

Space Group Informed Transformer for Crystalline Materials Generation

Zhendong Cao,^{1,2} Xiaoshan Luo,^{3,4} Jian Lv,^{3,*} and Lei Wang^{1,5,†}

¹*Beijing National Laboratory for Condensed Matter Physics and Institute of Physics, Chinese Academy of Sciences, Beijing 100190, China*

²*School of Physical Sciences, University of Chinese Academy of Sciences, Beijing 100190, China*

³*Key Laboratory of Material Simulation Methods and Software of Ministry of Education, College of Physics, Jilin University, Changchun 130012, P. R. China*

⁴*State Key Laboratory of Superhard Materials, College of Physics, Jilin University, Changchun 130012, P. R. China*

⁵*Songshan Lake Materials Laboratory, Dongguan, Guangdong 523808, China*

(Dated: August 19, 2024)

We introduce **CrystalFormer**, a transformer-based autoregressive model specifically designed for space group-controlled generation of crystalline materials. The incorporation of space group symmetry significantly simplifies the crystal space, which is crucial for data and compute efficient generative modeling of crystalline materials. Leveraging the prominent discrete and sequential nature of the Wyckoff positions, **CrystalFormer** learns to generate crystals by directly predicting the species and locations of symmetry-inequivalent atoms in the unit cell. We demonstrate the advantages of **CrystalFormer** in standard tasks such as symmetric structure initialization and element substitution compared to conventional methods implemented in popular crystal structure prediction software. Moreover, we showcase the application of **CrystalFormer** of property-guided materials design in a plug-and-play manner. Our analysis shows that **CrystalFormer** ingests sensible solid-state chemistry knowledge and heuristics by compressing the material dataset, thus enabling systematic exploration of crystalline materials. The simplicity, generality, and flexibility of **CrystalFormer** position it as a promising architecture to be the foundational model of the entire crystalline materials space, heralding a new era in materials modeling and discovery.

CONTENTS

I. Introduction	1	A. More details of CrystalFormer	17
II. CrystalFormer	3	1. Model architectures	17
A. Model	3	2. Sampling algorithm	19
B. Training	4		
C. Sampling	5		
III. CrystalFormer learns chemical intuition by compressing materials database	5	B. Validity and novelty of generated samples	21
A. Atom embeddings and chemical similarity	6	C. Discovered crystal samples with symmetric structure initialization	24
B. Atom number distributions	6	D. Discovered crystal samples with element substitution	24
C. Wyckoff-Atom gram	7	E. Details of plug-and-play materials design	24
D. Crystal likelihoods	7		
IV. Applications	8		
A. Symmetry-conditioned random structure initialization	8	I. INTRODUCTION	
B. Structure-conditioned element substitution	9		
C. Plug-and-play materials design	10		
V. Related works	11	Machine learning methods are playing an increasingly important role in material discovery, complementing conventional computational approaches [1, 2]. Generative machine learning, in particular, has been a promising step for matter inverse design [3, 4] which goes beyond machine learning accelerated structure search [5] and property screening [6]. Generative models learn the underlying distribution of training data and generate new samples from the learned distribution. In addition, the generation process can also be controlled by conditions such as desired material properties or experiment observations. Amazing programming abilities of generative models have been demonstrated in large language model [7], text-to-image generation [8, 9], and protein design [10].	11
VI. Outlook	12		
Acknowledgments	12		
References	12		

* lvjian@jlu.edu.cn

† wanglei@iphy.ac.cn

It is anticipated that generative model-based approaches will introduce groundbreaking changes to the traditional workflows of material discovery. A generative pre-trained

<i>P1</i> world	With space group symmetry
$(100 \times 100^3)^{20} \approx 10^{160}$	$(100 \times 10 \times 100)^5 \approx 10^{25}$

TABLE I. A back-of-envelope estimate of the size of the crystalline material space. In the "P1 world", one treats crystals as if they were in the first and the least symmetric *P1* space group. For the estimate, we consider 100 possible chemical elements and 20 atoms in the unit cell with a coordinate grid size of 100 in each direction. In the case of utilizing the symmetry of a typical space group, we consider 5 symmetry inequivalent atoms occupying 10 possible Wyckoff positions. The additional factor of 100 accounts for the remaining degree of freedom for the fractional coordinates. See Refs. [27, 28] for alternate estimates of the materials space in the context of crystal structure prediction.

foundation model for crystalline materials is a key step towards such a lofty goal. However, despite intensive efforts [11–22], the current generative models for crystalline materials **fall short to match the success of other domains**. Simply scaling the compute and model size of the current crystal generative model may not be feasible because the amount of **high-quality data for crystalline materials is much less than compared to language and image domains**. Therefore, leveraging the inherent inductive biases specific to crystalline structures for more data-efficient generative modeling is essential, as has been pursued in some of recent works [23–26].

The space group symmetry due to the joint outcome of the rotational and translational symmetry in space is arguably the most important inductive bias in the modeling of crystalline materials. There are in total 230 space groups [29] for three-dimensional crystal structures. Nature exhibits a preference for symmetric crystal structures, a tendency that may be attributed to the symmetry inherent in the interatomic interactions, which, in turn, are governed by the fundamental forces acting between elementary particles. As a result, the appearance of crystalline materials in the first and the least symmetric space group *P1* is rare [30], with many instances potentially even being misclassified [31]. Failing to match the space group distribution of nature in machine learning-generated materials is regarded as a matter of serious concern [32].

Space group symmetry imposes significant constraints on a crystal. First of all, the space group identifies the crystal system to which a crystal belongs, thereby limiting the permissible values for the lattice parameters that define the length and angles of the crystal's unit cell. Moreover, the symmetry operations associated with a given space group ensure that symmetry equivalent atoms are consistently mapped among themselves in the crystal. This requirement enforces strict conditions regarding the types of chemical elements present, their specific locations within the crystal, and the number of each chemical species in the unit cell. A key concept to express these constraints is the Wyckoff positions, which delineate unique areas within a unit cell that are defined by the symmetry operations of the crystal's space group. These positions are represented as fractional coordinates, enabling precise definition relative to the unit cell's axes. For example,

Fig. 1(a) shows the Wyckoff positions for the space group $R\bar{3}c$ (No. 167). The Wyckoff positions are labeled by letters in the alphabet, starting from special points in the bottom to general positions in the top. The multiplicity counts the number of equivalent positions connected by the space group symmetry operations. All of them should be occupied by the same type of atoms to uphold the space group symmetry. For example, the top row of the table in Fig. 1(a) contains the general position (x, y, z) that can be mapped to 36 positions under the symmetry operations of the $R\bar{3}c$ space group.

Nature tends to place atoms in those special Wyckoff positions at the bottom of the table. For example, we highlight the occupied Wyckoff positions of calcite (CaCO_3) crystal in Fig. 1, associated with the $R\bar{3}c$ space group. One sees that the Wyckoff letter '6a' and '6b' deterministically define the locations of the carbon and calcium atoms within the unit cell. In addition, it follows that $a = b$, and $\alpha = \beta = 90^\circ, \gamma = 120^\circ$ as the $R\bar{3}c$ space group belongs to the trigonal crystal system. Ultimately, despite having 30 atoms in the unit cell, there are only three continuous degrees of freedom for the CaCO_3 structure: the x-coordinate of oxygen atom 0.257 and the lattice constants $a = b = 4.99\text{\AA}$ and $c = 17.07\text{\AA}$. All other information about the crystal structure can be specified via discrete data such as the Wyckoff letters and chemical species.

The prominent discrete and sequential features illustrated in Figure 1 are ubiquitous in crystalline materials. The Wyckoff positions not only specify possible locations of atoms in the unit cell, but their associated multiplicities also put strong constraints on the number of atoms. Therefore, space group symmetry significantly reduces the degrees of freedom of crystalline materials. Failing to exploit this information in generative modeling not only renders learning inefficient, it also severely impairs the generalization ability of the model. For example, the performance of the generative model quickly deteriorates as the number of atoms increases due to it is challenging to generate highly symmetric crystal structures [16]. On the other hand, statistical analysis shows that the Wyckoff sequences of known inorganic compounds [33] are far from being exhausted, implying there are statistical correlations to be exploited to compress the materials database.

In this paper, we introduce **CrystalFormer**, an autoregressive transformer for generative modeling of crystalline materials. **CrystalFormer** models the joint probability distribution of Wyckoff positions, chemical species, and lattice parameters of crystals with a given space group. By treating the Wyckoff positions as the first class citizen in the model, **CrystalFormer** seamlessly and rigorously integrates the space group symmetry into crystal probabilistic modeling. As shown in Table I, explicit modeling of the Wyckoff positions greatly reduces the space of crystalline materials. The space group-informed transformer exploits this fundamental inductive bias to greatly simplify the learning and generation of crystals.

The organization of the present paper is as follows. Section II introduces the model architecture, training, and sampling of the **CrystalFormer**. Section III reveals the chemical intuition encoded in the trained model by inspecting its weights and generated crystal samples. Such inspection also

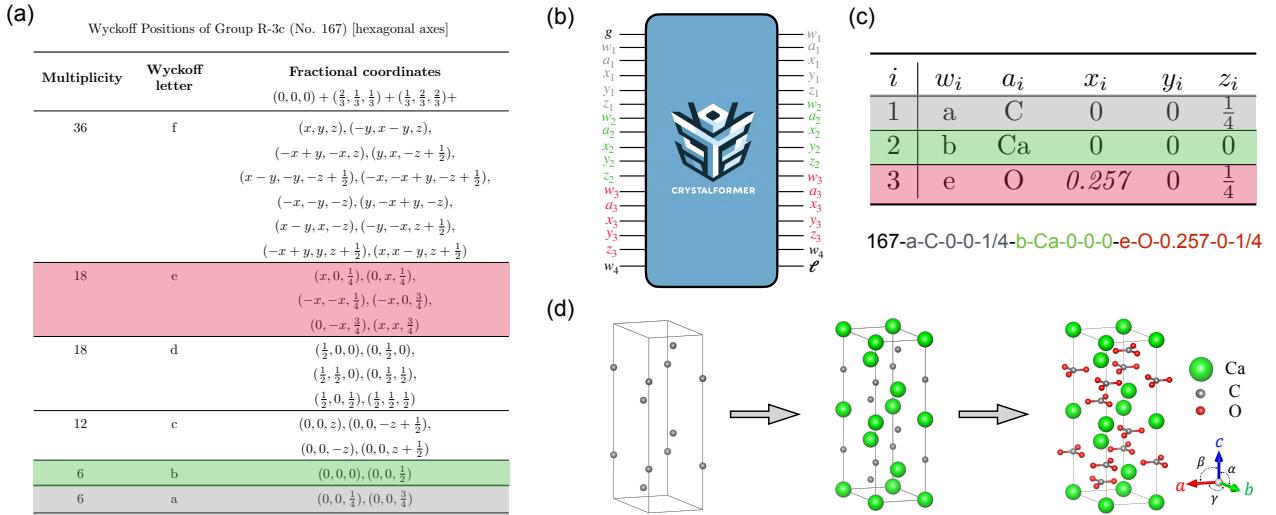


FIG. 1. (a) The Wyckoff positions of the $R\bar{3}c$ space group (No. 167). We highlight the occupied Wyckoff positions of calcite CaCO_3 crystal which belongs to this space group. Carbon, calcium, and oxygen atoms occupy the '6a', '6b', and '18e' positions, respectively. (b) The CrystalFormer is a decoder-only autoregressive transformer that models the space group controlled crystal structures by predicting probabilities of the Wyckoff letter w_i , chemical element a_i , and fractional coordinates (x_i, y_i, z_i) of each symmetry inequivalent atom, and finally, the lattice parametrized by ℓ sequentially. (c) The crystal data of CaCO_3 is summarized in a table. In the table, the x-coordinate of oxygen atom $x_3 = 0.257$ is the only continuous variable that needs to be predicted. All other fractional coordinates are fixed by discrete data like the space group number and Wyckoff letters. The string below the table shows the sequential representation of the CaCO_3 crystal with space group, Wyckoff letter, and atom species as the input to the CrystalFormer model. (d) Autoregressive generation of the crystal. One first places carbon atoms at the '6a' position, then places calcium atoms at '6b' position, and finally places oxygen atoms at '18e' position. In each step of the sampling procedure, there is a choice of the Wyckoff positions, atom species, and the fractional coordinates if they are still unspecified.

builds up an understanding of the strength of the model. Section IV demonstrates applications CrystalFormer for symmetric structure initialization and element substitution, where it exhibits great efficiency and generality compared to existing approaches. Furthermore, we showcase CrystalFormer's ability for property-guided materials design in a plug-and-play manner. Section V puts the key contribution of present work in the broad context of crystal generative models. Finally, Section VI outlooks for possible extension and future of the present line of research. We have released the codes and trained model at [34].

II. CRYSTALFORMER

We will introduce the design, training, and sampling of the CrystalFormer model.

A. Model

To exploit the space group symmetry of the crystal, we focus on the Wyckoff positions of symmetry-inequivalent atoms. Wyckoff letters follow the alphabetical order, where "a" stands for the positions with the highest order of symmetry for the given space group. Later letters in the alphabet indicate

more general positions with reduced site symmetries. Note that the information of the space group number and Wyckoff letter fully determine the multiplicities. In cases where the atom positions are not fully fixed by the Wyckoff letter, we will also consider the remaining fractional coordinates, e.g. the x -coordinate of the oxygen atoms in the CaCO_3 example shown in Fig. 1. To generate crystals, one samples the Wyckoff letter, chemical element, and fractional coordinates of each atom sequentially. The sampling procedure starts from special higher symmetry sites with smaller multiplicities and then goes on to general lower symmetry regions with larger multiplicities.

With these considerations, we define a crystal data as $\mathbf{C} = \{\mathbf{W}, \mathbf{A}, \mathbf{X}, \mathbf{L}\}$. Here $\mathbf{W} = [w_1, w_2, \dots, w_n]$ are Wyckoff letters and $\mathbf{A} = [a_1, a_2, \dots, a_n]$ are chemical species. Here, n stands for the number of symmetrically inequivalent atoms in the conventional unit cell. For example, as shown in Fig. 1(b) one has $n = 3$ for CaCO_3 . Explicitly including the Wyckoff letter in the generative modeling is the key of the present work. Next, $\mathbf{X} = [(x_i, y_i, z_i)] \in \mathbb{R}^{n \times 3}$ are the fractional coordinates of symmetrically inequivalent atoms. Lastly, $\mathbf{L} = [a, b, c, \alpha, \beta, \gamma]$ denotes the lattice parameters of the conventional unit cell of the crystal.

The central quantity to focus on is the conditional probability of a crystal \mathbf{C} given the space group number $g \in [1, 230]$: $p(\mathbf{C}|g)$. Since the space group is a fundamental characteriza-

tion for crystalline materials, g is a key control variable that greatly simplifies the distribution over the entire crystal materials space. In practical applications of crystal structure prediction and material design, the space group can either be considered separately as a control variable or predicted based on material composition [35–38].

