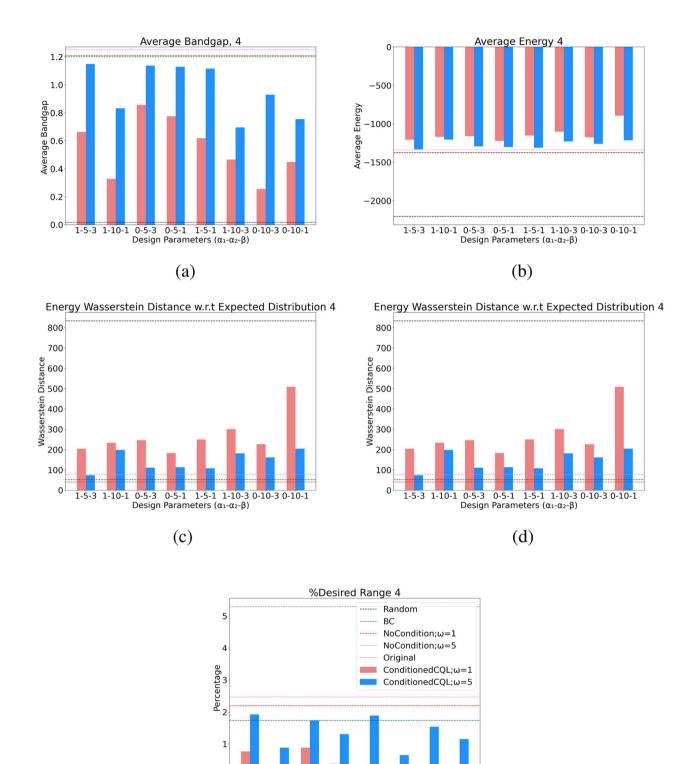


Fig. 8 Full design parameter values for all learned policies for the band gap design case of 4.0 eV. Nomenclature of the table is  $(\alpha_1 - \alpha_2 - \beta)$ corresponding to the terms of the reward function in eqn (16).

Design Parameters  $(\alpha_1-\alpha_2-\beta)$ 

(e)

0-5-1 1-5-1 1-10-3 0-10-3 0-10-1

1-5-3 1-10-1

0-5-3

procedure by using the same DFT configuration for all crystals (except for the crystal-specific parameters like number of atoms, species, and pseudopotentials directory). However, this resulted in multiple crystals failing DFT simulation. Some of the errors are explained below.

- Charge is wrong. Smearing is needed: this error mainly occurs because of unpaired electrons in the system, and can be resolved by changing the occupation to 'smearing' instead of 'fixed'. However, doing so will not help in determining the band gap of crystals, as it will only output the Fermi energy. Another way is to set the 'nspin' parameter to 2 and specify the total magnetization value as an additional input to Quantum Espresso. This helped us resolve most of the failures for the MP-20 crystals in the training and validation set because the total magnetization value is retrievable from the Materials Project, but for the newly generated crystals, we had to ignore those that failed because of this error. The error could also occur if the generated crystal is metallic, and this property is also difficult to identify directly from the structure and composition.
- NOT converged in 1000 iterations: for some crystals, the DFT simulation did not converge even after 1000 iterations. These crystals were ignored while constructing the offline dataset, and also when evaluating the policy-generated crystals.
- Time limit exceeded: for constructing the offline dataset using known crystals, we used a flexible time limit to ensure none of the crystals were discarded because of time restrictions. However, while performing DFT simulation for the policy-generated crystals, due to the high-throughput nature of our evaluation pipeline, we had to ignore crystals that did not converge in 15 minutes.
- Too few bands: this error occurs when the number of bands specified, through 'nbnds' parameter is insufficient for the crystal system being simulated. This error was largely resolved by specifying a higher number of bands. In our case, we used 256 bands for all crystals.

Overall, during evaluation of generated crystals, only 50–70% of the valid crystals successfully underwent DFT simulation to output the energy and band gap (Table 2), and the rest failed because of the above errors.

A.1.6.2 % DFT success. Table 2 shows the percentage of policy-generated crystals that successfully underwent DFT simulation based on failure handling strategies discussed in Appendix A.1.6.1.

