

---

# Crystal Structure Prediction by Joint Equivariant Diffusion

---

Rui Jiao<sup>1,2</sup> Wenbing Huang<sup>3,4\*</sup> Peijia Lin<sup>5</sup> Jiaqi Han<sup>6</sup> Pin Chen<sup>5</sup> Yutong Lu<sup>5</sup> Yang Liu<sup>1,2\*</sup>

<sup>1</sup>Dept. of Comp. Sci. & Tech., Institute for AI, Tsinghua University

<sup>2</sup>Institute for AIR, Tsinghua University

<sup>3</sup>Gaoling School of Artificial Intelligence, Renmin University of China

<sup>4</sup> Beijing Key Laboratory of Big Data Management and Analysis Methods, Beijing, China

<sup>5</sup> National Supercomputer Center in Guangzhou,

School of Computer Science and Engineering, Sun Yat-sen University

<sup>6</sup> Stanford University

## Abstract

Crystal Structure Prediction (CSP) is crucial in various scientific disciplines. While CSP can be addressed by employing currently-prevailing generative models (*e.g.* diffusion models), this task encounters unique challenges owing to the symmetric geometry of crystal structures—the invariance of translation, rotation, and periodicity. To incorporate the above symmetries, this paper proposes DiffCSP, a novel diffusion model to learn the structure distribution from stable crystals. To be specific, DiffCSP jointly generates the lattice and atom coordinates for each crystal by employing a periodic-E(3)-equivariant denoising model, to better model the crystal geometry. Notably, different from related equivariant generative approaches, DiffCSP leverages fractional coordinates other than Cartesian coordinates to represent crystals, remarkably promoting the diffusion and the generation process of atom positions. Extensive experiments verify that our DiffCSP significantly outperforms existing CSP methods, with a much lower computation cost in contrast to DFT-based methods. Moreover, the superiority of DiffCSP is also observed when it is extended for *ab initio* crystal generation. Code is available at <https://github.com/jiaor17/DiffCSP>.

## 1 Introduction

Crystal Structure Prediction (CSP), which returns the stable 3D structure of a compound based solely on its composition, has been a goal in physical sciences since the 1950s [1]. As crystals are the foundation of various materials, estimating their structures in 3D space determines the physical and chemical properties that greatly influence the application to various academic and industrial sciences, such as the design of drugs, batteries, and catalysts [2]. Conventional methods towards CSP mostly apply the Density Functional Theory (DFT) [3] to compute the energy at each iteration, guided by optimization algorithms (such as random search [4], Bayesian optimization [5], etc.) to iteratively search for the stable state corresponding to the local minima of the energy surface [6].

The DFT-based approaches are computationally-intensive. Recent attention has been paid to deep generative models, which directly learn the distribution from the training data consisting of stable structures [7, 8]. More recently, diffusion models, a special kind of deep generative models are employed for crystal generation [9], encouraged by their better physical interpretability and enhanced performance than other generative models. Intuitively, by conducting diffusion on stable structures,

---

\*Wenbing Huang and Yang Liu are corresponding authors.

the denoising process in diffusion models acts like a force field that drives the atom coordinates towards the energy local minimum and thus is able to increase stability. Indeed, the success of diffusion models is observed in broad scientific domains, including molecular conformation generation [10], protein structure prediction [11] and protein docking [12].

However, designing diffusion models for CSP is challenging. From the perspective of physics, any  $E(3)$  transformation, including translation, rotation, and reflection, of the crystal coordinates does not change the physical law and thus keeps the crystal distribution invariant. In other words, the generation process we design should yield  $E(3)$  invariant samples. Moreover, in contrast to other types of structures such as small molecules [13] and proteins [14], CSP exhibits unique challenges, mainly incurred by the periodicity of the atom arrangement in crystals. Figure 1 displays a crystal where the atoms in a unit cell are repeated infinitely in space. We identify such unique symmetry, jointly consisting of  $E(3)$  invariance and periodicity, as *periodic  $E(3)$  invariance* in this paper. Generating such type of structures requires not only modeling the distribution of the atom coordinates within every cell, but also inferring how their bases (*a.k.a.* lattice vectors) are placed in 3D space. Interestingly, as we will show in § 4.1, such view offers a natural disentanglement for fulfilling the periodic  $E(3)$  invariance by separately enforcing constraints on fractional coordinates and lattice vectors, which permits a feasible implementation to encode the crystal symmetry.

In this work, we introduce DiffCSP, an equivariant diffusion method to address CSP. Considering the specificity of the crystal geometry, our DiffCSP jointly generates the lattice vectors and the fractional coordinates of all atoms, by employing a proposed denoising model that is theoretically proved to generate periodic- $E(3)$ -invariant samples. A preferable characteristic of DiffCSP is that it leverages the fractional coordinate system (defined in § 3) other than the Cartesian system used in previous methods to represent crystals [9, 15], which encodes periodicity intrinsically. In particular, the fractional representation not only allows us to consider Wrapped Normal (WN) distribution [16] to better model the periodicity, but also facilitates the design of the denoising model via the Fourier transformation, compared to the traditional multi-graph encoder in crystal modeling [15].

CDVAE [9] is closely related with our paper. It adopts an equivariant Variational Auto-Encoder (VAE) based framework to learn the data distribution and then generates crystals in a score-matching-based diffusion process. However, CDVAE focuses mainly on ab initio crystal generation where the composition is also randomly sampled, which is distinct from the CSP task in this paper. Moreover, while CDVAE first predicts the lattice and then updates the coordinates with the fixed lattice, we jointly update the lattice and coordinates to better model the crystal geometry. Besides, CDVAE represents crystals by Cartesian coordinates upon multi-graph modeling, whereas our DiffCSP applies fractional coordinates without multi-graph modeling as mentioned above.