We express the space group conditioned probability distribution of crystals as an autoregressive product of conditional probabilities

$$\begin{aligned} p(\mathbf{C}|g) = & p(w_1|g) \times \\ & p(a_1|g, w_1) \times \\ & p(x_1|g, w_1, a_1) \times \\ & p(y_1|g, w_1, a_1, x_1) \times \\ & p(z_1|g, w_1, a_1, x_1, y_1) \times \dots \times \\ & p(\mathbf{L}|g, w_1, a_1, x_1, y_1, z_1 \dots, w_n, a_n, x_n, y_n, z_n). \end{aligned} \quad (1)$$

At first sight, it may appear unnatural to employ an autoregressive model for crystals since there seems to be no obvious order for atoms in the unit cell. However, the sequential nature of Wyckoff positions suggests a natural way to arrange symmetrically inequivalent atoms in an alphabetical order of the Wyckoff letters. Following this key observation, we represent crystal data as sequences of space groups, Wyckoff letters, chemical species, and fractional coordinates of each symmetrically inequivalent atom. Together with the information lattice parameters, such sequence fully characterizes the compositional and structural information of crystalline material. Since statistical analysis reveals that anions are in less symmetric positions than cations for inorganic crystals [30], one would expect that anion atoms will typically appear after cation atoms in such a sequence. For example, CaCO3 is represented as a string "167-a-C-0-0-1/4-b-Ca-0-0-0-e-O-0.257-0-1/4". Autoregressive sampling of such a string means the model generates the crystal by placing the atoms sequentially into the unit cell, starting from the special position with high site symmetry to the general position with the lowest site symmetry, see Fig. 1(d).

We model the conditional probability of the Wyckoff letters \mathbf{W} and chemical species \mathbf{A} as categorical distributions. On the other hand, we model the conditional probability of the fractional coordinates \mathbf{X} as a mixture of von Mises distribution for continuous periodic variables. For Wyckoff positions with multiplicities greater than one, we only consider the first of fractional coordinates that appear in the international tables for crystallography [39]. Lastly, we model the conditioned distribution of lattice parameters as a Gaussian mixture model.

We build **CrystalFormer**, an autoregressive transformer [40] to model the space group conditioned-probability distribution of crystalline materials Eq. (1). The space group number g is the first input to **CrystalFormer**. The remaining inputs are the Wyckoff letter, chemical species, and fractional coordinates of each atom. One can go through the table of Fig. 1(b) in a raster order to collect these atomistic features. We feed vector embeddings of the space group number, Wyckoff letter, and the chemical species input to the **CrystalFormer**. In particular, we also concatenate the vector embedding of g to all other inputs since it is the key control

variable for the crystal generation. Moreover, we have also provided the multiplicity of each Wyckoff position as an additional feature. The multiplicity can be easily inferred from the space group and the Wyckoff letters. We feed the fractional coordinates as Fourier features into the transformer so that the model preserves the periodicity of the unit cell [13, 41]. We pad the atom sequence up to a maximum length and treat the output as parameters of the conditional probability distribution Eq. (1), see Fig. 1(b). At the location of the first padding atom, we predict the lattice parameters.

We implement a number of constraints in the model to further reduce its phase space. First, the Wyckoff letters should be valid for the given space group. For example, for the space group $R\bar{3}c$ (No. 167) the Wyckoff letters go from 'a' up to 'f'. Second, we require that the Wyckoff letters w_i follow alphabetical order in the sequence [42]. Lastly, the Wyckoff positions with no free fractional coordinates (such as 'a', 'b', and 'd' positions in the $R\bar{3}c$ space group) can only be occupied once. Those constraints are implemented by setting the logit biases of Wyckoff letters to mask out invalid sequences [43, 44].

The design of **CrystalFormer** focuses mostly on the space group symmetries which we believe to be the most important inductive bias for crystalline materials. This design decision significantly impacts the treatment of other symmetries. First, it is often possible to place the origin of the unit cell at the inversion center of the specified space group. The chosen origin naturally fixes the continuous translation invariance of fractional coordinates. Second, by only considering symmetry-inequivalent atoms and labeling them with Wyckoff letters, one fixes most of the permutation invariance in the representation. For those Wyckoff positions with continuous degrees of freedom, there may be multiple symmetry-inequivalent atoms with the same Wyckoff letters. We arrange these atoms according to the lexicographic order of fractional coordinates [45] in the sequence. Note that in a crystal environment, the same type of atoms occupying different Wyckoff positions could be regarded as distinguished particles as they generally have different site symmetry. Lastly, the periodicity of the fractional coordinates is respected in **CrystalFormer** since they are treated as periodic variables under the von Mises distribution.

B. Training

The **CrystalFormer** is trained by minimizing the negative log-likelihood over training dataset

$$\mathcal{L} = - \mathbb{E}_{\mathbf{C}, g} [\ln p(\mathbf{C}|g)]. \quad (2)$$

Writing out $p(\mathbf{C}|g)$ according to Eq. (1), the objective function contains the negative log-likelihood of discrete variables such as Wyckoff letters \mathbf{W} and chemical species \mathbf{A} , as well as continuous variables such as fractional coordinates \mathbf{X} and the lattice parameters \mathbf{L} . In the objective function, for continuous variables \mathbf{X}, \mathbf{L} we consider only active ones that are not fixed by the space groups and Wyckoff letters. In this way, those

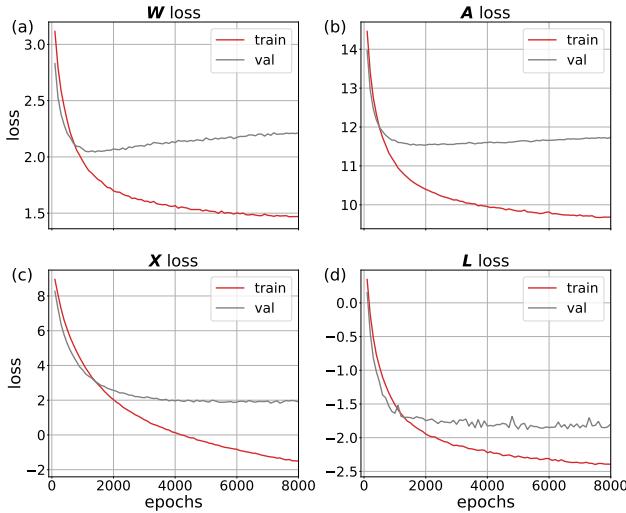


FIG. 2. Break up of the training and validation losses for (a) Wyckoff letters, (b) chemical species, (c) fractional coordinates, and (d) lattice parameters over training epochs.

special fractional coordinates (e.g. 0, $\frac{1}{4}$) and lattice parameters (e.g. 90° , 120°) do not contribute to the loss function.

In the present work, we **train the CrystalFormer using the MP-20 dataset** [11] which is a popular dataset that represents a majority of experimentally known crystalline materials at ambient conditions with no more than 20 atoms in the primitive unit cell. **The training dataset contains 27136 crystal structures.** The subdivision of the training samples according to the space group has greatly **reduced the number of samples in each space group category**. On top of that, the distribution of training samples is quite uneven among the space groups, which reflects the imbalanced distribution of crystals over space groups in nature [30]. In fact, **there is no training data in 61 out of 230 space groups** as shown in Fig. S2. Nevertheless, we still employ the MP-20 as the training set so that the performance of the model can be more easily **gauged with the others in the literature**, see appendix B. Figure 2 shows a breakup of the learning curves for the Wyckoff position, chemical species, fractional coordinates, and lattice parameters. We select the model checkpoint with the lowest total validation loss to generate crystal samples.

C. Sampling

To sample crystals from the CrystalFormer, one needs to **specify a space group number and a list of possible chemical elements**. The CrystalFormer samples the atoms one by one, starting **from more symmetric specific positions with lower multiplicities till less symmetric general positions with larger multiplicities**. We use the information of the space group and Wyckoff letter to control the sampling of fractional coordinates. By applying the symmetry projection to the sampled fractional coordinate, one rectifies it and ensures the generated fractional coordinates are compatible with the Wyck-

off positions. One can also mask out the logits of chemical species so that only a number of selected elements will be sampled. The number of symmetrically inequivalent atoms may fluctuate in the sampling procedure. Once one has sampled a padding atom, the model predicts the lattice parameters under the space group constraint. Moreover, we introduce a temperature parameter T in the sample distribution $p(C|g)^{1/T}$. With $T < 1$ we will draw samples from a sharper distribution, while $T > 1$ gives more diversity in the generated samples. In the present paper, we will generate crystals using temperature $T = 1$ unless mentioned explicitly.

Besides autoregressive sampling, one can also perform Markov chain Monte Carlo (MCMC) sampling based on the likelihood Eq. (1) of the CrystalFormer. MCMC sampling of crystals is a way to walk through the crystalline materials starting from an existing crystal structure. At each step of the random walk, one proposes an element substitution, atom position shift, or lattice deformation to change the crystal from C to C' , then accepts or rejects the proposal according to the model probability according to the Metropolis acceptance rate $\min\left[1, \frac{p(C'|g)}{p(C|g)}\right]$. MCMC sampling is particularly useful for incorporating additional constraints or guidance in the sampling procedure. Moreover, during the burn-in phase of such MCMC sampling, the generated samples will be similar to the starting material, which may be a desired feature in certain cases.

III. CRYSTALFORMER LEARNS CHEMICAL INTUITION BY COMPRESSING MATERIALS DATABASE

Nature tends to favor symmetrical crystal structures. Crystallographic space groups capture and quantify this inductive bias of nature, thereby significantly simplifying the spaces of crystal materials. In light of the space group symmetries, crystals also have an unexpected yet natural sequential and discrete representation, which derives from two tables in nature: the periodic table of elements determined by quantum mechanics and the table of Wyckoff positions of the 230 space groups determined by group theory. To construct a certain crystal, we only need to select atoms from the periodic table and place them sequentially into the Wyckoff positions in the unit cell. In this crystal language, the “word order” is determined by the alphabetical order of Wyckoff letters, the “grammar” corresponds to the solid-state chemistry rules, and the “synonyms” represent interchangeable elements (Sec. III A), the “sentence length” correspond to atom number in the unit cell (Sec. III B), and the “idioms” correspond to common chemical coordination (Sec. III C).

CrystalFormer employs an autoregressive transformer to learn the crystal language, thereby exploring yet-to-be-discovered crystalline materials. It compresses and internalizes the crystal materials database, expressing solid-state chemical knowledge through neural network parameters; reflecting the associative ability of material space through neural network activations; and describing chemical intuition through the model probability (Sec. III D). Similar to generative models used for generating text, images, and videos,

CrystalFormer can directly generate “realistic” crystal materials. However, rather than worrying about the fake contents of AI-generated media, these AI-generated crystal materials could potentially be synthesized and be useful to human civilization.

Next, we will inspect the learned features and sample statistics of the model to build up an understanding of the CrystalFormer. We carry out inspections for a few selected space groups. The findings are nevertheless general. These findings not only provide confidence a bit of understandings of the model, but also direct us to the suitable applications of CrystalFormer.

A. Atom embeddings and chemical similarity

Figure 3 visualizes the cosine similarity of the learned vector embedding of the chemical species. Red colors in the figure indicate similar chemical species identified by the model. One sees the chemical similarity within groups of elements show up as off-diagonal red stripes. Moreover, there are visible clusters for Lanthanide elements (La–Lu). The plot also suggests the similarity between the lanthanides and other rare-earth elements (Y and Sc). The features shown in Fig. 3 are strikingly similar to the similarity map constructed purposely based on substitution pattern [46, 47] which was later used for substitution-based material discoveries [5, 48]. In the context of language modeling, the chemical similarities correspond to synonyms of chemical species tokens. Having the ability to learn them from data [19, 46, 47, 49–53] is an encouraging signal that the model is picking up atomic physics to be able to generate reasonable crystal structures.

B. Atom number distributions

The number of atoms corresponds to the length of non-padding atoms in the sequential representation, which is captured well by CrystalFormer. Figure 4 presents the histogram of the total number of atoms in the conventional unit cell for several space groups. One sees a nice agreement between the atom number distribution in the test dataset and the generated samples. In addition, it appears that space group g is the key latent variable that decomposes the multi-modal atom number distribution of crystals. This is understandable because the number of atoms is determined by the sum of the multiplicities of occupied Wyckoff positions. Therefore, the space group symmetry is a key control variable for the atom number distribution. Incorporating Wyckoff positions information into the CrystalFormer model architecture removes the necessity of querying the training data to find out the number of atoms for a targeted space group [16].

Recently, Ref. [54] reports an abundance of inorganic compounds whose primitive unit cell contains a number of atoms that is a multiple of four. There are different ways to reason about the observed “rule of four” depending on one’s view of how a crystal is formed. For example, one can often break inorganic solids into polyhedra as building blocks.

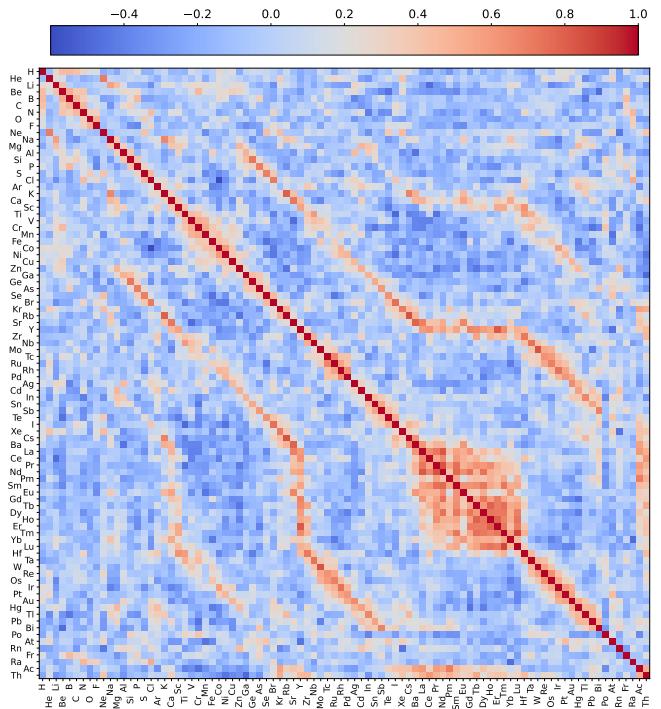


FIG. 3. The cosine similarity matrix for the chemical species based on the learned vector embeddings. The reddish color suggests similar chemical elements in the crystal environment.

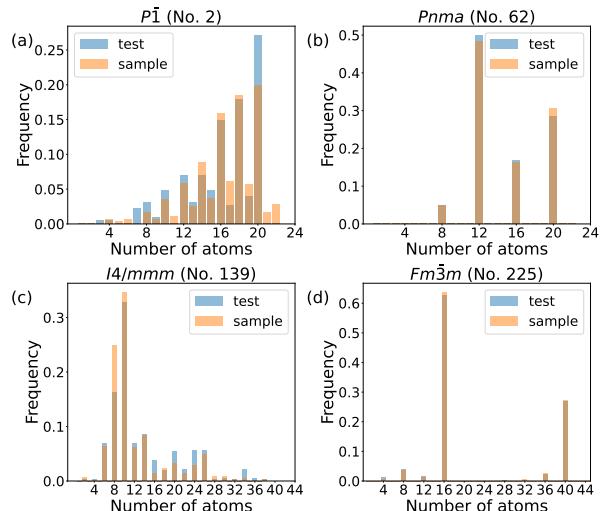


FIG. 4. The histogram for *total number* of atoms in the unit cell for several space groups in the test dataset and in the generated samples.