 $\mbox{ \begin{tabular}{ll} Table 4 & Band gap design case of 4 eV with similar nomenclature and conclusions as Table 1 \\ \end{tabular} }$ 

	Accuracy (%)		Similarity (%)		Validity (%)	
CQL weight	$\omega=1$	$\omega = 5$	$\omega=1$	$\omega = 5$	$\omega=1$	$\omega = 5$
Random	0.0115		0.1254		NaN	
BC	52.26		71.98		85.00	
uCQL	49.77	51.53	70.85	71.26	81.50	82.54
(0-5-1)	41.82	48.09	64.34	68.82	80.21	82.18
(0-5-3)	39.46	47.61	61.59	68.24	74.46	80.09
(0 - 10 - 1)	33.24	39.42	60.78	53.42	62.39	67.82
(0-10-3)	35.24	41.47	57.14	64.06	64.40	75.54
(1-5-1)	38.80	46.79	60.09	68.77	70.80	80.17
(1-5-3)	42.06	47.49	63.36	68.35	78.32	81.0
(1 - 10 - 1)	36.52	42.21	59.57	65.07	76.55	74.41
(1-10-3)	35.94	42.91	56.8	64.2	68.95	77.63

**A.1.7 Formation energy calculation.** The formation energy per atom was calculated using the total energies of the crystals and their constituent elements. The total energies of the isolated elements (88 in the action space) were calculated by performing SCF calculations on the most stable elemental crystals (*i.e.*, 0 formation energy) present in the Materials Project. For elements that do not have a stable elemental crystal (*e.g.* Lu) or those that have large number of atoms in the elemental crystal (*e.g.* P, Se), the total energies were calculated for a single atom inside a primary cubic cell of length 10 Å. For a crystal with *N* atoms, the formation energy (per atom) calculation is defined as follows.

$$E_{\text{form}} = \left(\frac{E_{\text{tot}} - \sum_{i} \frac{N_{i}}{n_{i}} E_{\text{tot}}^{i}}{N}\right) \times 13.6057039763 \text{ (eV atom)}$$
 (15)

Here,  $N_i$  is the number of atoms of the constituent element i present in the crystal,  $n_i$  is the number of atoms (sites) of i in the elemental crystal, and  $E_{\text{tot}}^i$  is the total energy of i in the most stable elemental crystal form. 13.6057039763 is the value of 1 Rydberg constant in eV.

#### A.1.8 Algorithm. A.2 True distributions of properties

This section shows the true distribution of the band gaps and total energies for both training and validation data (Fig. 5).

# Algorithm 1 Training Conditional CQL: DQN Version for Crystal Design with Target Property $\hat{p}$ Construct dataset $\mathcal{D}$ of size $N_{\mathcal{D}}$ consisting of transitions (s, a, s', r) using known crystals Load $\mathcal{D}$ in Replay Buffer $\mathcal{B}$ Initialize Q-network $Q_{\theta}$ and target network $Q_{\theta'}$ , batch size B for j=1 to max\_steps $\mathbf{do}$ Sample B transitions, $\{(s_i, a_i, s'_i, r_i)\}_{i=1}^B$ from $\mathcal{B}$ Compute TD loss $L_i^{TD}(\theta) = \begin{cases} (Q_{\theta}(s_i, a_i; \hat{p}) - (r_i + \gamma \max_{a} Q_{\theta'}(s'_i, a; \hat{p})))^2 & \text{if } s'_i \text{ is not terminal otherwise} \\ (Q_{\theta}(s_i, a_i; \hat{p}) - r_i)^2 & \text{otherwise} \end{cases}$ $L^{TD}(\theta) = \frac{1}{B} \sum_{i=1}^B L_i^{TD}(\theta)$ Compute conservative loss, $L^C(\theta) = \frac{1}{B} \sum_{i=1}^B [\log \sum_{a} \exp(Q_{\theta}(s_i, a; \hat{p})) - Q_{\theta}(s_i, a_i; \hat{p})]$ Compute total CQL loss $L^{CQL}(\theta) = \omega L^C(\theta) + \frac{1}{2}L^{TD}(\theta)$ Compute gradients and backpropogate: $\theta \leftarrow \theta - \eta \nabla L^{CQL}(\theta)$ , $\eta$ is the learning rate Update target network parameters $\theta'$ end for