To sum up, our contributions are as follows:

- To the best of our knowledge, we are the first to apply equivariant diffusion-based methods to address CSP. The proposed DiffCSP is more insightful than current learning-based approaches as the periodic  $E(3)$  invariance has been delicately considered.
- DiffCSP conducts joint diffusion on lattices and fractional coordinates, which is capable of capturing the crystal geometry as a whole. Besides, the usage of fractional coordinates in place of Cartesian coordinates used in previous methods (*e.g.* CDVAE [9]) remarkably promotes the diffusion and the generation process of atom positions.
- We verify the efficacy of DiffCSP on the CSP task against learning-based and DFT-based methods, and sufficiently ablate each proposed component in DiffCSP. We further extend DiffCSP into ab initio generation and show its effectiveness against related methods.

## 2 Related Works

**Crystal Structure Prediction** Traditional computation methods [4, 5, 17, 18] combine DFT [3] with optimization algorithms to search for local minima in the potential energy surface. However, DFT is computationally intensive, making it dilemmatic to balance efficiency and accuracy. With the improvement of crystal databases, machine-learning methods are applied as alternative energy predictors to DFT followed by optimization steps [19, 20, 21]. Apart from the predict-optimize paradigm, another line of approaches directly learns stable structures from data by deep generative models, which represents crystals by 3D voxels [7, 22, 23], distance matrices [8, 24, 25] or 3D

coordinates [26, 27, 28]. Unfortunately, these methods are unaware of the full symmetries of the crystal structure. CDVAE [9] has taken the required symmetries into account. However, as mentioned above, the initial version of CDVAE is for different task and utilizes different generation process.

**Equivariant Graph Neural Networks** Geometrically equivariant Graph Neural Networks (GNNs) that ensure  $E(3)$  symmetry are powerful tools to represent physical objects [29, 30, 31, 32, 33], and have showcased the superiority in modeling 3D structures [34, 35]. To further model the periodic materials, Xie and Grossman [15] propose the multi-graph edge construction to capture the periodicity by connecting the edges between adjacent lattices. Yan et al. [36] further introduce periodic pattern encoding into a Transformer-based backbone. In this work, we achieve the periodic invariance by introducing the Fourier transform on fractional coordinates.

**Diffusion Generative Models** Motivated by the non-equilibrium thermodynamics [37], diffusion models connect the data distribution with the prior distribution via forward and backward Markov chains [38], and have made remarkable progress in the field of image generation [39, 40]. Equipped with equivariant GNNs, diffusion models are capable of generating samples from the invariant distribution, which is desirable in conformation generation [10, 13], ab initio molecule design [41], protein generation [42], and so on. Recent works extend the diffusion models onto Riemann manifolds [43, 44], and enable the generation of periodic features like torsion angles [12, 16].

### 3 Preliminaries

**Representation of crystal structures** A 3D crystal can be represented as the infinite periodic arrangement of atoms in 3D space, and the smallest repeating unit is called a *unit cell*, as shown in Figure 1. A unit cell can be defined by a triplet  $\mathcal{M} = (\mathbf{A}, \mathbf{X}, \mathbf{L})$ , where  $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_N] \in \mathbb{R}^{h \times N}$  denotes the list of the one-hot representations of atom types,  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{3 \times N}$  consists of Cartesian coordinates of the atoms, and  $\mathbf{L} = [l_1, l_2, l_3] \in \mathbb{R}^{3 \times 3}$  represents the lattice matrix containing three basic vectors to describe the periodicity of the crystal. The infinite periodic crystal structure is represented by

$$\{(\mathbf{a}'_i, \mathbf{x}'_i) | \mathbf{a}'_i = \mathbf{a}_i, \mathbf{x}'_i = \mathbf{x}_i + \mathbf{L}\mathbf{k}, \forall \mathbf{k} \in \mathbb{Z}^{3 \times 1}\}, \quad (1)$$

where the  $j$ -th element of the integral vector  $\mathbf{k}$  denotes the integral 3D translation in units of  $l_j$ .

**Fractional coordinate system** The Cartesian coordinate system  $\mathbf{X}$  leverages three standard orthogonal bases as the coordinate axes. In crystallography, the fractional coordinate system is usually applied to reflect the periodicity of the crystal structure [26, 27, 28, 45], which utilizes the lattices  $(l_1, l_2, l_3)$  as the bases. In this way, a point represented by the fractional coordinate vector  $\mathbf{f} = [f_1, f_2, f_3]^T \in [0, 1]^3$  corresponds to the Cartesian vector  $\mathbf{x} = \sum_{i=1}^3 f_i l_i$ . This paper employs the fractional coordinate system, and denotes the crystal by  $\mathcal{M} = (\mathbf{A}, \mathbf{F}, \mathbf{L})$ , where the fractional coordinates of all atoms in a cell compose the matrix  $\mathbf{F} \in [0, 1]^{3 \times N}$ .

**Task definition** CSP predicts for each unit cell the lattice matrix  $\mathbf{L}$  and the fractional matrix  $\mathbf{F}$  given its chemical composition  $\mathbf{A}$ , namely, learning the conditional distribution  $p(\mathbf{L}, \mathbf{F} | \mathbf{A})$ .

## 4 The Proposed Method: DiffCSP

This section first presents the symmetries of the crystal geometry, and then introduces the joint equivariant diffusion process on  $\mathbf{L}$  and  $\mathbf{F}$ , followed by the architecture of the denoising function.