Otherwise, Ref. [55] considers the most probable values of the number of atoms in a formula unit and the number of formula units per primitive cell. In line with the discussion here, the “rule of four” is the combination of three factors 1) the distribution of crystalline materials among space groups [30]; 2) the distribution of atoms in Wyckoff positions [33] of a given space group; and 3) the multiplicities of Wyckoff positions and multiplicities of conventional versus cells. The first

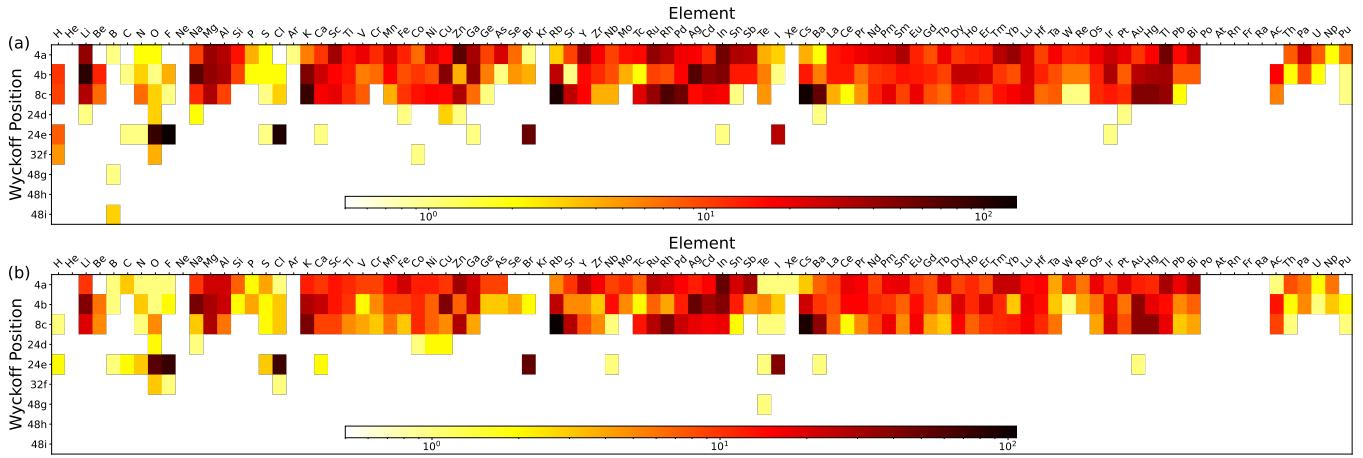


FIG. 5. The heat map for Wyckoff positions and atom species of (a) the test dataset and (b) generated samples for the $Fm\bar{3}m$ space group (No. 225). It is an analog of bigram frequency statistics of language modeling, which shows where atoms tend to occupy in the unit cell.

two are statistical rules determined by the inter-atomic interactions while the third one is a mathematical fact of space group theory. In the end, the point we want to make is that the design of **CrystalFormer** and its associated crystal representation allow it to learn the "rule of four" and many other to-be-discovered "rules", which manifest themselves as marginal statistics of learn probability distribution. Most importantly, **CrystalFormer** will utilize these "empirical rules" when generating novel yet reasonable crystal samples.

C. Wyckoff-Atom gram

Figure 5 shows heat maps of Wyckoff positions and chemical species for the $Fm\bar{3}m$ space group (No. 225). The heat map is analogous to bigram frequency statistics in language modeling. In the present context, it reveals interesting solid-state chemistry knowledge related to where each atom tends to appear in a unit cell. First of all, one sees that most atoms occupy special Wyckoff positions (Wyckoff letters at the beginning of the alphabet) with higher site symmetries. The distribution of generated data is in agreement with test data and recent statistics [33]. Moreover, there are vertical blanks at the locations of inert elements (He, Ne, Ar...) as they are rare in crystalline materials. Lastly, one sees that oxygen and halogen elements (F, Cl, Br, I) appear quite often in the Wyckoff position "24e", which is a consequence of their electronegativities [30]. Overall, we see the **CrystalFormer** has learned these key motifs for generating crystalline materials. Nevertheless, one sees that several Wyckoff locations of the hydrogen are missing in the generated samples compared to the test dataset. We believe that is due to that the hydrogen element takes only about 0.4% in the training data for the $Fm\bar{3}m$ space group. Collecting more data with better coverage of elements will be crucial to further boost the performance of the current model.

Along the same line of thoughts, coordination polyhedra [56] and lattice structure [57] manifest themselves as

higher-order n-gram correlations of Wyckoff position and atom species in the crystal language, which will be captured by the **CrystalFormer**. There has been a long history of mining empirical chemistry rules encoded in materials data and then using them to instruct the search of crystal structures [47, 58–62]. Our analysis shows that **CrystalFormer** ingests chemical intuition, be it speakable or unspeakable, in the training data for generating new materials.

D. Crystal likelihoods

CrystalFormer compresses chemistry knowledge stored in the material dataset into its parameters. In addition to generating crystal samples, **CrystalFormer** can also compute the likelihoods of crystals via Eq. (1). Therefore, it is possible to employ **CrystalFormer** in the Monte Carlo search of crystal structures besides sampling materials directly as a generative model.

Figure 6 shows the agreement of the likelihoods of generated samples and samples in the test dataset. We also visualize structures of a few generated samples which are deemed to be very likely, typical, and unlikely according to their likelihood values. We have checked that likelihood is related to the energy of the crystal by locally perturbing the fractional coordinates and lattice parameters. However, we did not observe a correlation between the likelihood of these crystals and their energies on a global scale. We envision the landscape of likelihoods is much less rough compared to the potential energy surface of crystalline materials. Intuitively, it means that the **CrystalFormer** compresses the materials space into a more compact space without many holes that correspond to infeasible high energy states. Therefore, likelihood-based exploration of the crystal space discussed in Sec. II C can be more efficient compared to traditional sampling approaches based on the physical energy functions.

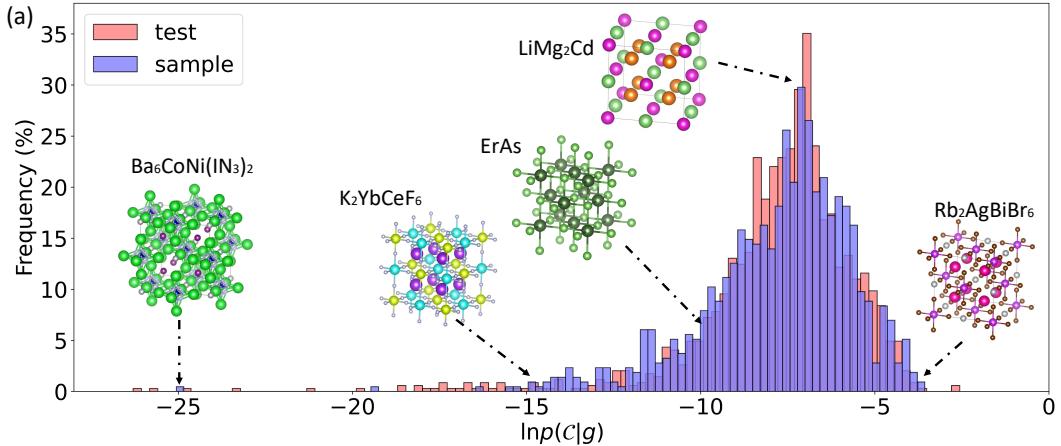


FIG. 6. The log-likelihood histogram for 1000 samples in the $Fm\bar{3}m$ (No. 225) space group and the test dataset. The insets visualize the crystal structure of a few generated samples. $Rb_2AgBiBr_6$ and $LiMg_2Cd$ are in the training dataset. $ErAs$ is in the validation dataset. K_2YbCeF_6 and $Ba_6CoNi(IN_3)_2$ are not in the MP-20 or Materials Project database.

IV. APPLICATIONS

We now move on to the practical applications of **CrystalFormer** to materials discovery and design. Compared to many existing materials generation models, **CrystalFormer** offers precise control over space group symmetry and enables efficient computation of model likelihood. These unique features open a wide range of possibilities for integrating it with existing computational software and machine-learning models in a flexible way as we demonstrate below. For these applications, we have excluded radioactive elements from the samples [32].

A. Symmetry-conditioned random structure initialization

Crystal structure prediction has long been the dream of solid-state chemistry and computational material science researchers [63]. Typical crystal structure prediction workflow consists of two steps. First, one randomly initializes a batch of diverse crystal structures as candidates. Second, one optimizes the crystal structures via local and global optimization strategies. Utilizing space group symmetries plays a crucial role in both steps, as symmetry enlarges the span of the energy distribution [64–66] and reduces the search space.

It is a common practice for crystal structure prediction software [65–70] and structure search [71–73] to exploit space group symmetry in the crystal structure initialization. However, such an initialization approach faces combinatorial difficulty as the number of chemical species and atoms in the unit cell grows. The **CrystalFormer** is ready to act as a drop-in replacement of random structure initialization for crystal structure prediction. In this way, one bypasses the curse-of-dimensionality of exact enumeration [65] with a data-driven probabilistic approach. Moreover, the ability of **CrystalFormer** to generate diverse and near-stable structures can greatly reduce the computational costs of down-

stream optimizations.

We select seven space groups $P\bar{1}$ (No. 2), $C2/m$ (No. 12), $Pnma$ (No. 62), $I4/mmm$ (No. 139), $R\bar{3}m$ (No. 166), $P6_3/mmc$ (No. 194), and $Fm\bar{3}m$ (No. 225) as representatives of the seven crystal systems. We randomly generate 100 crystals for each space group using **CrystalFormer**. On the other hand, we employ PyXtal [66] to generate crystal samples with the same stoichiometry in the same space groups. We then carry out structure relaxation using density functional (DFT) calculations.

Figure 7(a)–(g) shows the average energy difference to the energy of final structures versus DFT relaxation steps. We neglected the structures whose energy changes and energy change intervals per step during relaxations exceeded 10 eV/atom to eliminate the impact of erroneous steps. One sees that **CrystalFormer** samples generally reach lower energies in fewer relaxation steps. This is especially true for space groups with higher symmetries. The ability to initialize diverse and high-quality crystal structures enables one to discover more stable materials faster. Figure 7(h) shows the histogram of energy above the convex hull constructed by the Materials Project database. The dashed line denotes the criterion $E_{\text{hull}} < 0.1$ eV/atom [74] for selecting stable materials. Among these candidates, we found 34 and 12 relaxed structures with **CrystalFormer** and PyXtal initializations that are not contained in the MP-20 dataset. We summarize them in Table S3 and Table S4 of appendix C.

Table II lists detailed statistics of structure-relaxed samples in seven representative space groups. Overall, we find that the **CrystalFormer** generated structures are of higher quality, especially for those space groups with higher symmetry. This observation is supported by the fact that the DFT relaxation often retains the space group symmetry. The average root mean squared displacement (RMSD) computed for these converged structures is 0.961 Å and 1.743 Å for **CrystalFormer** and PyXtal initialization, respectively. In addition, the average energy above the convex hull also confirms the samples

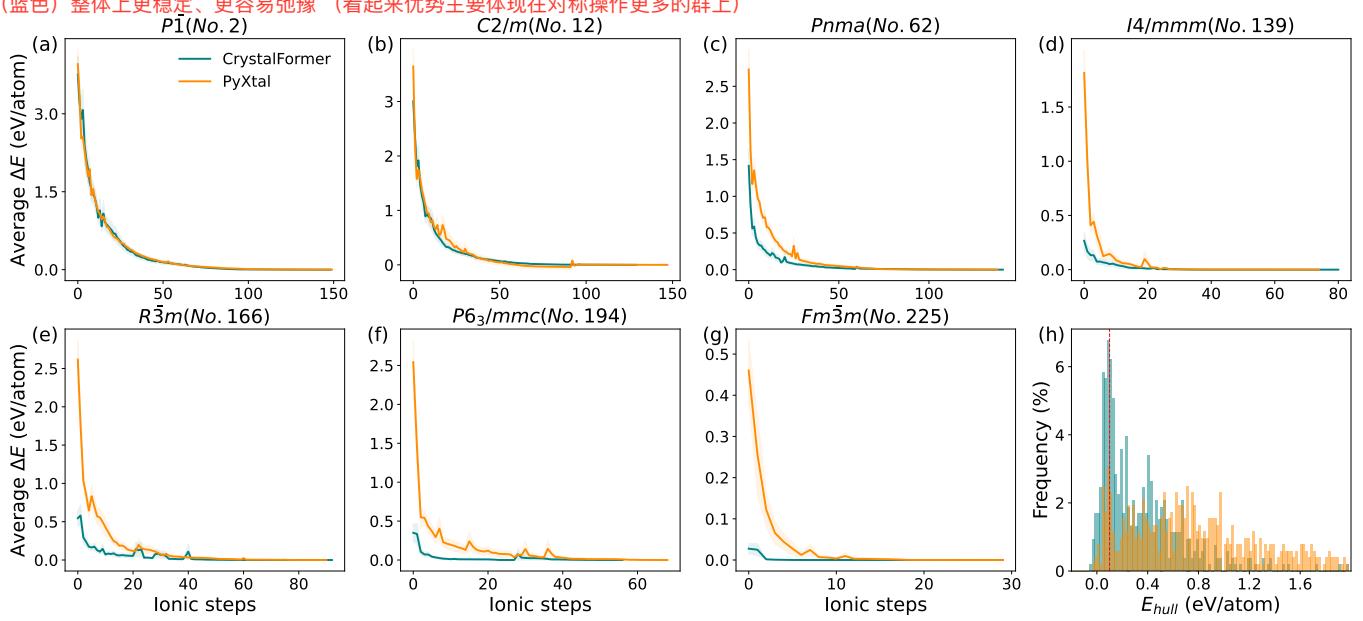


FIG. 7. (a)-(g) Average energy difference versus relaxation steps for seven representative space groups. (h) The histograms of energy above the convex hull for relaxed crystal structures. The dashed line indicates the criterion for selecting candidates for stable materials listed in Appendix C since materials with $E_{\text{hull}} < 0.1$ eV/atom are usually metastable and have the potential to be synthesized [74].

TABLE II. For each space group we randomly generate 100 crystal structures with the same composition using CrystalFormer and PyXtal. We carry out energy relaxation using DFT calculations and report the number of converged samples, the number of structures that maintain the original space group symmetry, the average RMSD between generated and relaxed structures, and the averaged energy above the convex hull.