#### A.3 Experiments with total energy

As part of our initial analysis, we performed the experiments with total energy  $(E_{tot})$  in the reward formulation instead of formation energy, with the aim of designing crystals that are generally considered stable (in an absolute sense), so they can be used for practical purposes. However, total energy is less meaningful when it comes to comparing the stability of different crystals, while energy above hull is the best-known metric to compare thermodynamic stability.

A.3.1 Reward formulation. Since the units of total energy are in Rydberg (Ry), our reward function in eqn (7) can be redefined as follows.

$$r_N = \alpha_1 \log_{10}(-E_{\text{tot}}) + \alpha_2 \exp\left[-\frac{(p-\hat{p})^2}{\beta}\right]. \tag{16}$$

A.3.2 Full experimental metrics with total energy. We provide full experimental for our reward function design parameters for both the 1.12 eV design case (Table 3 and Fig. 7) and 4 eV case (Table 4 and Fig. 8) below. The tables and figures include evaluation of both the pre-simulation (except Novelty) and post-simulation metrics (except  $\kappa$ ) described in Section 5. With  $\alpha_1 = 1$ ,  $\alpha_2 = 5$ ,  $\beta = 1$ , the post-simulation results for all four band gap targets are shown in Fig. 6. All models in highlighted in this section were trained for 500 000 steps.

# Note added after first publication

This article replaces the version published on 22 March 2024. The caption for Fig. 2 contains additional details regarding parts (a)-(e).

# Acknowledgements

We would like to thank Pierre-Paul De Breuck and Rajesh Raju for their valuable domain-related inputs. We acknowledge and thank the people responsible for the smooth functioning of the compute resources of Mila and Compute Canada. We also wish to acknowledge funding from Intel and CIFAR. Janarthanan Rajendran is supported by IVADO postdoctoral fellowship. Sarath Chandar is supported by the Canada CIFAR AI Chairs program, the Canada Research Chair in Lifelong Machine Learning, and the NSERC Discovery Grant.

#### References

- 1 H. Wang, T. Fu, Y. Du, W. Gao, K. Huang, Z. Liu, et al., Scientific discovery in the age of artificial intelligence, Nature, 2023, 620(7972), 47-60.
- 2 M. Jain, S. C. Raparthy, A. Hernandez-Garcia, J. Rector-Brooks, Y. Bengio, S. Miret, et al., Multi-objective gflownets, in International Conference on Machine Learning, PMLR, 2023, pp. 14631-14653.
- 3 M. Xu, X. Yuan, S. Miret and J. Tang, ProtST: Multi-Modality Learning of Protein Sequences and Biomedical Texts, in Proceedings of the 40th International Conference on Machine