### 4.1 Symmetries of Crystal Structure Distribution

While various generative models can be utilized to address CSP, this task encounters particular challenges, including constraints arising from symmetries of crystal structure distribution. Here, we consider the three types of symmetries in the distribution of  $p(\mathbf{L}, \mathbf{F} | \mathbf{A})$ : permutation invariance,  $O(3)$  invariance, and periodic translation invariance. Their detailed definitions are provided as follows.

**Definition 1** (Permutation Invariance). *For any permutation  $\mathbf{P} \in S_N$ ,  $p(\mathbf{L}, \mathbf{F} | \mathbf{A}) = p(\mathbf{L}, \mathbf{F}\mathbf{P} | \mathbf{A}\mathbf{P})$ , i.e., changing the order of atoms will not change the distribution.*

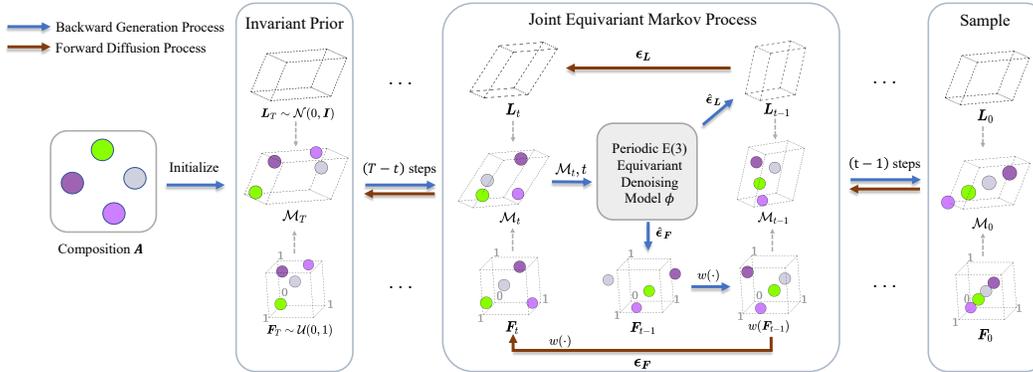


Figure 2: Overview of DiffCSP. Given the composition  $\mathbf{A}$ , we denote the crystal, its lattice and fractional coordinate matrix at time  $t$  as  $\mathcal{M}_t$ ,  $\mathbf{L}_t$  and  $\mathbf{F}_t$ , respectively. The terms  $\epsilon_{\mathbf{L}}$  and  $\epsilon_{\mathbf{F}}$  are Gaussian noises,  $\hat{\epsilon}_{\mathbf{L}}$  and  $\hat{\epsilon}_{\mathbf{F}}$  are predicted by the denoising model  $\phi$ .

**Definition 2** ( $O(3)$  Invariance). For any orthogonal transformation  $\mathbf{Q} \in \mathbb{R}^{3 \times 3}$  satisfying  $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$ ,  $p(\mathbf{Q}\mathbf{L}, \mathbf{F} \mid \mathbf{A}) = p(\mathbf{L}, \mathbf{F} \mid \mathbf{A})$ , namely, any rotation/reflection of  $\mathbf{L}$  keeps the distribution unchanged.

**Definition 3** (Periodic Translation Invariance). For any translation  $\mathbf{t} \in \mathbb{R}^{3 \times 1}$ ,  $p(\mathbf{L}, w(\mathbf{F} + \mathbf{t}\mathbf{1}^T) \mid \mathbf{A}) = p(\mathbf{L}, \mathbf{F} \mid \mathbf{A})$ , where the function  $w(\mathbf{F}) = \mathbf{F} - \lfloor \mathbf{F} \rfloor \in [0, 1)^{3 \times N}$  returns the fractional part of each element in  $\mathbf{F}$ , and  $\mathbf{1} \in \mathbb{R}^{3 \times 1}$  is a vector with all elements set to one. It explains that any periodic translation of  $\mathbf{F}$  will not change the distribution<sup>2</sup>.

The permutation invariance is tractably encapsulated by using GNNs as the backbone for generation [47]. We mainly focus on the other two kinds of invariance (see Figure 1), since GNNs are our default choices. For simplicity, we compactly term the  $O(3)$  invariance and periodic translation invariance as *periodic  $E(3)$  invariance* henceforth. Previous approaches (e.g. [9, 36]) adopt Cartesian coordinates  $\mathbf{X}$  other than fractional coordinates  $\mathbf{F}$ , hence their derived forms of the symmetry are different. Particularly, in Definition 2, the orthogonal transformation additionally acts on  $\mathbf{X}$ ; in Definition 3, the periodic translation  $w(\mathbf{F} + \mathbf{t}\mathbf{1}^T)$  becomes the translation along the lattice bases  $\mathbf{X} + \mathbf{L}\mathbf{t}\mathbf{1}^T$ ; besides,  $\mathbf{X}$  should also maintain  $E(3)$  translation invariance, that is  $p(\mathbf{L}, \mathbf{X} + \mathbf{t}\mathbf{1}^T \mid \mathbf{A}) = p(\mathbf{L}, \mathbf{X} \mid \mathbf{A})$ . With the help of the fractional system, the periodic  $E(3)$  invariance is made tractable by fulfilling  $O(3)$  invariance *w.r.t.* the orthogonal transformations on  $\mathbf{L}$  and periodic translation invariance *w.r.t.* the periodic translations on  $\mathbf{F}$ , respectively. In this way, such approach, as detailed in the next section, facilitates the application of diffusion methods to the CSP task.