Space group	Crystal system	Converged structures \uparrow		Retain symmetry \uparrow		RMSD $^1\text{\AA}$ \downarrow		E_{hull}^1 (eV/atom) \downarrow	
		CrystalFormer	PyXtal	CrystalFormer	PyXtal	CrystalFormer	PyXtal	CrystalFormer	PyXtal
$P\bar{1}$ (No. 2)	Triclinic	46	67	45	67	1.929	2.038	1.034	0.913
$C2/m$ (No. 12)	Monoclinic	55	72	53	67	1.846	2.235	1.233	1.660
$Pnma$ (No. 62)	Orthorhombic	77	83	76	66	1.086	1.940	0.313	1.633
$I4/mmm$ (No. 139)	Tetragonal	91	81	88	63	0.405	1.160	0.240	1.100
$R\bar{3}m$ (No. 166)	Trigonal	83	74	80	71	1.028	2.183	0.352	2.489
$P6_3/mmc$ (No. 194)	Hexagonal	97	77	96	60	0.355	2.129	0.324	4.100
$Fm\bar{3}m$ (No. 225)	Cubic	98	96	95	92	0.076	0.518	0.214	0.483
Average		78.1	78.6	76.1	69.4	0.961	1.743	0.530	1.769

¹ Calculated on the converged structures.

generated by CrystalFormer are indeed much closer to the DFT local minimum than PyXtal initialization.

B. Structure-conditioned element substitution

Mutation of known crystals is a prominent approach to materials discovery. For example, one can employ a machine-learned force field to relax crystal structures [5, 75–77] after element substitutions. In the lens of generative modeling, the machine learning force field can be regarded as the energy-based model or Boltzmann machines. A potential drawback of exploring materials space with an energy-based model is the slow mixing or even ergodicity issue posed by the rough land-

scape of the potential energy surface. In this sense, element substitutions provide a variety of initial seeds, compensating for the limitation of energy-based exploration. Having an alternative measure of crystal likelihood other than the potential energy surface opens a way to employ the model likelihood as a guide for structure search.

Many crystal structures can be traced back to a few simple, highly symmetrical types. Numerous crystals share the same structural prototype but differ in composition, such as perovskite (ABX_3), spinel (AB_2X_4), fluorite (AX_2), and so on. Figure 8(a) shows double perovskite crystal structures $A_2BB'X_6$ which belong to the $Fm\bar{3}m$ (No. 225) space group. There are hundreds of known double perovskites with significant interests in their semiconducting, ferroelectric, thermo-

electric, and superconducting properties [78]. Finding more stable materials with this structure prototype using brute force enumeration and high-throughput calculation is a computationally demanding task [79]. We will generate new double perovskites with **CrystalFormer** and demonstrate its advantage of over standard element substitution methods.

Figure 8(a) shows the string representation of double perovskites. To generate candidates of double perovskites, we use **CrystalFormer** to carry out string in-filling tasks. Since the autoregressive sampling of the atoms is insufficient to take into account non-causal information in the sequence, we employ MCMC to sweep through the sequence and update chemical species and fractional coordinates [80]. The acceptance rate for these MCMC updates makes use of the marginalized probability for elements and fractional coordinates as the lattice parameter that appears at the end of the sequence can be integrated. Only after the MCMC sampling has been thermalized, we sample the lattice parameters autoregressive to account for the adjustment of the unit cell for given atoms and occupations. We use **CrystalFormer** to generate 100 candidates as the initial DFT relaxation.

As a comparison, we also employ the **SubstitutionPredictorTransformation** function [46] implemented in pymatgen [81] to perform element substitution for the crystals with double perovskite structures in the training dataset. The substitution probabilities come from data-mining of ICSD dataset [46]. After the substitution, we use **DLSVolumePredictor** [82] function of pymatgen to predict the volume of the structure. This lattice scaling scheme relies on data-mined bond lengths to predict the crystal volume of a given structure. To collect 100 candidates in the ionic substitution approach we have set the probability threshold of **SubstitutionPredictorTransformation** to 0.01, which is smaller than the typical values adopted in Ref. [48].

The RMSD computed for the DFT-relaxed structures is 0.211 Å and 0.342 Å for **CrystalFormer** and ionic substitution [46], respectively. This observation confirms that the structures generated by **CrystalFormer** are closer to the DFT local minimum than ionic substitution [46]. Moreover, Figure 8(b) shows the histogram of energy above the convex hull of the Materials Project database. Overall, **CrystalFormer** and ionic substitution [46] found 9 and 3 double perovskites with $E_{\text{hull}} < 0.1\text{eV/atom}$ which are not contained in the MP-20 dataset, details in appendix D. The superior performance of **CrystalFormer**-guided MCMC is understandable since its likelihood takes into account the context of space group and atomic environment rather than marginal two-body correlation [46] in ionic substitution. The ionic substitution approach also shows two limitations in practical applications. First, some of the ions in the compound can not be substituted as they are missing in the probability table. Second, the approach relies on the calculability of the elements' valence states.

As a final remark, although the discussion here focuses on generating crystals with given prototype structures, the generation of crystals with a given crystal lattice [83] is also feasible with **CrystalFormer**. This is because the crystal lattice

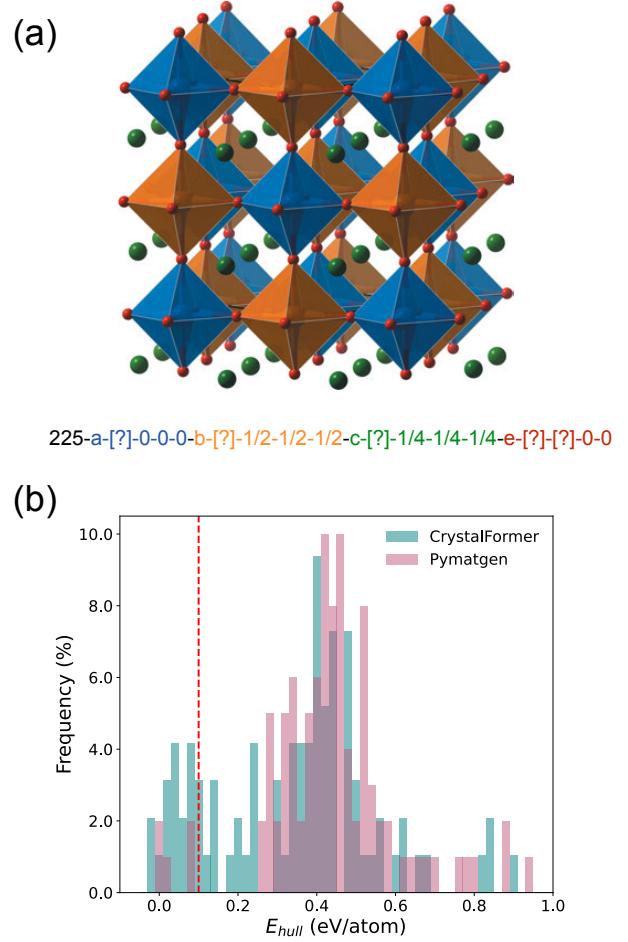


FIG. 8. (a) Double perovskites crystal structure. The crystal string representation of double perovskites with blank spaces for chemical elements and the x-coordinate of the atom resides in the 'e' position. **CrystalFormer** generates crystals with double perovskite structures via sequence infilling. (b) The histograms of energy above the convex hull for the relaxed crystal structures. The dashed line indicates the criterion for selecting candidates for the stable materials.

can be straightforwardly expressed as constraints on the space group and occupied Wyckoff letters [57].

C. Plug-and-play materials design

Finally, we demonstrate **CrystalFormer**'s ability to aid **property-guided exploration** of crystalline materials in a versatile and flexible manner. The trained **CrystalFormer** captures the space group conditioned crystal probability $p(\mathbf{C}|g)$, which we treat as a prior probability for stable crystals. By **combining it with a crystal property prediction model that provides the forward likelihood probability $p(y|\mathbf{C})$** , one can carry out **property-guided materials generation in a plug-and-play manner**. According to Bayes' rule, the posterior for crystals given property y reads

$$p(\mathbf{C}|g, y) \propto p(y|\mathbf{C})p(\mathbf{C}|g). \quad (3)$$

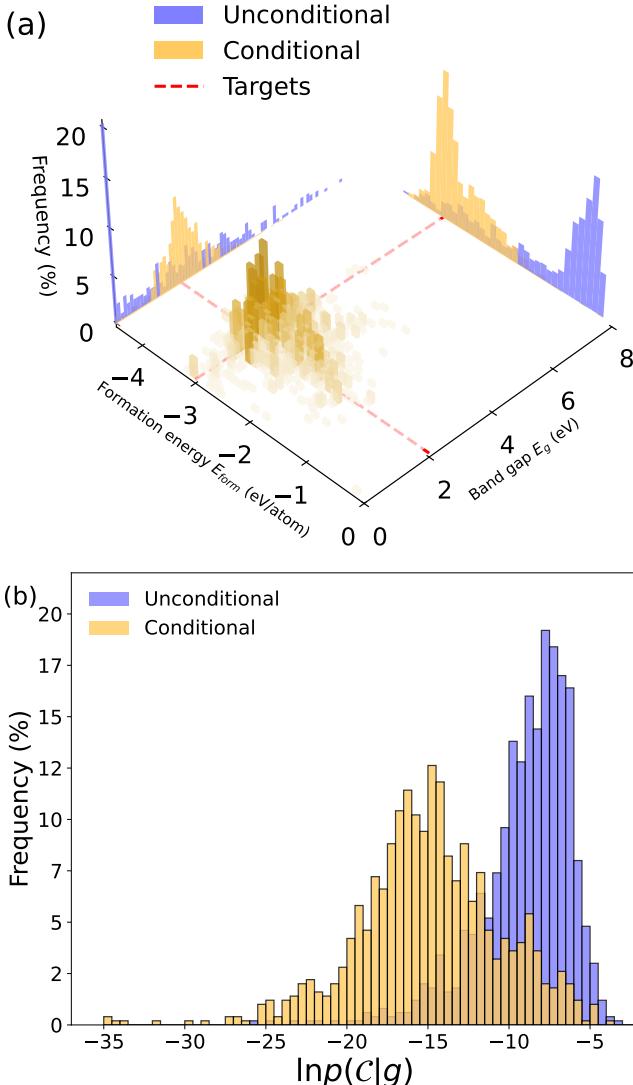


FIG. 9. (a) The histogram of band gap and formation energy for crystal samples generated in the $Fm\bar{3}m$ space group (No. 225). The dashed red lines in the plane indicate target values. The marginals on the side show the shift of the property distributions with respect to the unconditionally generated samples. Note that we scale the 3d histograms for better visualization. (b) The likelihoods of conditioned generated samples compared to the unconditional samples.

By sampling from this posterior distribution, one can generate crystal samples with property guidance. Since the posterior probability Eq. (3) typically does not process autoregressive property with respect to C , we carry out MCMC sampling to sample from the posterior distribution [84]. The plug-and-play feature makes designing crystalline materials in this way particularly appealing because it is possible to apply multiple conditions by simply adding log-likelihoods from multiple predictors. The framework applies to the inverse problem of solving crystal structures based on experimentally observed diffraction spectra equally well [85, 86], where the goal is to

simultaneously optimizing the matching probability to experimental observation and stability of the crystal.

Any property prediction model can be used in conjunction with CrystalFormer for property-guided material generation. We train two crystal property prediction models using the MP-20 dataset and achieve mean absolute error (MAE) of 0.325 eV in band gap and MAE of 0.110 eV/atom in formation energy, which is comparable to the crystal graph convolutional neural network [87] trained on the same dataset [88]. We use the output of these two property prediction models as the forward probability $p(y|C)$ of Eq. (3). More details are in appendix E.

Figure 9(a) demonstrates the controlled generation of materials with target band gap at $E_g = 2$ eV and the formation of energy $E_{\text{form}} = -3$ eV/atom crystals [89, 90]. The conditional probability Eq. (3) contains both the model likelihood and the property regression MAE. Therefore, the generated samples will strike a balance between the two. To draw samples from the conditional probability distribution, we randomly generate a batch of 1000 crystal samples and sweep through the crystal sequence to update the atom species, fractional coordinates, and lattice parameters [91]. Monte Carlo update of chemical species so to achieve the desired properties can be regarded as a systematic data-informed way of carrying out cation-transmutation for materials inverse design [92]. After reaching equilibrium, the histogram of band gap and formation energy is centered around the target values, which is shifted significantly away from the value of unconditionally generated samples. On the other hand, the likelihood shown in Fig. 9(b) indicates these conditional-generated crystals are not typical samples with respect to the unconditional distribution. Nevertheless, they are still probable samples according to the crystal prior given by the CrystalFormer. Given a myriad of materials property prediction models developed over the years and the inconvenience of re-training or fine-tuning the foundational generative model [16], we envision the plug-and-play generation approach demonstrated here to be a scalable way for materials design. We have exposed an interface of CrystalFormer in the code repository [34] for users to plug in arbitrary conditioners for guided materials generation.

V. RELATED WORKS

Crystal generative models have been explored using variational autoencoder [23, 93], generative adversarial networks [12, 94], normalizing flows [95–98], diffusion models [11–16, 21, 22, 26], GFlowNet [24, 25], and autoregressive models [17–20, 99–101]. In these autoregressive models, one either uses atomistic features [17, 99–101] or uses pure text descriptions [18–20]. Nevertheless, with the introduction of specialized tokens for crystals, the boundary between the two is blurred.

The CrystalFormer is most closely related to the autoregressive generative model originally designed for molecules [99–101]. However, instead of predicting the relative distances of atoms, we predict the Wyckoff positions of symmetry-inequivalent atoms in the unit cell. Having the

luxury of the space group symmetry for crystals provides strong hints on where to put the atoms in the unit cell and greatly simplifies the design around spatial symmetries. On the other hand, compared to Refs. [19] which treats text descriptions of crystals using autoregressive language models, we are only dealing with a more concise and essential atomistic representation of crystals, which leads to a smaller model size and faster sampling speed. Fast generation speed is not only a welcoming feature but also will be crucial for further exploration of materials space based on combinations of probabilistic generation and post-selection, Monte Carlo sampling, backtracking, and searching techniques [102]. More importantly, by baking in the space group symmetry in the model rather than learning them as statistical correlation from texts [18, 20], **CrystalFormer** guarantees space group constraints and cherishes the precious data and computing time. In a sense, the present work employs intrinsic mathematical (as opposed to natural) language to incorporate the symmetry principle in the generative modeling of crystals.

As a side remark, the Wyckoff position features have been used in machine learning models for materials property prediction [49, 103]. Incorporating space group information in the encoder-only transformer models may also enhance their property prediction performance [104–106] as suggested by Ref. [107].

VI. OUTLOOK

Precisely controlling the space group in the generative model of crystalline materials not only greatly simplifies the task but also is a highly desired feature for materials discovery and design. **CrystalFormer** integrates exact symmetry principles from math and empirical chemical intuitions from data into one unified framework. Probabilistic generative modeling of crystalline materials using **CrystalFormer** opens the way to many future innovations in materials discovery.