- Learning, ed. A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato and J. Scarlett, Proceedings of Machine Learning Research, 2023, vol. 202, pp. 38749-38767, available from: https://proceedings.mlr.press/v202/xu23t.html.
- 4 A. M. Bran, S. Cox, A. D. White and P. Schwaller, ChemCrow: Augmenting large-language models with chemistry tools, arXiv, 2023, preprint, arXiv:230405376, DOI: 10.48550/ arXiv.2304.05376.
- 5 M. Sim, M. G. Vakili, F. Strieth-Kalthoff, H. Hao, R. Hickman, S. Miret, et al., ChemOS 2.0: an orchestration architecture for chemical self-driving laboratories, 2023.
- 6 S. Miret, M. Skreta, B. Sanchez-Lengelin, S. P. Ong, Z. Morgan-Chan and A. Aspuru-Guzik, AI4Mat: AI for Accelerated Materials Design NeurIPS 2022 Workshop, 2022, available from: https://sites.google.com/view/ai4mat.
- 7 Y. Song, S. Miret and B. Liu, MatSci-NLP: Evaluating Scientific Language Models on Materials Science Language Tasks Using Text-to-Schema Modeling, in *Proceedings of the* 61st Annual Meeting of the Association for Computational Linguistics, ACL, 2023.
- 8 K. L. K. Lee, C. Gonzales, M. Nassar, M. Spellings, M. Galkin and S. Miret, MatSciML: A Broad, Multi-Task Benchmark for Solid-State Materials Modeling, in AI for Accelerated Materials Design - NeurIPS 2023 Workshop, 2023, available from: https://openreview.net/forum?id=josIqIStKs.
- 9 S. Miret, K. L. K. Lee, C. Gonzales, M. Nassar and M. Spellings, The Open MatSci ML Toolkit: A Flexible Framework for Machine Learning in Materials Science, Transactions on Machine Learning Research, 2023, available from: https://openreview.net/forum?id=QBMyDZsPMd.
- 10 J. S. Rutherford, Crystal Structure, in Encyclopedia of Condensed Matter Physics, ed. F. Bassani, G. L. Liedl and P. Wyder, Elsevier, Oxford, 2005, pp. 289-294, available https://www.sciencedirect.com/science/article/pii/ B0123694019006860.
- 11 Y. Zhao, E. M. D. Siriwardane, Z. Wu, N. Fu, M. Al-Fahdi, M. Hu, et al., Physics guided deep learning for generative design of crystal materials with symmetry constraints, npj Comput. Mater., 2023, 9(1), 38.
- 12 A. Nouira, N. Sokolovska and J. Crivello, Crystalgan: learning to discover crystallographic structures with generative adversarial networks. arXiv. 2018, preprint, arXiv:181011203, DOI: 10.48550/arXiv.1810.11203.
- 13 T. Xie, X. Fu, O. E. Ganea, R. Barzilay and and T. Jaakkola, Crystal Diffusion Variational Autoencoder for Periodic Material Generation, in International Conference on Learning Representations, 2022, available from: https:// openreview.net/forum?id=03RLpj-tc\_.
- 14 R. Jiao, W. Huang, P. Lin, J. Han, P. Chen, Y. Lu, et al., Crystal Structure Prediction by Joint Equivariant Diffusion, in Thirty-seventh Conference on Neural Information Processing Systems, 2023, available from: https://openreview.net/ forum?id=DNdN26m2Jk.
- 15 S. A. Meldgaard, H. L. Mortensen, M. S. Jørgensen and Structure prediction reconstructions by deep reinforcement learning, J. Phys.: Condens. Matter, 2020, 32(40), 404005.