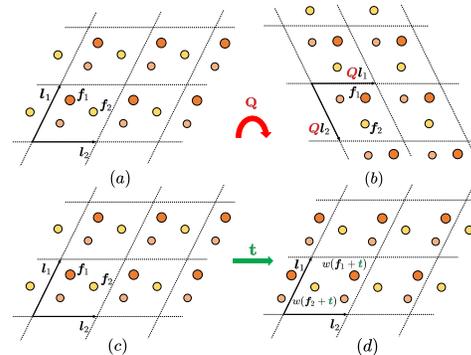


Figure 1: (a)→(b): The orthogonal transformation of the lattice vectors. (c)→(d): The periodic translation of the fractional coordinates. Both cases do not change the structure.

## 4.2 Joint Equivariant Diffusion

Our method DiffCSP addresses CSP by simultaneously diffusing the lattice  $\mathbf{L}$  and the fractional coordinate matrix  $\mathbf{F}$ . Given the atom composition  $\mathbf{A}$ ,  $\mathcal{M}_t$  denotes the intermediate state of  $\mathbf{L}$  and  $\mathbf{F}$  at time step  $t$  ( $0 \leq t \leq T$ ). DiffCSP defines two Markov processes: the forward diffusion process gradually adds noise to  $\mathcal{M}_0$ , and the backward generation process iteratively samples from the prior distribution  $\mathcal{M}_T$  to recover the origin data  $\mathcal{M}_0$ .

<sup>2</sup>Previous works (e.g. [36]) further discuss the scaling invariance of a unit cell formed by periodic boundaries, allowing  $\mathbf{L} \rightarrow \alpha\mathbf{L}, \forall \alpha \in \mathbb{N}_+^3$ . In this paper, the scaling invariance is unnecessary since we apply the Niggli reduction [46] on the primitive cell as a canonical scale representation of the lattice vectors where we fix  $\alpha = (1, 1, 1)^T$ . Additionally, periodic translation invariance in our paper is equivalent to the invariance of shifting periodic boundaries in [36]. We provide more discussions in Appendix A.4.

Joining the statements in § 4.1, the recovered distribution from  $\mathcal{M}_T$  should meet periodic E(3) invariance. Such requirement is satisfied if the prior distribution  $p(\mathcal{M}_T)$  is invariant and the Markov transition  $p(\mathcal{M}_{t-1} | \mathcal{M}_t)$  is equivariant, according to the diffusion-based generation literature [10]. Here, an equivariant transition is specified as  $p(g \cdot \mathcal{M}_{t-1} | g \cdot \mathcal{M}_t) = p(\mathcal{M}_{t-1} | \mathcal{M}_t)$  where  $g \cdot \mathcal{M}$  refers to any orthogonal/translational transformation  $g$  acts on  $\mathcal{M}$  in the way presented in Definitions 2-3. We separately explain the derivation details of  $\mathbf{L}$  and  $\mathbf{F}$  below. The detailed flowcharts are summarized in Algorithms 1 and 2 in Appendix B.3.

**Diffusion on  $\mathbf{L}$**  Given that  $\mathbf{L}$  is a continuous variable, we exploit Denoising Diffusion Probabilistic Model (DDPM) [38] to accomplish the generation. We define the forward process that progressively diffuses  $\mathbf{L}_0$  towards the Normal prior  $p(\mathbf{L}_T) = \mathcal{N}(0, \mathbf{I})$  by  $q(\mathbf{L}_t | \mathbf{L}_{t-1})$  which can be devised as the probability conditional on the initial distribution:

$$q(\mathbf{L}_t | \mathbf{L}_0) = \mathcal{N}\left(\mathbf{L}_t | \sqrt{\bar{\alpha}_t} \mathbf{L}_0, (1 - \bar{\alpha}_t) \mathbf{I}\right), \quad (2)$$

where  $\beta_t \in (0, 1)$  controls the variance, and  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s = \prod_{s=1}^t (1 - \beta_s)$  is valued in accordance to the cosine scheduler [48].

The backward generation process is given by:

$$p(\mathbf{L}_{t-1} | \mathcal{M}_t) = \mathcal{N}(\mathbf{L}_{t-1} | \mu(\mathcal{M}_t), \sigma^2(\mathcal{M}_t) \mathbf{I}), \quad (3)$$

where  $\mu(\mathcal{M}_t) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{L}_t - \frac{\beta_t}{\sqrt{1-\alpha_t}} \hat{\epsilon}_{\mathbf{L}}(\mathcal{M}_t, t) \right)$ ,  $\sigma^2(\mathcal{M}_t) = \beta_t \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}$ . The denoising term  $\hat{\epsilon}_{\mathbf{L}}(\mathcal{M}_t, t) \in \mathbb{R}^{3 \times 3}$  is predicted by the model  $\phi(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t)$ .

As the prior distribution  $p(\mathbf{L}_T) = \mathcal{N}(0, \mathbf{I})$  is already O(3)-invariant, we require the generation process in Eq. (3) to be O(3)-equivariant, which is formally stated below.

**Proposition 1.** *The marginal distribution  $p(\mathbf{L}_0)$  by Eq. (3) is O(3)-invariant if  $\hat{\epsilon}_{\mathbf{L}}(\mathcal{M}_t, t)$  is O(3)-equivariant, namely  $\hat{\epsilon}_{\mathbf{L}}(\mathbf{Q}\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t) = \mathbf{Q}\hat{\epsilon}_{\mathbf{L}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t), \forall \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$ .*

To train the denoising model  $\phi$ , we first sample  $\epsilon_{\mathbf{L}} \sim \mathcal{N}(0, \mathbf{I})$  and reparameterize  $\mathbf{L}_t = \sqrt{\bar{\alpha}_t} \mathbf{L}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_{\mathbf{L}}$  based on Eq. (2). The training objective is defined as the  $\ell_2$  loss between  $\epsilon_{\mathbf{L}}$  and  $\hat{\epsilon}_{\mathbf{L}}$ :

$$\mathcal{L}_{\mathbf{L}} = \mathbb{E}_{\epsilon_{\mathbf{L}} \sim \mathcal{N}(0, \mathbf{I}), t \sim \mathcal{U}(1, T)} [\|\epsilon_{\mathbf{L}} - \hat{\epsilon}_{\mathbf{L}}(\mathcal{M}_t, t)\|_2^2]. \quad (4)$$