Note that the MP-20 dataset has by no means exhausted all available crystalline material [16, 19]. An obvious future direction is to scale up the model as well as the training dataset, especially curating a dataset with better coverage of space groups. In particular, extending the dataset to include both inorganic and organic crystals [108] may be beneficial as it improves the data coverage of low symmetric space groups. The transformer-based generative model is ready to be scaled up to work with much larger and more diverse training data, in the same fashion as large language models [109]. Given

similar model architectures, the idea of generative pretraining of a foundational model for material generation is appealing. When scaling up the model it will be interesting to note the possible appearance of neural scaling law [110] as it has also been showing up in other contexts of atomistic modeling [111].

The model architecture and sampling strategy are both open to further refinement to better serve the purpose of material discovery. First of all, to better facilitate data efficiency learning and structure phase transitions-related applications, it will be useful to further exploit the Euclidean normalizer [112] and group-subgroup relation [113] in the model architecture or training procedure. Second, it is worth exploring using **CrystalFormer** as the base distribution in the flow model and employs symmetry-persevering transportation to further adjust the atoms coordinates and unit cells [10, 26], which mimics a symmetry-constrained relaxation process [114]. Lastly, it may be worth employing more advanced constrained and guided sequence generation methods [115–118] for more flexible control on the elements, structure, or stoichiometry of generated materials.

Conditioned materials generation depending on properties [16, 21, 22, 101] and experimental measurements [119] are highly desired features of materials generative model. Although it is straightforward to extend **CrystalFormer** (e.g. extend the space group embedding or employ the encoder-decoder transformer architecture [40]) to incorporate these conditions, we are particularly excited about the plug-and-play routine demonstrated in Sec. IV C. Along this line, we envision an ecosystem [120] where the foundational generative model for $p(C|g)$ and more specialized discriminative models $p(y|C)$ are developed separately but brought together via the Bayes rule.

ACKNOWLEDGMENTS

We thank Han Wang, Lin Yao, Linfeng Zhang, Chen Fang, Yanchao Wang, Zhenyu Wang, Qi Yang, Shigang Ou, Xinyang Dong, Wenbing Huang, Quansheng Wu, Wanjin Yin, Xi Dai, Shuang Jia, Hangtian Zhu, Jiangang Guo, and Hongjian Zhao for useful discussions. This project is supported by the National Natural Science Foundation of China under Grants No. T2225018, No. 92270107, No. 12188101, No. T2121001, and No. 12034009 and the Strategic Priority Research Program of Chinese Academy of Sciences under Grants No. XDB05000000 and No. XDB30000000.

-
- [1] S. M. Woodley and R. Catlow, Crystal structure prediction from first principles, *Nature Materials* **7**, 937 (2008).
 - [2] A. R. Oganov, C. J. Pickard, Q. Zhu, and R. J. Needs, Structure prediction drives materials discovery, *Nature Reviews Materials* **4**, 331 (2019).
 - [3] R. Gómez-Bombarelli, J. N. Wei, D. Duvenaud, J. M. Hernández-Lobato, B. Sánchez-Lengeling, D. Sheberla, J. Aguilera-Iparraguirre, T. D. Hirzel, R. P. Adams, and A. Aspuru-Guzik, Automatic chemical design using a data-driven continuous representation of molecules, *ACS Central Science* **4**, 268 (2018), pMID: 29532027.
 - [4] B. Sanchez-Lengeling and A. Aspuru-Guzik, Inverse molecular design using machine learning: Generative models for matter engineering, *Science* **361**, 360 (2018).
 - [5] A. Merchant, S. Batzner, S. S. Schoenholz, M. Aykol, G. Cheon, and E. D. Cubuk, Scaling deep learning for ma-

- terials discovery, *Nature* **624**, 80 (2023).
- [6] C. Chen, D. T. Nguyen, S. J. Lee, N. A. Baker, A. S. Karakoti, L. Lauw, C. Owen, K. T. Mueller, B. A. Bilodeau, V. Murugesan, and M. Troyer, Accelerating computational materials discovery with artificial intelligence and cloud high-performance computing: from large-scale screening to experimental validation, (2024), [arXiv:2401.04070 \[cond-mat.mtrl-sci\]](https://arxiv.org/abs/2401.04070).
- [7] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, *et al.*, Gpt-4 technical report, (2023), [arXiv:2303.08774 \[cs.CL\]](https://arxiv.org/abs/2303.08774).
- [8] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, Zero-shot text-to-image generation, in *Proceedings of the 38th International Conference on Machine Learning*, Proceedings of Machine Learning Research, Vol. 139, edited by M. Meila and T. Zhang (PMLR, 2021) pp. 8821–8831.
- [9] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, High-resolution image synthesis with latent diffusion models, in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2022) pp. 10684–10695.
- [10] J. B. Ingraham, M. Baranov, Z. Costello, K. W. Barber, W. Wang, A. Ismail, V. Frappier, D. M. Lord, C. Ng-Thow-Hing, E. R. Van Vlack, S. Tie, V. Xue, S. C. Cowles, A. Leung, J. V. Rodrigues, C. L. Morales-Perez, A. M. Ayoub, R. Green, K. Puentes, F. Oplinger, N. V. Panwar, F. Obermeyer, A. R. Root, A. L. Beam, F. J. Poelwijk, and G. Grigoryan, Illuminating protein space with a programmable generative model, *Nature* **623**, 1070 (2023).
- [11] T. Xie, X. Fu, O.-E. Ganea, R. Barzilay, and T. Jaakkola, Crystal diffusion variational autoencoder for periodic material generation, (2021), [arXiv:2110.06197 \[cs.LG\]](https://arxiv.org/abs/2110.06197).
- [12] Y. Luo, C. Liu, and S. Ji, Towards symmetry-aware generation of periodic materials, (2023), [arXiv:2307.02707 \[cs.LG\]](https://arxiv.org/abs/2307.02707).
- [13] R. Jiao, W. Huang, P. Lin, J. Han, P. Chen, Y. Lu, and Y. Liu, Crystal structure prediction by joint equivariant diffusion, (2023), [arXiv:2309.04475 \[cond-mat.mtrl-sci\]](https://arxiv.org/abs/2309.04475).
- [14] S. Zheng, J. He, C. Liu, Y. Shi, Z. Lu, W. Feng, F. Ju, J. Wang, J. Zhu, Y. Min, H. Zhang, S. Tang, H. Hao, P. Jin, C. Chen, F. Noé, H. Liu, and T.-Y. Liu, Towards predicting equilibrium distributions for molecular systems with deep learning, (2023), [arXiv:2306.0544 \[physics.chem-ph\]](https://arxiv.org/abs/2306.0544).
- [15] M. Yang, K. Cho, A. Merchant, P. Abbeel, D. Schuurmans, I. Mordatch, and E. D. Cubuk, Scalable diffusion for materials discovery, (2023), [arXiv:2311.09235 \[cs.LG\]](https://arxiv.org/abs/2311.09235).
- [16] C. Zeni, R. Pinsky, D. Zigner, A. Fowler, M. Horton, X. Fu, S. Shysheya, J. Crabbé, L. Sun, J. Smith, *et al.*, Mattergen: a generative model for inorganic materials design, (2023), [arXiv:2312.03687 \[cond-mat.mtrl-sci\]](https://arxiv.org/abs/2312.03687).
- [17] H. Xiao, R. Li, X. Shi, Y. Chen, L. Zhu, X. Chen, and L. Wang, An invertible, invariant crystal representation for inverse design of solid-state materials using generative deep learning, *Nature Communications* **14**, 7027 (2023).
- [18] D. Flam-Shepherd and A. Aspuru-Guzik, Language models can generate molecules, materials, and protein binding sites directly in three dimensions as xyz, cif, and pdb files, (2023), [arXiv:2305.05708 \[cs.LG\]](https://arxiv.org/abs/2305.05708).
- [19] L. M. Antunes, K. T. Butler, and R. Grau-Crespo, Crystal structure generation with autoregressive large language modeling, (2023), [arXiv:2307.04340 \[cond-mat.mtrl-sci\]](https://arxiv.org/abs/2307.04340).
- [20] N. Gruver, A. Sriram, A. Madotto, A. G. Wilson, C. L. Zitnick, and Z. Ulissi, Fine-tuned language models generate stable inorganic materials as text, (2024), [arXiv:2402.04379 \[cs.LG\]](https://arxiv.org/abs/2402.04379).
- [21] X. Luo, Z. Wang, P. Gao, J. Lv, Y. Wang, C. Chen, and Y. Ma, Deep learning generative model for crystal structure prediction, (2024), [arXiv:2403.10846 \[cond-mat.mtrl-sci\]](https://arxiv.org/abs/2403.10846).
- [22] C.-Y. Ye, H.-M. Weng, and Q.-S. Wu, Con-cdvae: A method for the conditional generation of crystal structures, (2024), [arXiv:2403.12478 \[cond-mat.mtrl-sci\]](https://arxiv.org/abs/2403.12478).
- [23] R. Zhu, W. Nong, S. Yamazaki, and K. Hippalgaonkar, Wyccrst: Wyckoff inorganic crystal generator framework, *Matter* doi.org/10.1016/j.matt.2024.05.042 (2024).
- [24] M. AI4Science, A. Hernandez-Garcia, A. Duval, A. Volokhova, Y. Bengio, D. Sharma, P. L. Carrier, M. Koziarski, and V. Schmidt, Crystal-gfn: sampling crystals with desirable properties and constraints, (2023), [arXiv:2310.04925 \[cs.LG\]](https://arxiv.org/abs/2310.04925).
- [25] T. M. Nguyen, S. A. Tawfik, T. Tran, S. Gupta, S. Rana, and S. Venkatesh, *Hierarchical GFlownet for crystal structure generation* (2024).
- [26] R. Jiao, W. Huang, Y. Liu, D. Zhao, and Y. Liu, Space group constrained crystal generation, (2024), [arXiv:2402.03992 \[cs.LG\]](https://arxiv.org/abs/2402.03992).
- [27] A. R. Oganov and C. W. Glass, Crystal structure prediction using ab initio evolutionary techniques: Principles and applications, *The Journal of chemical physics* **124** (2006).
- [28] D. W. Davies, K. T. Butler, A. J. Jackson, A. Morris, J. M. Frost, J. M. Skelton, and A. Walsh, Computational screening of all stoichiometric inorganic materials, *Chem* **1**, 617 (2016).
- [29] M. Glazer, G. Burns, and A. Glazer, *Space Groups for Solid State Scientists* (Elsevier Science, 2012).
- [30] V. S. Urusov and T. N. Nadezhina, Frequency distribution and selection of space groups in inorganic crystal chemistry, *Journal of Structural Chemistry* **50**, 22 (2009).
- [31] R. E. Marsh, P1 or P1? Or something else?, *Acta Crystallographica Section B* **55**, 931 (1999).
- [32] A. K. Cheetham and R. Seshadri, Artificial intelligence driving materials discovery? perspective on the article: Scaling deep learning for materials discovery, *Chemistry of Materials* **36**, 3490 (2024).
- [33] W. Hornfeck, On the combinatorics of crystal structures: number of wyckoff sequences of given length, *Acta Crystallographica Section A: Foundations and Advances* **78**, 149 (2022).
- [34] See <https://github.com/deepmodeling/CrystalFormer> for code and model checkpoint.
- [35] Y. Zhao, Y. Cui, Z. Xiong, J. Jin, Z. Liu, R. Dong, and J. Hu, Machine learning-based prediction of crystal systems and space groups from inorganic materials compositions, *ACS Omega* **5**, 3596 (2020), pMID: 32118175.
- [36] H. Liang, V. Stanev, A. G. Kusne, and I. Takeuchi, Cryspnet: Crystal structure predictions via neural networks, *Phys. Rev. Mater.* **4**, 123802 (2020).
- [37] D.-Y. Wang, H.-F. Lv, and X.-J. Wu, Crystallographic groups prediction from chemical composition via deep learning, *Chinese Journal of Chemical Physics* **36**, 66 (2023).
- [38] V. Venkatraman and P. A. Carvalho, Accurate space-group prediction from composition, *Journal of Applied Crystallography* (2024).
- [39] T. Hahn, U. Shmueli, and J. W. Arthur, *International tables for crystallography*, Vol. 1 (Reidel Dordrecht, 1983).
- [40] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, Attention is all you need, in *Advances in Neural Information Processing Systems*, Vol. 30, edited by I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Curran Associates, Inc., 2017).
- [41] P. Wirnsberger, A. J. Ballard, G. Papamakarios, S. Abercrom-