- 16 E. Zamaraeva, C. M. Collins, D. Antypov, V. V. Gusev, R. Savani, M. S. Dyer, et al., Reinforcement Learning in Crystal Structure Prediction, 2023.
- 17 J. Damewood, J. Karaguesian, J. R. Lunger, A. R. Tan, M. Xie, J. Peng, *et al.*, Representations of materials for machine learning, *Annu. Rev. Mater. Res.*, 2023, **53**(1), 399–426.
- 18 A. Duval, V. Schmidt, S. Miret, Y. Bengio, A. Hernández-García and D. Rolnick, PhAST: Physics-Aware, Scalable, and Task-specific GNNs for Accelerated Catalyst Design, in AI for Accelerated Materials Design NeurIPS 2022 Workshop, 2022, available from: <a href="https://openreview.net/forum?id=hHercGKiXvP">https://openreview.net/forum?id=hHercGKiXvP</a>.
- 19 Y. Zhao, M. Al-Fahdi, M. Hu, E. M. Siriwardane, Y. Song, A. Nasiri, et al., High-throughput discovery of novel cubic crystal materials using deep generative neural networks, Advanced Science, 2021, 8(20), 2100566.
- 20 A. Kumar, A. Zhou, G. Tucker and S. Levine, Conservative q-learning for offline reinforcement learning, *Adv. Neural Inf. Process.*, 2020, **33**, 1179–1191.
- 21 E. O. Pyzer-Knapp, C. Suh, R. Gómez-Bombarelli, J. Aguilera-Iparraguirre and A. Aspuru-Guzik, What is high-throughput virtual screening? A perspective from organic materials discovery, *Annu. Rev. Mater. Res.*, 2015, 45, 195–216.
- 22 P. Giannozzi, S. Baroni, N. Bonini, M. Calandra, R. Car, C. Cavazzoni, et al., QUANTUM ESPRESSO: a modular and open-source software project for quantum simulations of materials, J. Phys.: Condens. Matter, 2009, 21(39), 395502.
- 23 C. W. Glass, A. R. Oganov and N. Hansen, USPEX—Evolutionary crystal structure prediction, *Comput. Phys. Commun.*, 2006, 175(11–12), 713–720.
- 24 K. Doll, J. Schön and M. Jansen, Structure prediction based on *ab initio* simulated annealing for boron nitride, *Phys. Rev. B: Condens. Matter Mater. Phys.*, 2008, **78**(14), 144110.
- 25 Y. Wang, J. Lv, L. Zhu and Y. Ma, CALYPSO: A method for crystal structure prediction, *Comput. Phys. Commun.*, 2012, 183(10), 2063–2070.
- 26 J. Li, K. Lim, H. Yang, Z. Ren, S. Raghavan, P. Y. Chen, *et al.*, AI applications through the whole life cycle of material discovery, *Matter*, 2020, 3(2), 393–432.
- 27 S. Miret, K. L. K. Lee, C. Gonzales, M. Nassar and M. Spellings, The Open MatSci ML Toolkit: A Flexible Framework for Machine Learning in Materials Science, Transactions on Machine Learning Research, 2023, available from: https://openreview.net/forum? id=QBMyDZsPMd.
- 28 L. Chanussot, A. Das, S. Goyal, T. Lavril, M. Shuaibi, M. Riviere, et al., Open catalyst 2020 (OC20) dataset and community challenges, ACS Catal., 2021, 11(10), 6059–6072.
- 29 C. Chen and S. P. Ong, A universal graph deep learning interatomic potential for the periodic table, *Nat. Comput. Sci.*, 2022, 2(11), 718–728.
- 30 A. A. Duval, V. Schmidt, A. Hernandez-Garcia, S. Miret, F. D. Malliaros, Y. Bengio, et al., FAENet: Frame Averaging Equivariant GNN for Materials Modeling, in Proceedings of the 40th International Conference on Machine Learning, ed. A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato and J. Scarlett, PMLR, vol. 202, 2023, pp., pp. 9013–9033,