**Diffusion on  $\mathbf{F}$**  The domain of fractional coordinates  $[0, 1)^{3 \times N}$  forms a quotient space  $\mathbb{R}^{3 \times N} / \mathbb{Z}^{3 \times N}$  induced by the crystal periodicity. It is not suitable to apply the above DDPM fashion to generate  $\mathbf{F}$ , as the normal distribution used in DDPM is unable to model the cyclical and bounded domain of  $\mathbf{F}$ . Instead, we leverage Score-Matching (SM) based framework [49, 50] along with Wrapped Normal (WN) distribution [43] to fit the specificity here. Note that WN distribution has been explored in generative models, such as molecular conformation generation [16].

During the forward process, we first sample each column of  $\epsilon_{\mathbf{F}} \in \mathbb{R}^{3 \times N}$  from  $\mathcal{N}(0, \mathbf{I})$ , and then acquire  $\mathbf{F}_t = w(\mathbf{F}_0 + \sigma_t \epsilon_{\mathbf{F}})$  where the truncation function  $w(\cdot)$  is already defined in Definition 3. This truncated sampling implies the WN transition:

$$q(\mathbf{F}_t | \mathbf{F}_0) \propto \sum_{\mathbf{Z} \in \mathbb{Z}^{3 \times N}} \exp\left(-\frac{\|\mathbf{F}_t - \mathbf{F}_0 + \mathbf{Z}\|_{\mathbf{F}}^2}{2\sigma_t^2}\right). \quad (5)$$

Basically, this process ensures the probability distribution over  $[z, z + 1)^{3 \times N}$  for any integer  $z$  to be the same to keep the crystal periodicity. Here, the noise scale  $\sigma_t$  obeys the exponential scheduler:  $\sigma_0 = 0$  and  $\sigma_t = \sigma_1 \left(\frac{\sigma_T}{\sigma_1}\right)^{\frac{t-1}{T-1}}$ , if  $t > 0$ . Desirably,  $q(\mathbf{F}_t | \mathbf{F}_0)$  is periodic translation equivariant, and approaches a uniform distribution  $\mathcal{U}(0, 1)$  if  $\sigma_T$  is sufficiently large.

For the backward process, we first initialize  $\mathbf{F}_T$  from the uniform distribution  $\mathcal{U}(0, 1)$ , which is periodic translation invariant. With the denoising term  $\hat{\epsilon}_{\mathbf{F}}$  predicted by  $\phi(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t)$ , we combine the ancestral predictor [38, 50] with the Langevin corrector [49] to sample  $\mathbf{F}_0$ . We immediately have:

**Proposition 2.** *The marginal distribution  $p(\mathbf{F}_0)$  is periodic translation invariant if  $\hat{\epsilon}_{\mathbf{F}}(\mathcal{M}_t, t)$  is periodic translation invariant, namely  $\hat{\epsilon}_{\mathbf{F}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t) = \hat{\epsilon}_{\mathbf{F}}(\mathbf{L}_t, w(\mathbf{F}_t + \mathbf{t}\mathbf{1}^\top), \mathbf{A}, t), \forall \mathbf{t} \in \mathbb{R}^3$ .*

The training objective for score matching is:

$$\mathcal{L}_{\mathbf{F}} = \mathbb{E}_{\mathbf{F}_t \sim q(\mathbf{F}_t | \mathbf{F}_0), t \sim \mathcal{U}(1, T)} [\lambda_t \|\nabla_{\mathbf{F}_t} \log q(\mathbf{F}_t | \mathbf{F}_0) - \hat{\epsilon}_{\mathbf{F}}(\mathcal{M}_t, t)\|_2^2],$$

where  $\lambda_t = \mathbb{E}_{\mathbf{F}_t}^{-1} [\|\nabla_{\mathbf{F}_t} \log q(\mathbf{F}_t | \mathbf{F}_0)\|_2^2]$  is approximated via Monte-Carlo sampling. More details are deferred to Appendix B.1.

**Extension to ab initio crystal generation** Although our method is proposed to address CSP where the composition  $\mathbf{A}$  is fixed, our method is able to be extended for the ab initio generation task by further generating  $\mathbf{A}$ . We achieve this by additionally optimizing the one-hot representation  $\mathbf{A}$  with a DDPM-based approach. We provide more details in Appendix G.

### 4.3 The Architecture of the Denoising Model

This subsection designs the denoising model  $\phi(\mathbf{L}, \mathbf{F}, \mathbf{A}, t)$  that outputs  $\hat{\epsilon}_{\mathbf{L}}$  and  $\hat{\epsilon}_{\mathbf{F}}$  satisfying the properties stated in Proposition 1 and 2.

Let  $\mathbf{H}^{(s)} = [\mathbf{h}_1^{(s)}, \dots, \mathbf{h}_N^{(s)}]$  denote the node representations of the  $s$ -th layer. The input feature is given by  $\mathbf{h}_i^{(0)} = \rho(f_{\text{atom}}(\mathbf{a}_i), f_{\text{pos}}(t))$ , where  $f_{\text{atom}}$  and  $f_{\text{pos}}$  are the atomic embedding and sinusoidal positional encoding [38, 51], respectively;  $\rho$  is a multi-layer perceptron (MLP).