- bie, S. Racanière, A. Pritzel, D. J. Rezende, and C. Blundell, Targeted free energy estimation via learned mappings, *Journal of Chemical Physics* **153**, 144112 (2020).
- [42] We follow the convention that capital letters appear *after* lower case letters. This handles the edge case of the *Pmmm* space group (No. 47) whose Wyckoff letters used up 26 lowercase letters and reached 'A' for the generic position. In addition, we use the letter 'X' to indicate the Wyckoff position of padding atoms that appear at the end of the sequence, see e.g. w_4 of Fig. 1(b).
- [43] OpenAI, [Using logit bias to alter token probability with the openai api](#), OpenAI Help Center, retrieved March 14, 2024.
- [44] H. Xie, L. Zhang, and L. Wang, m^* of two-dimensional electron gas: A neural canonical transformation study, *SciPost Phys.* **14**, 154 (2023).
- [45] E. Parthé and L. M. Gelato, The standardization of inorganic crystal-structure data, *Acta Crystallographica Section A* **40**, 169 (1984).
- [46] G. Hautier, C. Fischer, V. Ehrlacher, A. Jain, and G. Ceder, Data mined ionic substitutions for the discovery of new compounds, *Inorganic chemistry* **50**, 656 (2011).
- [47] H. Glawe, A. Sanna, E. K. U. Gross, and M. A. L. Marques, The optimal one dimensional periodic table: a modified petti-for chemical scale from data mining, *New Journal of Physics* **18**, 093011 (2016).
- [48] H.-C. Wang, S. Botti, and M. A. L. Marques, Predicting stable crystalline compounds using chemical similarity, *npj Computational Materials* **7**, 12 (2021).
- [49] A. Jain and T. Bligaard, Atomic-position independent descriptor for machine learning of material properties, *Phys. Rev. B* **98**, 214112 (2018).
- [50] Q. Zhou, P. Tang, S. Liu, J. Pan, Q. Yan, and S.-C. Zhang, Learning atoms for materials discovery, *Proceedings of the National Academy of Sciences* **115**, E6411 (2018).
- [51] C. Chen, W. Ye, Y. Zuo, C. Zheng, and S. P. Ong, Graph networks as a universal machine learning framework for molecules and crystals, *Chemistry of Materials* **31**, 3564 (2019).
- [52] D. Zhang, H. Bi, F.-Z. Dai, W. Jiang, L. Zhang, and H. Wang, Dpa-1: Pretraining of attention-based deep potential model for molecular simulation, (2022), [arXiv:2208.08236 \[physics.chem-ph\]](#).
- [53] A. Y.-T. Wang, M. S. Mahmoud, M. Czasny, and A. Gurlo, Crabnet for explainable deep learning in materials science: bridging the gap between academia and industry, *Integrating Materials and Manufacturing Innovation* **11**, 41 (2022).
- [54] E. Gazzarrini, R. K. Cersonsky, M. Bercx, C. S. Adorf, and N. Marzari, The rule of four: anomalous distributions in the stoichiometries of inorganic compounds, *Npj Computational Materials* **10**, 73 (2024).
- [55] R. Palgrave, An explanation for the rule of four in inorganic materials [10.26434/chemrxiv-2024-sxqwh](#) (2024).
- [56] S. Alvarez, Polyhedra in (inorganic) chemistry, *Dalton Transactions* , 2209 (2005).
- [57] N. Regnault, Y. Xu, M.-R. Li, D.-S. Ma, M. Jovanovic, A. Yazdani, S. S. Parkin, C. Felser, L. M. Schoop, N. P. Ong, *et al.*, Catalogue of flat-band stoichiometric materials, *Nature* **603**, 824 (2022).
- [58] L. Pauling, The principles determining the structure of complex ionic crystals, *Journal of the American Chemical Society* **51**, 1010 (1929).
- [59] V. Goldschmidt, Crystal structure and chemical constitution, *Transactions of the Faraday Society* **25**, 253 (1929).
- [60] D. Pettifor, Structure maps for pseudobinary and ternary phases, *Materials Science and Technology* **4**, 675 (1988).
- [61] C. C. Fischer, K. J. Tibbetts, D. Morgan, and G. Ceder, Predicting crystal structure by merging data mining with quantum mechanics, *Nature Materials* **5**, 641 (2006).
- [62] Z. Allahyari and A. R. Oganov, Coevolutionary search for optimal materials in the space of all possible compounds, *npj Computational Materials* **6**, 55 (2020).
- [63] J. Maddox, Crystals from first principles, *Nature* **335**, 201 (1988).
- [64] D. J. Wales, Symmetry, near-symmetry and energetics, *Chemical physics letters* **285**, 330 (1998).
- [65] P. Avery and E. Zurek, Randspg: An open-source program for generating atomistic crystal structures with specific space-groups, *Computer Physics Communications* **213**, 208 (2017).
- [66] S. Fredericks, K. Parrish, D. Sayre, and Q. Zhu, Pyxtal: A python library for crystal structure generation and symmetry analysis, *Computer Physics Communications* **261**, 107810 (2021).
- [67] Y. Wang, J. Lv, L. Zhu, and Y. Ma, Crystal structure prediction via particle-swarm optimization, *Phys. Rev. B* **82**, 094116 (2010).
- [68] C. J. Pickard and R. J. Needs, Ab initio random structure searching, *Journal of Physics: Condensed Matter* **23**, 053201 (2011).
- [69] A. O. Lyakhov, A. R. Oganov, H. T. Stokes, and Q. Zhu, New developments in evolutionary structure prediction algorithm uspex, *Computer Physics Communications* **184**, 1172 (2013).
- [70] Z. Falls, P. Avery, X. Wang, K. P. Hilleke, and E. Zurek, The xtalopt evolutionary algorithm for crystal structure prediction, *The Journal of Physical Chemistry C* **125**, 1601 (2020).
- [71] G. Cheng, X.-G. Gong, and W.-J. Yin, Crystal structure prediction by combining graph network and optimization algorithm, *Nature Communications* **13**, 1492 (2022).
- [72] H.-C. Wang, J. Schmidt, M. A. L. Marques, L. Wirtz, and A. H. Romero, Symmetry-based computational search for novel binary and ternary 2d materials, *2D Materials* **10**, 035007 (2023).
- [73] Q. Zhang, A. Choudhury, and A. Chernatynskiy, A symmetry-oriented crystal structure prediction method for crystals with rigid bodies, (2024), [arXiv:2407.21337 \[cond-mat.mtrl-sci\]](#).
- [74] W. Sun, S. T. Dacek, S. P. Ong, G. Hautier, A. Jain, W. D. Richards, A. C. Gamst, K. A. Persson, and G. Ceder, The thermodynamic scale of inorganic crystalline metastability, *Science Advances* **2**, e1600225 (2016).
- [75] C. Chen and S. P. Ong, A universal graph deep learning interatomic potential for the periodic table, *Nature Computational Science* **2**, 718 (2022).
- [76] D. Zhang, X. Liu, X. Zhang, C. Zhang, C. Cai, H. Bi, Y. Du, X. Qin, J. Huang, B. Li, Y. Shan, J. Zeng, Y. Zhang, S. Liu, Y. Li, J. Chang, X. Wang, S. Zhou, J. Liu, X. Luo, Z. Wang, W. Jiang, J. Wu, Y. Yang, J. Yang, M. Yang, F.-Q. Gong, L. Zhang, M. Shi, F.-Z. Dai, D. M. York, S. Liu, T. Zhu, Z. Zhong, J. Lv, J. Cheng, W. Jia, M. Chen, G. Ke, W. E, L. Zhang, and H. Wang, Dpa-2: Towards a universal large atomic model for molecular and material simulation, (2023), [arXiv:2312.15492 \[physics.chem-ph\]](#).
- [77] I. Batatia, P. Benner, Y. Chiang, A. M. Elena, D. P. Kovács, J. Riebesell, X. R. Advincula, M. Asta, M. Avaylon, W. J. Baldwin, F. Berger, N. Bernstein, A. Bhowmik, S. M. Blau, V. Cărare, J. P. Darby, S. De, F. D. Pia, V. L. Deringer, R. Elijošius, Z. El-Machachi, F. Falcioni, E. Fako, A. C. Ferrari, A. Genreith-Schriever, J. George, R. E. A. Goodall, C. P. Grey, P. Grigorev, S. Han, W. Handley, H. H. Heenen, K. Hermansson, C. Holm, J. Jaafar, S. Hofmann, K. S. Jakob,

- H. Jung, V. Kapil, A. D. Kaplan, N. Karimitari, J. R. Kermode, N. Kroupa, J. Kullgren, M. C. Kuner, D. Kuryla, G. Liepuoniute, J. T. Margraf, I.-B. Magdäu, A. Michaelides, J. H. Moore, A. A. Naik, S. P. Niblett, S. W. Norwood, N. O'Neill, C. Ortner, K. A. Persson, K. Reuter, A. S. Rosen, L. L. Schaaf, C. Schran, B. X. Shi, E. Sivonxay, T. K. Stenczel, V. Svahn, C. Sutton, T. D. Swinburne, J. Tilly, C. van der Oord, E. Varga-Umbrich, T. Vegge, M. Vondrák, Y. Wang, W. C. Witt, F. Zills, and G. Csányi, A foundation model for atomistic materials chemistry, (2024), [arXiv:2401.00096 \[physics.chem-ph\]](https://arxiv.org/abs/2401.00096).
- [78] S. Vasala and M. Karppinen, A2b'b"o6 perovskites: a review, *Progress in solid state chemistry* **43**, 1 (2015).
- [79] Y. Wang, B. Baldassarri, J. Shen, J. He, and C. Wolverton, Landscape of thermodynamic stabilities of a2bb'o6 compounds, *Chemistry of Materials* (2024).
- [80] N. Miao, H. Zhou, L. Mou, R. Yan, and L. Li, CGMH: Constrained sentence generation by metropolis-hastings sampling, in *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33 (2019) pp. 6834–6842.
- [81] S. P. Ong, W. D. Richards, A. Jain, G. Hautier, M. Kocher, S. Cholia, D. Gunter, V. L. Chevrier, K. A. Persson, and G. Ceder, Python materials genomics (pymatgen): A robust, open-source python library for materials analysis, *Computational Materials Science* **68**, 314 (2013).
- [82] I.-H. Chu, S. Roychowdhury, D. Han, A. Jain, and S. P. Ong, Predicting the volumes of crystals, *Computational Materials Science* **146**, 184 (2018).
- [83] R. Okabe, M. Cheng, A. Chottrattanapituk, N. T. Hung, X. Fu, B. Han, Y. Wang, W. Xie, R. J. Cava, T. S. Jaakkola, Y. Cheng, and M. Li, Structural constraint integration in generative model for discovery of quantum material candidates, (2024), [arXiv:2407.04557 \[cond-mat.mtrl-sci\]](https://arxiv.org/abs/2407.04557).
- [84] R. Verkuil, O. Kabeli, Y. Du, B. I. Wicky, L. F. Milles, J. Dauparas, D. Baker, S. Ovchinnikov, T. Sercu, and A. Rives, Language models generalize beyond natural proteins, *BioRxiv* , 2022 (2022).
- [85] B. Meredig and C. Wolverton, A hybrid computational-experimental approach for automated crystal structure solution, *Nature materials* **12**, 123 (2013).
- [86] A. S. Parackal, R. E. A. Goodall, F. A. Faber, and R. Armiento, Identifying crystal structures beyond known prototypes from x-ray powder diffraction spectra, (2023), [arXiv:2309.16454 \[cond-mat.mtrl-sci\]](https://arxiv.org/abs/2309.16454).
- [87] T. Xie and J. C. Grossman, Crystal graph convolutional neural networks for an accurate and interpretable prediction of material properties, *Phys. Rev. Lett.* **120**, 145301 (2018).
- [88] C.-Y. Ye, H.-M. Weng, and Q.-S. Wu, Con-cdvae: A method for the conditional generation of crystal structures, *Computational Materials Today* **1**, 100003 (2024).
- [89] A. Franceschetti and A. Zunger, The inverse band-structure problem of finding an atomic configuration with given electronic properties, *Nature* **402**, 60 (1999).
- [90] G. Cheng, X.-G. Gong, and W.-J. Yin, Global optimization in the discrete and variable-dimension conformational space: The case of crystal with the strongest atomic cohesion, (2023), [arXiv:2302.13537 \[cond-mat.mtrl-sci\]](https://arxiv.org/abs/2302.13537).
- [91] For simplicity of the Wyckoff sequence are kept unchanged in the MCMC sampling.
- [92] X.-G. Zhao, J.-H. Yang, Y. Fu, D. Yang, Q. Xu, L. Yu, S.-H. Wei, and L. Zhang, Design of lead-free inorganic halide perovskites for solar cells via cation-transmutation, *Journal of the American Chemical Society* **139**, 2630 (2017).
- [93] Z. Ren, S. I. P. Tian, J. Noh, F. Oviedo, G. Xing, J. Li, Q. Liang, R. Zhu, A. G. Aberle, S. Sun, *et al.*, An invertible crystallographic representation for general inverse design of inorganic crystals with targeted properties, *Matter* **5**, 314 (2022).
- [94] Y. Zhao, M. Al-Fahdi, M. Hu, E. M. Siriwardane, Y. Song, A. Nasiri, and J. Hu, High-throughput discovery of novel cubic crystal materials using deep generative neural networks, *Advanced Science* **8**, 2100566 (2021).
- [95] R. Ahmad and W. Cai, Free energy calculation of crystalline solids using normalizing flows, *Modelling and Simulation in Materials Science and Engineering* **30**, 065007 (2022).
- [96] P. Wirnsberger, G. Papamakarios, B. Ibarz, S. Racanière, A. J. Ballard, A. Pritzel, and C. Blundell, Normalizing flows for atomic solids, *Machine Learning: Science and Technology* **3**, 025009 (2022).
- [97] J. Köhler, M. Invernizzi, P. De Haan, and F. Noé, Rigid body flows for sampling molecular crystal structures, (2023), [arXiv:2301.11355 \[cs.LG\]](https://arxiv.org/abs/2301.11355).
- [98] B. K. Miller, R. T. Chen, A. Sriram, and B. M. Wood, Flowmm: Generating materials with riemannian flow matching, (2024), [arXiv:2406.04713 \[cs.CL\]](https://arxiv.org/abs/2406.04713).
- [99] N. W. Gebauer, M. Gastegger, and K. T. Schütt, Generating equilibrium molecules with deep neural networks, (2018), [arXiv:1810.11347 \[stat.ML\]](https://arxiv.org/abs/1810.11347).
- [100] N. Gebauer, M. Gastegger, and K. Schütt, Symmetry-adapted generation of 3d point sets for the targeted discovery of molecules, in *Advances in Neural Information Processing Systems*, Vol. 32 (Curran Associates, Inc., 2019).
- [101] N. W. A. Gebauer, M. Gastegger, S. S. P. Hessmann, K.-R. Müller, and K. T. Schütt, Inverse design of 3d molecular structures with conditional generative neural networks, *Nature Communications* **13**, 973 (2022).
- [102] S. Yao, D. Yu, J. Zhao, I. Shafran, T. L. Griffiths, Y. Cao, and K. Narasimhan, Tree of thoughts: Deliberate problem solving with large language models, (2023), [arXiv:2305.10601 \[cs.CL\]](https://arxiv.org/abs/2305.10601).
- [103] R. E. A. Goodall, A. S. Parackal, F. A. Faber, R. Armiento, and A. A. Lee, Rapid discovery of stable materials by coordinate-free coarse graining, *Science Advances* **8**, eabn4117 (2022).
- [104] K. Yan, Y. Liu, Y. Lin, and S. Ji, Periodic graph transformers for crystal material property prediction, in *Advances in Neural Information Processing Systems*, Vol. 35, edited by S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (Curran Associates, Inc., 2022) pp. 15066–15080.
- [105] T. Taniai, R. Igarashi, Y. Suzuki, N. Chiba, K. Saito, Y. Ushiku, and K. Ono, Crystalfomer: Infinitely connected attention for periodic structure encoding, (2024), [arXiv:2403.11686 \[cs.LG\]](https://arxiv.org/abs/2403.11686).
- [106] H. Xu, D. Qian, and J. Wang, Predicting many properties of crystals by a single deep learning model, (2024), [arXiv:2405.18944 \[cond-mat.mtrl-sci\]](https://arxiv.org/abs/2405.18944).
- [107] A. N. Rubungo, C. Arnold, B. P. Rand, and A. B. Ding, Llm-prop: Predicting physical and electronic properties of crystalline solids from their text descriptions, (2023), [arXiv:2310.14029 \[cs.CL\]](https://arxiv.org/abs/2310.14029).
- [108] C. R. Groom, I. J. Bruno, M. P. Lightfoot, and S. C. Ward, The cambridge structural database, *Structural Science* **72**, 171 (2016).
- [109] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, *et al.*, Language models are few-shot learners, *Advances in neural information processing systems* **33**, 1877 (2020).
- [110] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei, Scaling laws for neural language models, (2020),