- available from: https://proceedings.mlr.press/v202/duval23a.html.
- 31 S. Kim, J. Noh, G. H. Gu, A. Aspuru-Guzik and Y. Jung, Generative Adversarial Networks for Crystal Structure Prediction, *ACS Cent. Sci.*, 2020, **6**(8), 1412–1420, DOI: **10.1021/acscentsci.0c00426**.
- 32 S. Zheng, J. He, C. Liu, Y. Shi, Z. Lu, W. Feng, *et al.*, Towards Predicting Equilibrium Distributions for Molecular Systems with Deep Learning, *arXiv*, 2023, preprint, arXiv:230605445, DOI: 10.48550/arXiv.2306.05445.
- 33 E. Pan, C. Karpovich and E. Olivetti, Deep Reinforcement Learning for Inverse Inorganic Materials Design, in *AI for Accelerated Materials Design NeurIPS 2022 Workshop*, 2022, available from: <a href="https://openreview.net/forum?id=V-DQd">https://openreview.net/forum?id=V-DQd</a> iX1xJ.
- 34 F. Sui, R. Guo, Z. Zhang, G. X. Gu and L. Lin, Deep reinforcement learning for digital materials design, *ACS Mater. Lett.*, 2021, 3(10), 1433–1439.
- 35 J. N. Law, S. Pandey, P. Gorai and P. C. St John, Upper-Bound Energy Minimization to Search for Stable Functional Materials with Graph Neural Networks, *JACS Au*, 2022, 3(1), 113–123.
- 36 B. Zheng, Z. Zheng and G. X. Gu, Designing mechanically tough graphene oxide materials using deep reinforcement learning, *npj Comput. Mater.*, 2022, **8**(1), 225.
- 37 S. Levine, A. Kumar, G. Tucker and J. Fu, Offline Reinforcement Learning: Tutorial, Review. and Perspectives on Open Problems, 2020, vol. 5.
- 38 R. F. Prudencio, M. R. Maximo and E. L. Colombini, *A survey on offline reinforcement learning: Taxonomy, review, and open problems*, IEEE Transactions on Neural Networks and Learning Systems, 2023.
- 39 A. Nair, M. Dalal, A. Gupta and S. Levine, {AWAC}:
  Accelerating Online Reinforcement Learning with Offline
  Datasets, 2021, available from: https://openreview.net/
  forum?id=OJiM1R3jAtZ.
- 40 I. Kostrikov, A. Nair and S. Levine, Offline Reinforcement Learning with Implicit Q-Learning, in *International Conference on Learning Representations*, 2022, available from: https://openreview.net/forum?id=68n2s9ZJWF8.
- 41 T. Yu, A. Kumar, R. Rafailov, A. Rajeswaran, S. Levine and C. Finn, Combo: Conservative offline model-based policy optimization, *Adv. Neural Inf. Process.*, 2021, **34**, 28954–28967.
- 42 T. Xie and J. C. Grossman, Crystal Graph Convolutional Neural Networks for an Accurate and Interpretable Prediction of Material Properties, *Phys. Rev. Lett.*, 2018, 120(14), 145301.
- 43 V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, *et al.*, Human-level control through deep reinforcement learning, *nature*, 2015, 518(7540), 529–533.
- 44 S. Kurth, M. A. L. Marques and E. K. U. Gross, Density-Functional Theory, in *Encyclopedia of Condensed Matter Physics*, ed. F. Bassani, G. L. Liedl and P. Wyder, Elsevier, Oxford, 2005, pp. , pp. 395–402, available from: https://www.sciencedirect.com/science/article/pii/B0123694019004459.

- 45 A. Seidl, A. Görling, P. Vogl, J. A. Majewski and M. Levy, Generalized Kohn-Sham schemes and the band-gap problem, Phys. Rev. B: Condens. Matter Mater. Phys., 1996, 53(7), 3764.
- 46 J. Heyd, G. E. Scuseria and M. Ernzerhof, Hybrid functionals based on a screened Coulomb potential, J. Chem. Phys., 2003, **118**(18), 8207-8215.
- 47 F. Aryasetiawan and O. Gunnarsson, The GW method, Rep. Prog. Phys., 1998, 61(3), 237.
- 48 A. Jain, S. P. Ong, G. Hautier, W. Chen, W. D. Richards, S. Dacek, et al., Commentary: The Materials Project: A materials genome approach to accelerating materials innovation, APL Mater., 2013, 1(1), 011002.
- 49 C. Chen, W. Ye, Y. Zuo, C. Zheng and S. P. Ong, Graph networks as a universal machine learning framework for molecules and crystals, Chem. Mater., 2019, 31(9), 3564-3572.

- 50 H. Pan, A. M. Ganose, M. Horton, M. Aykol, K. A. Persson, N. E. Zimmermann, et al., Benchmarking coordination number prediction algorithms on inorganic crystal structures, Inorg. Chem., 2021, 60(3), 1590-1603.
- 51 D. W. Davies, K. T. Butler, A. J. Jackson, J. M. Skelton, K. Morita and A. Walsh, SMACT: Semiconducting materials by analogy and chemical theory, J. Open Source Softw., 2019, 4(38), 1361.
- 52 A. Li, D. Misra, A. Kolobov and C. A. Cheng, Survival Instinct in Offline Reinforcement Learning, in Thirty-seventh Conference on Neural Information Processing Systems, 2023, available from: https://openreview.net/forum?id=shePL2nbwl.
- 53 O. Vinyals, S. Bengio and M. Kudlur, Order matters: Sequence to sequence for sets, arXiv, 2015, preprint, arXiv:151106391, DOI: 10.48550/arXiv.1511.06391.