Built upon EGNN [32], the  $s$ -th layer message-passing is unfolded as follows:

Here  $\varphi_m$  and  $\varphi_h$  are MLPs. The function  $\psi_{\text{FT}} : (-1, 1)^3 \rightarrow [-1, 1]^{3 \times K}$  is Fourier Transformation of the relative fractional coordinate  $\mathbf{f}_j - \mathbf{f}_i$ . Specifically, suppose the input to be  $\mathbf{f} = [f_1, f_2, f_3]^T$ , then the  $c$ -th row and  $k$ -th column of the output is calculated by

$$\mathbf{m}_{ij}^{(s)} = \varphi_m(\mathbf{h}_i^{(s-1)}, \mathbf{h}_j^{(s-1)}, \mathbf{L}^T \mathbf{L}, \psi_{\text{FT}}(\mathbf{f}_j - \mathbf{f}_i)), \quad (6)$$

$$\mathbf{m}_i^{(s)} = \sum_{j=1}^N \mathbf{m}_{ij}^{(s)}, \quad (7)$$

$$\mathbf{h}_i^{(s)} = \mathbf{h}_i^{(s-1)} + \varphi_h(\mathbf{h}_i^{(s-1)}, \mathbf{m}_i^{(s)}). \quad (8)$$

$\psi_{\text{FT}}(\mathbf{f})[c, k] = \sin(2\pi m f_c)$ , if  $k = 2m$  (even); and  $\psi_{\text{FT}}(\mathbf{f})[c, k] = \cos(2\pi m f_c)$ , if  $k = 2m + 1$  (odd).  $\psi_{\text{FT}}$  extracts various frequencies of all relative fractional distances that are helpful for crystal structure modeling, and more importantly,  $\psi_{\text{FT}}$  is periodic translation invariant, namely,  $\psi_{\text{FT}}(w(\mathbf{f}_j + \mathbf{t}) - w(\mathbf{f}_i + \mathbf{t})) = \psi_{\text{FT}}(\mathbf{f}_j - \mathbf{f}_i)$  for any translation  $\mathbf{t}$ , which is proved in Appendix A.3.

After  $S$  layers of message passing conducted on the fully connected graph, the lattice noise  $\hat{\epsilon}_{\mathbf{L}}$  is acquired by a linear combination of  $\mathbf{L}$ , with the weights given by the final layer:

$$\hat{\epsilon}_{\mathbf{L}} = \mathbf{L} \varphi_L \left( \frac{1}{N} \sum_{i=1}^N \mathbf{h}_i^{(S)} \right), \quad (9)$$

where  $\varphi_L$  is an MLP with output shape as  $3 \times 3$ . The fractional coordinate score  $\hat{\epsilon}_{\mathbf{F}}$  is output by:

$$\hat{\epsilon}_{\mathbf{F}}[:, i] = \varphi_F(\mathbf{h}_i^{(S)}), \quad (10)$$

where  $\hat{\epsilon}_{\mathbf{F}}[:, i]$  defines the  $i$ -th column of  $\hat{\epsilon}_{\mathbf{F}}$ , and  $\varphi_F$  is an MLP on the final representation.

We apply the inner product term  $\mathbf{L}^T \mathbf{L}$  in Eq. (6) to achieve O(3)-invariance, as  $(\mathbf{Q}\mathbf{L})^T (\mathbf{Q}\mathbf{L}) = \mathbf{L}^T \mathbf{L}$  for any orthogonal matrix  $\mathbf{Q} \in \mathbb{R}^{3 \times 3}$ . This leads to the O(3)-invariance of  $\varphi_L$  in Eq. (10), and we further left-multiply  $\mathbf{L}$  with  $\varphi_L$  to ensure the O(3)-equivariance of  $\hat{\epsilon}_{\mathbf{L}}$ . Therefore, the above formulation of the denoising model  $\phi(\mathbf{L}, \mathbf{F}, \mathbf{A}, t)$  ensures the following property.

**Proposition 3.** *The score  $\hat{\epsilon}_{\mathbf{L}}$  by Eq. (9) is O(3)-equivariant, and the score  $\hat{\epsilon}_{\mathbf{F}}$  from Eq. (10) is periodic translation invariant. Hence, the generated distribution by DiffCSP is periodic E(3) invariant.*

**Comparison with multi-graph representation** Previous methods [9, 15, 29, 52] utilize Cartesian coordinates, and usually describe crystals with multi-graph representation to encode the periodic structures. They create multiple edges to connect each pair of nodes where different edges refer to different integral cell translations. Here, we no longer require multi-graph representation, since we employ fractional coordinates that naturally encode periodicity and the Fourier transform  $\psi_{\text{FT}}$  in our message passing is already periodic translation invariant. We will ablate the benefit in Table 3.

Table 1: Results on stable structure prediction task.

	# of samples	Perov-5		Carbon-24		MP-20		MPTS-52	
		Match rate $\uparrow$	RMSE $\downarrow$						
RS [21]	20	29.22	0.2924	14.63	0.4041	8.73	0.2501	2.05	0.3329
	5,000	36.56	0.0886	14.63	0.4041	11.49	0.2822	2.68	0.3444
BO [21]	20	21.03	0.2830	0.44	0.3653	8.11	0.2402	2.05	0.3024
	5,000	55.09	0.2037	12.17	0.4089	12.68	0.2816	6.69	0.3444
PSO [21]	20	20.90	0.0836	6.40	0.4204	4.05	0.1567	1.06	0.2339
	5,000	21.88	0.0844	6.50	0.4211	4.35	0.1670	1.09	0.2390
P-cG-SchNet [53]	1	48.22	0.4179	17.29	0.3846	15.39	0.3762	3.67	0.4115
	20	97.94	0.3463	55.91	0.3551	32.64	0.3018	12.96	0.3942
CDVAE [9]	1	45.31	0.1138	17.09	0.2969	33.90	0.1045	5.34	0.2106
	20	88.51	0.0464	88.37	0.2286	66.95	0.1026	20.79	0.2085
DiffCSP	1	52.02	0.0760	17.54	0.2759	51.49	0.0631	12.19	0.1786
	20	<b>98.60</b>	<b>0.0128</b>	<b>88.47</b>	<b>0.2192</b>	<b>77.93</b>	<b>0.0492</b>	<b>34.02</b>	<b>0.1749</b>