- arXiv:2001.08361 [cs.LG].
- [111] N. C. Frey, R. Soklaski, S. Axelrod, S. Samsi, R. Gómez-Bombarelli, C. W. Coley, and V. Gadepally, Neural scaling of deep chemical models, *Nature Machine Intelligence* **5**, 1297 (2023).
- [112] U. Müller, *Symmetry Relationships between Crystal Structures: Applications of Crystallographic Group Theory in Crystal Chemistry*, International Union of Crystallography Texts on Crystallography (OUP Oxford, 2013).
- [113] H. T. Stokes and D. M. Hatch, *Isotropy Subgroups of the 230 Crystallographic Space Groups* (WORLD SCIENTIFIC, 1989).
- [114] S. Cox and A. D. White, Symmetric molecular dynamics, *Journal of Chemical Theory and Computation* **18**, 4077 (2022).
- [115] S. Dathathri, A. Madotto, J. Lan, J. Hung, E. Frank, P. Molino, J. Yosinski, and R. Liu, Plug and play language models: A simple approach to controlled text generation, (2019), arXiv:1912.02164 [cs.CL].
- [116] M. Zhang, N. Jiang, L. Li, and Y. Xue, Language generation via combinatorial constraint satisfaction: A tree search enhanced monte-carlo approach, (2020), arXiv:2011.12334 [cs.LG].
- [117] L. Qin, S. Welleck, D. Khashabi, and Y. Choi, Cold decoding: Energy-based constrained text generation with langevin dynamics, *Advances in Neural Information Processing Systems* **35**, 9538 (2022).
- [118] A. K. Lew, T. Zhi-Xuan, G. Grand, and V. K. Mansinghka, Sequential monte carlo steering of large language models using probabilistic programs, (2023), arXiv:2306.03081 [cs.AI].
- [119] Q. Lai, L. Yao, Z. Gao, S. Liu, H. Wang, S. Lu, D. He, L. Wang, C. Wang, and G. Ke, End-to-end crystal structure prediction from powder x-ray diffraction, (2024), arXiv:2401.03862 [physics.chem-ph].
- [120] E. S. Raymond, *The Cathedral and the Bazaar: Musings on Linux and Open Source by an Accidental Revolutionary* (O'Reilly Media, Sebastopol, CA, 1999).
- [121] C. J. Court, B. Yildirim, A. Jain, and J. M. Cole, 3-d inorganic crystal structure generation and property prediction via representation learning, *Journal of Chemical Information and Modeling* **60**, 4518 (2020), pMID: 32866381.
- [122] 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, *Journal of Open Source Software* **4**, 1361 (2019).
- [123] <https://github.com/sparks-baird/matbench-genmetrics>.
- [124] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning* (MIT Press, 2016) <http://www.deeplearningbook.org>.
- [125] J. P. Perdew, K. Burke, and M. Ernzerhof, Generalized Gradient Approximation Made Simple, *Phys. Rev. Lett.* **77**, 3865 (1996).
- [126] P. E. Blöchl, Projector augmented-wave method, *Phys. Rev. B* **50**, 17953 (1994).
- [127] G. Kresse and J. Furthmüller, Efficient iterative schemes for ab initio total-energy calculations using a plane-wave basis set, *Phys. Rev. B* **54**, 11169 (1996).
- [128] M. Chen, A. Radford, R. Child, J. Wu, H. Jun, D. Luan, and I. Sutskever, Generative pretraining from pixels, in *Proceedings of the 37th International Conference on Machine Learning*, Proceedings of Machine Learning Research, Vol. 119, edited by H. D. III and A. Singh (PMLR, 2020) pp. 1691–1703.

Appendix A: More details of `CrystalFormer`

To recap, the space group information plays a key role in the architecture, training, and sampling of `CrystalFormer`. First of all, the vector embedding of space group number g controls all subsequent outputs of the transformer corresponding to the Wyckoff letters, chemical species, fractional coordinates, and lattice parameters. Second, the information of the space group and Wyckoff letter are used to select active components in the fractional coordinates and lattice parameters in the loss function during training. Lastly, the space group determines the concrete meaning of Wyckoff letters in terms of multiplicities and fractional coordinates, which are used to place the right number of atoms precisely in the unit cell during sampling.

1. Model architectures

Algorithm 1 summarized the model architecture of `CrystalFormer`. Training the model for 3,800 epochs with the hyperparameters shown in Table S1 takes about 13 hours on a single A100 GPU.

Algorithm 1 The `CrystalFormer` architecture

Input: Space group number g , Wyckoff letters $\mathbf{W} = [w_i]$, multiplicity of Wyckoff positions $\mathbf{M} = [m_i]$, chemical elements $\mathbf{A} = [a_i]$, fractional coordinates $\mathbf{X} = [(x_i, y_i, z_i)]$ of each atom in the unit cell.

Output: Parameters for the conditional probability of Wyckoff letters ω_i , chemical element α_i , and fractional coordinates χ_i, ν_i, ζ_i of atoms and the lattice parameters ℓ .

```

1:  $\omega_1 = \text{Net}(g)$                                  $\triangleright$  the logit of the first Wyckoff position is implemented as a standalone neural network.
2: # prepare input features
3:  $\mathbf{h}_W = [\text{Embed}(g), \text{Embed}(w_i), m_i]$ .
4:  $\mathbf{h}_A = [\text{Embed}(g), \text{Embed}(a_i)]$ .
5:  $\mathbf{h}_X = [\text{Embed}(g), \cos(2\pi x_i), \sin(2\pi x_i), \dots, \cos(2\pi x_i N_f), \sin(2\pi x_i N_f)]$ 
6:  $\mathbf{h}_Y = \dots$ 
7:  $\mathbf{h}_Z = \dots$ 
8: # concatenate along particle dimension
9:  $\mathbf{h} = \text{Concatenate}(\mathbf{h}_W, \mathbf{h}_A, \mathbf{h}_X, \mathbf{h}_Y, \mathbf{h}_Z)$ 
10: Project  $\mathbf{h}$  feature size to  $d_{\text{model}}$  and add position embedding
11:  $\mathbf{h} = \text{MaskedTransformer}(\mathbf{h})$ 
12: Project  $\mathbf{h}$  feature size to desired dimensions
13: # split along particle dimension
14:  $\omega_i, \alpha_i, \chi_i, \nu_i, \zeta_i, \ell = \text{Split}(\mathbf{h})$ 
15: Mask  $\omega_i$  to ensure the Wyckoff letters are valid for the given space group  $g$  and appear in alphabetical order.
16: return  $[\omega_1, \alpha_1, \chi_1, \nu_1, \zeta_1, \omega_2, \alpha_2, \chi_2, \dots, \ell]$ 

```

TABLE S1. A table of hyperparameters used in this work.

Hyperparameters	Value	Remarks
The length of atom sequence including the padding atoms	21	
Number of chemical species	119	'H' to 'Og', plus padding atom
Number of possible Wyckoff letters	28	'a-z'+'A', plus padding atom
Number of modes in von-Mises mixture distribution K_x	16	
Number of modes in lattice Gaussian mixture distribution K_l	16	
Hidden layer dimension for the composite type of the first atom	256	
Transformer number of layers	16	
Transformer number of heads	16	
Transformer key size	64	
Transformer model size d_{model}	32	
Embedding dimension of discrete input	32	
Number of Fourier frequency N_f	5	
Learning rate	0.0001	
Learning rate decay	0.0	
Weight decay	0.0	
Clip grad	1.0	
Batch Size	100	
Optimizer	Adam	
Dropout rate	0.5	
Total number of parameters: 4840295		

2. Sampling algorithm

Algorithm 2 summarizes the autoregressive sampling method of CrystalFormer. It takes 520 seconds to generate a batch size 13,000 crystal samples on a single A100 GPU, which translates to a generation speed of 40 milliseconds per sample.

Algorithm 2 Autoregressive sampling of CrystalFormer

Input: space group number g , a list of chemical elements `element_list`, length n of the atom sequence, sampling temperature T
Output: Wyckoff letters W , chemical species A , fractional coordinates X of atoms, and lattice parameters L of the unit cell.

```

1: Initialize  $W = [\emptyset]$ ,  $A = [\emptyset]$ ,  $X = [\emptyset]$ 
2: for  $i = 1 \dots, n$  do
3:   # sample Wyckoff letter  $w$ 
4:   Get the last  $\omega$  from CrystalFormer( $g, W, A, X$ )
5:    $w \sim \text{Categorical}(\omega)^{1/T}$ 
6:    $W[i] = w$ 
7:   # sample atom species  $a$ 
8:   Get the last  $\alpha$  from CrystalFormer( $g, W, A, X$ )
9:   Mask the logits in  $\alpha$  according to element_list
10:   $a \sim \text{Categorical}(\alpha)^{1/T}$ 
11:   $A[i] = a$ 
12:  # sample fractional coordinate  $x$ 
13:  Get the last  $\chi$  from CrystalFormer( $g, W, A, X$ )
14:   $x \sim \text{vonMisesMix}(\chi)^{1/T}$ 
15:  Project  $x$  to Wyckoff positions according to the Wyckoff letter  $w$ 
16:  update  $X$  with  $x$ 
17:  # sample fractional coordinate  $y$ 
18:  Get the last  $\nu$  from CrystalFormer( $g, W, A, X$ )
19:  ...
20:  update  $X$  with  $y$ 
21:  # sample fractional coordinate  $z$ 
22:  Get the last  $\zeta$  from CrystalFormer( $g, W, A, X$ )
23:  ...
24:  update  $X$  with  $z$ 
25: end for
26: # sample  $L$ 
27: Get  $\ell$  from CrystalFormer( $g, W, A, X$ )
28:  $L \sim \text{GaussianMix}(\ell)^{1/T}$ 
29: Symmetrize  $L$  according to space group  $g$ 
30: return  $W, A, X, L$ 

```

Algorithm 3 summarizes the Markov chain Monte Carlo sampling of `CrystalFormer`. It is used in the structure constrained generation of crystals in Sec. IV B. The property-guided materials design discussed in Sec. IV C employs a similar sampling strategy for the posterior probability distribution Eq. (3).

Algorithm 3 Markov chain Monte Carlo sampling of `CrystalFormer`

Input: Space group number g , Wyckoff letters W , chemical species A , fractional coordinates X of atoms, and lattice parameters L of the unit cell, length n of the atom sequence, a list of chemical elements `element_list`, step size ϵ ,

Output: Wyckoff letters W , chemical species A , fractional coordinates X of atoms, and lattice parameters L of the unit cell.

```

1:  $C = (W, A, X, L)$ 
2: for  $i = 1 \dots, \text{steps}$  do
3:   for  $j = 1 \dots, n$  do
4:     # update element  $a$ 
5:      $a' \sim \text{element\_list}$ 
6:     # update coordinate  $x$ 
7:      $\Delta x \sim \text{vonMises}$ 
8:     Mask the  $x$  according to the space group  $g$  and Wyckoff letter  $w$ 
9:     # update  $C$ 
10:    Propose an update  $C \rightarrow C'$  with  $a \rightarrow a'$  and  $x \rightarrow x + \epsilon \Delta x$ 
11:    Update according to the Metropolis acceptance probability  $\min\left[1, \frac{p(C'|g)}{p(C|g)}\right]$ 
12:   end for
13: end for
14: # update lattice  $L$ 
15: Get  $\ell$  from CrystalFormer( $g, W, A, X$ )
16:  $L \sim \text{GaussianMix}(\ell)^{1/T}$ 
17: Symmetrize the  $L$  according to the space group  $g$ 
18: return  $W, A, X, L$ 

```

TABLE S2. Validity of generated crystal structure for representative space groups. Training samples count the number of samples in the training set.

Space group	Crystal system	Training samples	Validity (%) ↑	
			Struc.	Comp.
2	Triclinic	676	83.10	83.0
12	Monoclinic	1273	87.70	81.80
62	Orthorhombic	1187	95.50	87.20
139	Tetragonal	1233	97.70	83.40
166	Trigonal	1076	98.50	85.0
194	Hexagonal	1129	99.40	89.90
225	Cubic	3960	99.60	93.50
1	Triclinic	27136	91.40	80.20
Autoregressive models				
PGSchNet [100]			99.65	75.96
LM-CH (character-level tokenization) [18]			84.81	83.55
LM-AC (atom coordinate-level tokenization) [18]			95.81	88.87
Crystal-text-LLM [20]			96.5	86.3
Diffusion models				
CDVAE [11]			100.0	86.70
DiffCSP [13]			100.0	83.25
DiffCSP++ [26]			99.94	85.12
UniMat-Large [15]			97.2	89.4

Appendix B: Validity and novelty of generated samples

Figure S2 illustrates the structure and compositional validity of generated samples across all 230 space groups. Following the Ref. [121], a structure meets the validity criteria if the shortest atomic distance exceeds 0.5 Å, a lenient standard. Composition validity requires charge neutrality as computed by SMACT [122]. This is, however, an overly stringent criterion since the composition validity of the training set is only around 90% by this measure [11]. Note that the CrystalFormer can generate reasonable samples even for those space groups without any training data. This because the model can exploit knowledge learned from other space groups to place suitable atoms in the Wyckoff positions due to weight sharing. Moreover, since the sampling process makes use of Wyckoff position table. The three dimensional coordinates of atoms are not completely random even for unseen space groups.

Table S2 reports the validity of generated samples for selected space groups in each of the seven crystal systems. To ensure the numbers are representative, we chose the space group to be the one with the most training data for each crystal system. One sees that the model performs better for more symmetric space groups. This is a nice feature that complements existing crystal generative models, which mostly have difficulties in generating highly symmetric structures. As a reference, we also list the validity of the generated samples suppose one treats the crystals as if they are all in the *P1* space group (No. 1) with only translational symmetry. One sees the structure validity scores of *P1* space group improve compared to the one shown in Figure S2 due to increased training samples.

The second part of Table S2 shows reference results in the literature for the same validity test. In principle, the performance of the present model should fall back to the language model approaches [18, 20]. The remaining gap may be due to details such as including the header line in the crystallographic information files (CIF), specific sampling strategy of language models, or the additional post-selection of samples [20]. In the table, the DiffCSP++ [26] is the only alternative model that exploits the space group symmetry in a rigorous manner similar to ours. However, the DiffCSP++ does not predict Wyckoff positions in the generation process. Instead, one needs to search for template structures in the training set for generation, which may limit its generality. Besides works listed in Table S2 that reported validity scores in comparable settings, Ref. [19] has conditioned the generation of CIF on the space group symbols in a language model setting. Ref. [16] considered space group conditioned crystal generation using a fine-tuned generative model with space group labels. Neither approach provides exact constraints on the space group, which could yield problematic structures for large systems and highly symmetric space groups. Ref. [23] considered generating symmetric crystals in their Wyckoff representation. However, the model does not consider fractional coordinates and

lattice parameters, so it requires a subsequent computational search to completely determine the crystal structure.

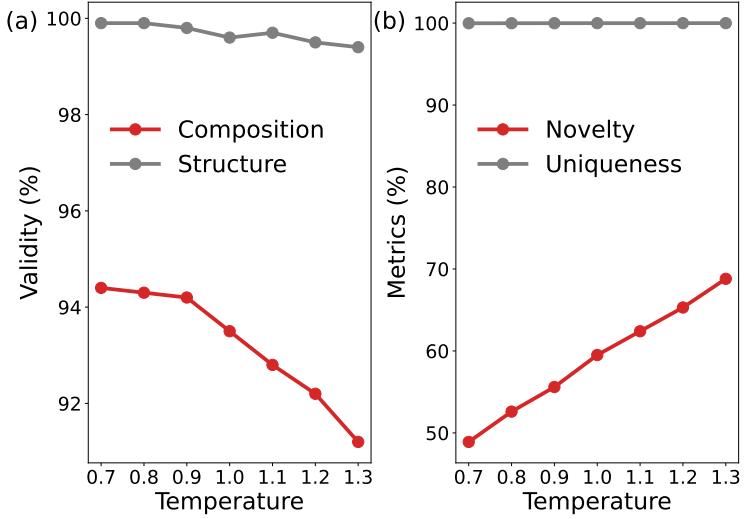


FIG. S1. (a) Structure and composition validity, (b) novelty and uniqueness of generated samples evaluated according to [11] for the $Fm\bar{3}m$ space group.