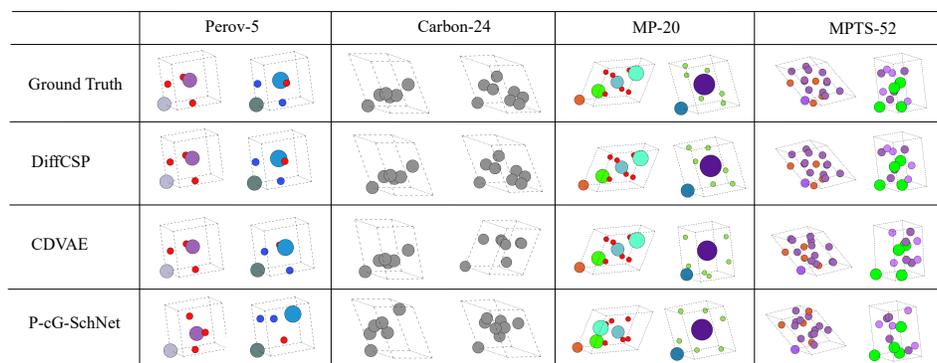


Figure 3: Visualization of the predicted structures from different methods. We select the structure of the lowest RMSE over 20 candidates. We translate the same predicted atom by all methods to the origin for better comparison. Our DiffCSP accurately delivers high quality structure predictions.

## 5 Experiments

In this section, we evaluate the efficacy of DiffCSP on a diverse range of tasks, by showing that it can generate high-quality structures of different crystals in § 5.1, with lower time cost comparing with DFT-based optimization method in § 5.2. Ablations in § 5.3 exhibit the necessity of each designed component. We further showcase the capability of DiffCSP in the ab initio generation task in § 5.4.

### 5.1 Stable Structure Prediction

**Dataset** We conduct experiments on four datasets with distinct levels of difficulty. **Perov-5** [54, 55] contains 18,928 perovskite materials with similar structures. Each structure has 5 atoms in a unit cell. **Carbon-24** [56] includes 10,153 carbon materials with 6~24 atoms in a cell. **MP-20** [57] selects 45,231 stable inorganic materials from Material Projects [57], which includes the majority of experimentally-generated materials with at most 20 atoms in a unit cell. **MPTS-52** is a more challenging extension of MP-20, consisting of 40,476 structures up to 52 atoms per cell, sorted according to the earliest published year in literature. For Perov-5, Carbon-24 and MP-20, we apply the 60-20-20 split in line with Xie et al. [9]. For MPTS-52, we split 27,380/5,000/8,096 for training/validation/testing in chronological order.

**Baselines** We contrast two types of previous works. The first type follows the predict-optimize paradigm, which first trains a predictor of the target property and then utilizes certain optimization algorithms to search for optimal structures. Following Cheng et al. [21], we apply MEGNet [52] as the predictor of the formation energy. For the optimization algorithms, we choose Random Search (**RS**), Bayesian Optimization (**BO**), and Particle Swarm Optimization (**PSO**), all iterated over 5,000 steps. The second type is based on deep generative models. We follow the modification in Xie et al. [9], and leverage cG-SchNet [53] that utilizes SchNet [29] as the backbone and additionally

considers the ground-truth lattice initialization for encoding periodicity, yielding a final model named **P-cG-SchNet**. Another baseline **CDVAE** [9] is a VAE-based framework for pure crystal generation, by first predicting the lattice and the initial composition and then optimizing the atom types and coordinates via annealed Langevin dynamics [49]. To adapt CDVAE into the CSP task, we replace the original normal prior for generation with a parametric prior conditional on the encoding of the given composition. More details are provided in Appendix B.2.

**Evaluation metrics** Following the common practice [9], we evaluate by matching the predicted candidates with the ground-truth structure. Specifically, for each structure in the test set, we first generate  $k$  samples of the same composition and then identify the matching if at least one of the samples matches the ground truth structure, under the metric by the StructureMatcher class in pymatgen [58] with thresholds stol=0.5, angle\_tol=10, ltol=0.3. The **Match rate** is the proportion of the matched structures over the test set. **RMSE** is calculated between the ground truth and the best matching candidate, normalized by  $\sqrt[3]{V/N}$  where  $V$  is the volume of the lattice, and averaged over the matched structures. For optimization methods, we select 20 structures of the lowest energy of all 5,000 structures from all iterations during testing as candidates. For generative baselines and our DiffCSP, we let  $k = 1$  and  $k = 20$  for evaluation. We provide more details in Appendix B, C.1 and I.

**Results** Table 1 conveys the following observations. **1.** The optimization methods encounter low Match rates, signifying the difficulty of locating the optimal structures within the vast search space. **2.** In comparison to other generative methods that construct structures atom by atom or predict the lattice and atom coordinates in two stages, our method demonstrates superior performance, highlighting the effectiveness of jointly refining the lattice and coordinates during generation. **3.** All methods struggle with performance degradation as the number of atoms per cell increases, on the datasets from Perov-5 to MPTS-52. For example, the Match rates of the optimization methods are less than 10% in MPTS-52. Even so, our method consistently outperforms all other methods.