Figure S1(a) shows the validity of generated samples as a function of sampling temperature. One clearly sees that reduced temperature $T < 1$ increases the validity of samples at the cost of reducing the diversity [20]. Figure S1(b) shows the novelty and uniqueness evaluated on 1000 generated samples with temperature. Novelty quantifies the proportion of new structures in the generated samples that were unseen in the training dataset. Uniqueness represents the percentage of distinct, non-redundant structures among the generated samples [123]. One sees that across different temperatures the uniqueness remains high, indicating the model does not collapse to a mode that produces duplicated samples. On the other hand, the model produces close to 70% novel material as temperature increases, which nicely demonstrates modal covering behavior of the maximum likelihood estimation training [124]. Having a model distribution broader than the span of the dataset is crucial for material discovery.

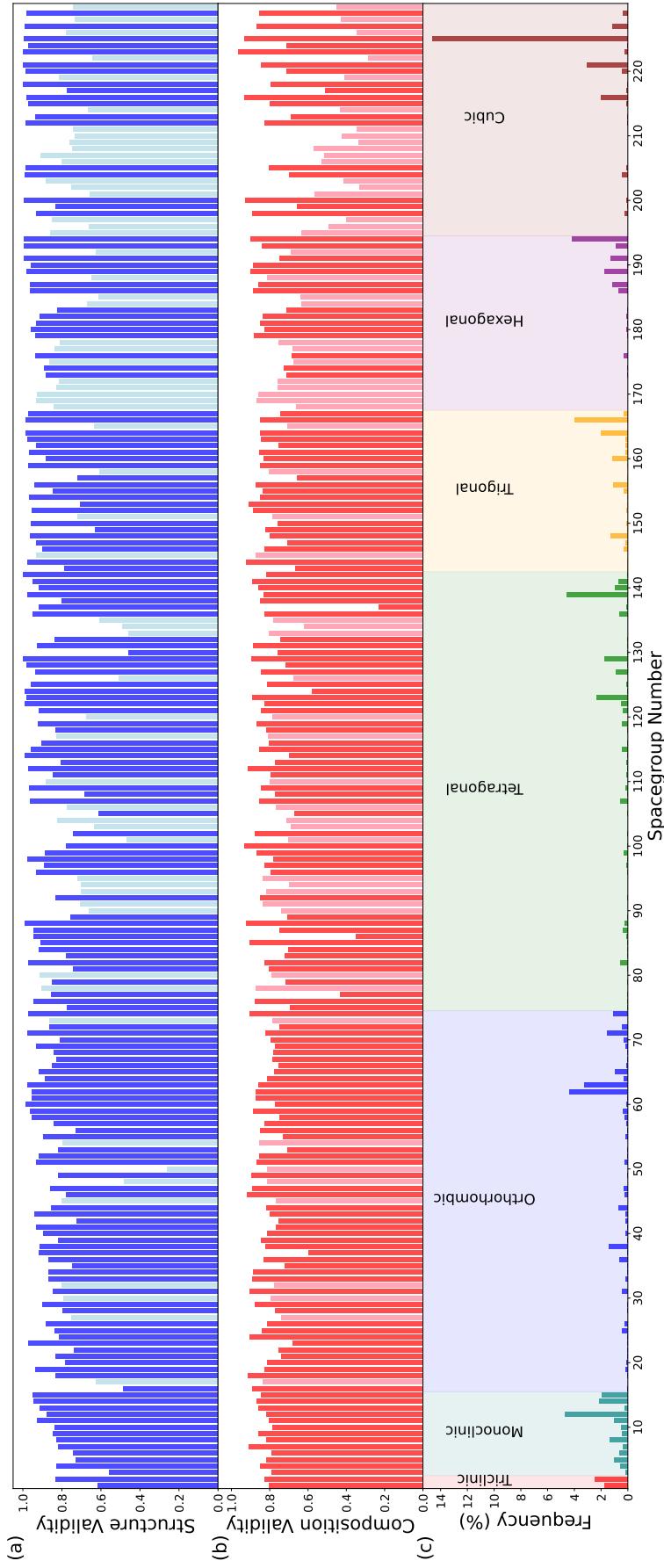


FIG. S2. (a) Structure validity histogram, (b) Composition validity histogram (c) frequency of space groups in the training dataset. The histograms in light colors in (a) and (b) are for those space groups where there is no training data. A dataset of 1000 generated samples per space group was curated for evaluation.

Appendix C: Discovered crystal samples with symmetric structure initialization

Table S3 and Table S4 list novel samples with $E_{\text{hull}} < 0.1$ eV/atom with CrystalFormer and PyXtal initializations respectively. Among them, Pd_3Pt , MgAgPd_2 , TaNbRu_2 , Gd_2HgAu , and TbMgPd_2 are both discovered by CrystalFormer and PyXtal initialization. These compounds share identical crystal structures in the two tables, with lattice constants differing by less than 0.02 Å.

To relax samples and estimate the energy above the hull, the DFT calculations were performed with the Perdew-Burke-Ernzerhof (PBE) exchange-correlation functional [125] and all-electron projector-augmented wave method [126], as implemented in the VASP code [127]. All parameters of the calculations including settings of PBE functional, Hubbard U corrections, and ferromagnetic initialization are chosen to be consistent with Materials Project by using of MPRelaxSet function in pymatgen [81]. A double relaxation strategy was employed. The maximum optimization ionic step and the maximum running time were constrained to 150 steps and 20 hours, respectively. All structures containing Yb element are ignored when calculating energy above the hull due to they are unavailable from the Materials Project at the time of writing.

Appendix D: Discovered crystal samples with element substitution

Tables S5 and S6 list double perovskites crystals in the $Fm\bar{3}m$ (No. 225) space group with $E_{\text{hull}} < 0.1$ eV/atom found by CrystalFormer and element substitution of pymatgen respectively. We found that CrystalFormer has found three times more perovskites crystals compared to ionic substitution [46]. The DFT calculation setup is the same as in the previous section.

Appendix E: Details of plug-and-play materials design

We build two independent predictive models of band gap and formation energy following the generative pretraining strategy [128]. Generative pretrain of CrystalFormer allows the model to extract crystal representations. Then, by attaching a regression neural network to its feature and training on property labels, the model can be used to predict the properties of crystals.

For the regression task, the mean absolute error corresponds to a predictive model $f(\mathbf{C})$ characterized by the Laplace distribution $p(\mathbf{y}|\mathbf{C}) \propto \exp(-\alpha|\mathbf{y} - f(\mathbf{C})|)$, where α is a positive scale parameter. According to Bayes' rule, the posterior probability for crystals given the property \mathbf{y} is:

$$p(\mathbf{C}|\mathbf{g}, \mathbf{y}) \propto p(\mathbf{y}|\mathbf{C})p(\mathbf{C}|\mathbf{g}) \propto e^{-\alpha|\mathbf{y} - f(\mathbf{C})|} p(\mathbf{C}|\mathbf{g}). \quad (\text{S1})$$

The crystal log-likelihood with a given property condition is then given by:

$$\ln p(\mathbf{C}|\mathbf{g}, \mathbf{y}) \propto \ln p(\mathbf{C}|\mathbf{g}) - \alpha|\mathbf{y} - f(\mathbf{C})|. \quad (\text{S2})$$

It is instructive to see that property-guided samples are balanced between the CrystalFormer prior and the property prediction model likelihood. Here α plays the role of guidance strength. For multi-property prediction, the likelihood function can be extended to:

$$\ln p(\mathbf{C}|\mathbf{g}, \mathbf{y}_1, \dots, \mathbf{y}_n) \propto \ln p(\mathbf{C}|\mathbf{g}) - \sum_{i=1}^n \alpha_i |\mathbf{y}_i - f_i(\mathbf{C})|. \quad (\text{S3})$$

Note that MCMC sampling via the Metropolis algorithm does not involve the normalization factor of the conditional distribution. In our experiments, we set $\alpha_1 = 3$ for band gap prediction and $\alpha_2 = 10$ for formation energy prediction, corresponding to the inverses of the MAE values obtained from the regression model.

TABLE S3. Discovered crystalline materials by CrystalFormer with $E_{\text{hull}} < 0.1$ eV/atom which are not in the MP-20 dataset.

Formula	Space group	Wyckoff-Atom sequence ¹	E_{hull} (eV/atom)
LuSiNi	62	c-Lu-c-Si-c-Ni	0.0134
Eu ₃ Mg	62	c-Eu-c-Mg-d-Eu	0.0661
EuAgAu	62	c-Eu-c-Ag-c-Au	0.0111
SmSbPd	62	c-Sm-c-Sb-c-Pd	0.0368
Nd ₃ Rh	62	c-Nd-c-Rh-d-Nd	0.0673
SrSnPd	62	c-Sr-c-Sn-c-Pd	-0.0531
SrCdSn	62	c-Sr-c-Cd-c-Sn	0.0863
SrMgHg	62	c-Sr-c-Mg-c-Hg	0.0517
Eu ₂ Ru	62	c-Eu-c-Eu-c-Ru	0.0890
PrHgAu ₂	62	c-Pr-c-Hg-c-Au-c-Au	0.0868
Ce(SiPt) ₂	139	a-Ce-d-Si-e-Pt	0.0648
Dy(AlGa) ₂	139	a-Dy-d-Al-e-Ga	0.0902
Nd ₃ Dy	139	a-Dy-b-Nd-d-Nd	0.0690
Pd ₃ Pt	139	a-Pt-b-Pd-d-Pd	-0.0298
LaMgSn	139	c-Mg-e-La-e-Sn	0.0274
Rb ₃ Mn ₂ F ₇	139	a-F-b-Rb-e-Rb-e-Mn-e-F-g-F	0.0874
AlAgAu ₂	139	a-Al-b-Ag-d-Au	0.0912
Ba ₃ Tl	139	a-Tl-b-Ba-d-Ba	0.0933
CsPrTe ₂	166	a-Cs-b-Pr-c-Te	0.0977
Sr ₂ Ca	166	a-Ca-c-Sr	0.0771
Pr ₂ AlRu ₃	166	a-Al-c-Pr-d-Ru	0.0909
Dy ₃ Ho	194	d-Ho-h-Dy	0.0738
Zr ₃ Pb	194	d-Pb-h-Zr	0.0426
PrAu ₃	194	d-Pr-h-Au	0.0696
Pd ₃ Pt	194	d-Pt-h-Pd	0.0053
TbZnGa	194	a-Tb-c-Ga-d-Zn	0.0718
BaSr(GaGe) ₂	194	a-Ba-b-Sr-f-Ga-f-Ge	0.0855
Ho ₂ Er	194	b-Er-f-Ho	0.0790
NdDyHg ₂	225	a-Nd-b-Dy-c-Hg	0.0974
MgAgPd ₂	225	a-Mg-b-Ag-c-Pd	0.0540
PrEuIn ₂	225	a-Eu-b-Pr-c-In	0.0957
TaNbRu ₂	225	a-Ta-b-Nb-c-Ru	0.0823
PdAu	221	a-Au-b-Pd	0.0025
Gd ₂ HgAu	225	a-Hg-b-Au-c-Gd	0.0522

¹ The fractional coordinates and lattice parameters are omitted for brevity. See details at https://drive.google.com/file/d/1gOIkWkjSH_Ed0-wzPk8VgxLjJcidbrH3/view?usp=sharing

TABLE S4. Discovered crystalline materials by PyXtal initialization with $E_{\text{hull}} < 0.1$ eV/atom which are not in the MP-20 dataset.

Formula	Space group	Wyckoff-Atom sequence	${}^1E_{\text{hull}}$ (eV/atom)
Ba ₂ Si	62	b-Ba-c-Ba-c-Si	0.0766
GdIr	62	c-Gd-c-Ir	-0.2310
Nd ₃ Dy	139	a-Dy-b-Nd-d-Nd	0.0680
Pd ₃ Pt	139	a-Pd-b-Pt-d-Pd	0.0147
PdPt	221	a-Pd-b-Pt	0.0847
NdDyHg ₂	225	a-Dy-b-Nd-c-Hg	0.0973
MgAgPd ₂	225	a-Mg-b-Ag-c-Pd	0.0480
HfRe	221	a-Re-b-Hf	0.0926
TaNbRu ₂	225	a-Nb-b-Ta-c-Ru	0.0770
PdAu	221	a-Au-b-Pd	0.0056
Gd ₂ HgAu	225	a-Hg-b-Au-c-Gd	0.0476
TbMgPd ₂	225	a-Mg-b-Tb-c-Pd	0.0676

¹ The fractional coordinates and lattice parameters are omitted for brevity. See details at https://drive.google.com/file/d/18sdsp-6yRSBaNez1A0lr20Ru6H7OOzvJ/view?usp=drive_link.

TABLE S5. Discovered double perovskites by likelihood guided structure mutation of CrystalFormer with $E_{\text{hull}} < 0.1$ eV/atom which are not in the MP-20 dataset.

Formula	Space group	Wyckoff-Atom sequence	${}^1E_{\text{hull}}$ (eV/atom)	Lattice constant (Å)
Sr ₂ TaNiO ₆	225	a-Ni-b-Ta-c-Sr-e-O	0.0426	7.9762
Ba ₂ TaNiO ₆	225	a-Ni-b-Ta-c-Ba-e-O	0.0711	8.2129
Ba ₂ YMnO ₆	225	a-Y-b-Mn-c-Ba-e-O	0.0284	8.3575
Eu ₂ CoMoO ₆	225	a-Co-b-Mo-c-Eu-e-O	0.0720	7.9220
Ba ₂ TbFeO ₆	225	a-Tb-b-Fe-c-Ba-e-O	0.0343	8.3200
Ba ₂ HfWO ₆	225	a-W-b-Hf-c-Ba-e-O	0.0939	8.3568
Eu ₂ CoNiO ₆	225	a-Ni-b-Co-c-Eu-e-O	0.0646	7.6181
Sr ₂ NbVO ₆	225	a-V-b-Nb-c-Sr-e-O	0.0166	8.0375
Sr ₂ NaWO ₆	225	a-Na-b-W-c-Sr-e-O	0.0901	8.2959

¹ The fractional coordinates are omitted for brevity.

TABLE S6. Discovered double perovskites using the element substitution of Hautier et al. [46] with $E_{\text{hull}} < 0.1$ eV/atom which are not in the MP-20 dataset.

Formula	Space group	Wyckoff-Atom sequence	${}^1E_{\text{hull}}$ (eV/atom)	Lattice constant (Å)
Ba ₂ FeBiO ₆	225	a-Bi-b-Fe-c-Ba-e-O	-0.0030	8.3808
Ba ₂ NpCoO ₆	225	a-Np-b-Co-c-Ba-e-O	0.0241	8.3983
Ba ₂ BiMoO ₆	225	a-Bi-b-Mo-c-Ba-e-O	0.0706	8.6322

¹ The fractional coordinates are omitted for brevity.