**Visualization** Figure 3 provides qualitative comparisons. DiffCSP clearly makes the best predictions.

## 5.2 Comparison with DFT-based Methods

We further select 10 binary and 5 ternary compounds in MP-20 testing set and compare our model with USPEX [59], a DFT-based software equipped with the evolutionary algorithm to search for stable structures.

Table 2: Overall results over the 15 selected compounds.

	Match rate (%) $\uparrow$	Avg. RMSD $\downarrow$	Avg. Time $\downarrow$
USPEX [59]	53.33	<b>0.0159</b>	12.5h
DiffCSP	<b>73.33</b>	0.0172	<b>10s</b>

For our method, we sample 20 candidates for each compound following the setting in Table 1. We select the model trained on MP-20 for inference, with a training duration of 5.2 hours. For USPEX, we apply 20 generations, 20 populations for each compound, and select the best sample in each generation, leading to 20 candidates as well. We summarize the **Match rate** over the 15 compounds, the **Averaged RMSD** over the matched structures, and the **Averaged Inference Time** to generate 20 candidates for each compound in Table 10. The detailed results for each compound are listed in Appendix F. DiffCSP correctly predicts more structures with higher match rate, and more importantly, its time cost is much less than USPEX, allowing more potential for real applications.

## 5.3 Ablation Studies

We ablate each component of DiffCSP in Table 3, and probe the following aspects. **1.** To verify the necessity of jointly updating the lattice  $\mathbf{L}$  and fractional coordinates  $\mathbf{F}$ , we construct two variants that separate the joint optimization into two stages, denoted as  $\mathbf{L} \rightarrow \mathbf{F}$  and  $\mathbf{F} \rightarrow \mathbf{L}$ . Particularly,  $\mathbf{L} \rightarrow \mathbf{F}$  applies two networks to learn the reverse processes  $p_{\theta_1}(\mathbf{L}_{0:T-1}|\mathbf{A}, \mathbf{F}_T, \mathbf{L}_T)$  and  $p_{\theta_2}(\mathbf{F}_{0:T-1}|\mathbf{A}, \mathbf{F}_T, \mathbf{L}_0)$ . During inference, we first sample  $\mathbf{L}_T, \mathbf{F}_T$  from their prior distributions, acquiring  $\mathbf{L}_0$  via  $p_{\theta_1}$ , and then  $\mathbf{F}_0$  by  $p_{\theta_2}$  based on  $\mathbf{L}_0$ .  $\mathbf{F} \rightarrow \mathbf{L}$  is similarly executed but with the generation order of  $\mathbf{L}_0$  and  $\mathbf{F}_0$  exchanged. Results indicate that  $\mathbf{L} \rightarrow \mathbf{F}$  performs better than the  $\mathbf{F} \rightarrow \mathbf{L}$ , but both are inferior to the joint update in DiffCSP, which endorses our design. We conjecture that the joint diffusion fashion enables  $\mathbf{L}$  and  $\mathbf{F}$  to update synergistically, which makes the generation process more tractable to learn and thus leads to better performance. **2.** We explore the necessity of preserving the  $O(3)$  invariance when generating  $\mathbf{L}$ , which is ensured by the inner product  $\mathbf{L}^\top \mathbf{L}$  in Eq. (6).

When we replace it with  $\mathbf{L}$  and change the final output as  $\hat{\mathbf{e}}_{\mathbf{L}} = \varphi_{\mathbf{L}}\left(\frac{1}{N} \sum_{i=1}^N \mathbf{h}_i^{(S)}\right)$  in Eq. (9) to break the equivariance, the model suffers from extreme performance detriment. Only 1.66% structures are successfully matched, which obviously implies the importance of incorporating  $O(3)$  equivariance. Furthermore, we introduce the chirality into the denoising model by adding  $\text{sign}(|\mathbf{L}|)$ , the sign of the determinant of the lattice matrix, as an additional input in Eq. 6. The adapted model is  $SO(3)$ -invariant, but breaks the reflection symmetry and hence is NOT  $O(3)$ -invariant. There is no essential performance change, indicating the chirality is not quite crucial in distinguishing different crystal structures for the datasets used in this paper. **3.** We further assess the importance of periodic translation invariance from two perspectives. For the generation process, we generate  $\mathbf{F}$  via the score-based model with the Wrapped Normal (WN) distribution. We replace this module with DDPM under standard Gaussian as  $q(\mathbf{F}_t|\mathcal{M}_0) = \mathcal{N}(\mathbf{F}_t|\sqrt{\bar{\alpha}_t}\mathbf{F}_0, (1 - \bar{\alpha}_t)\mathbf{I})$  similarly defined as Eq. (3). A lower match rate and higher RMSE are observed for this variant. For the model architecture, we adopt Fourier Transformation(FT) in Eq. (6) to capture periodicity. To investigate its effect, we replace  $\psi_{\text{FT}}(\mathbf{f}_j - \mathbf{f}_i)$  with  $\mathbf{f}_j - \mathbf{f}_i$ , and the match rate drops from 51.49% to 29.15%. Both of the two observations verify the importance of retaining the periodic translation invariance. **4.** We further change the fully-connected graph into the multi-graph approach adopted in Xie and Grossman [15]. The multi-graph approach decreases the match rate, since the multi-graphs constructed under different intermediate structures may differ vibrantly during generation, leading to substantially higher training difficulty and lower sampling stability. We will provide more discussions in Appendix E.

#### 5.4 Ab Initio Crystal Generation

DiffCSP is extendable to ab initio crystal generation by further conducting discrete diffusion on atom types  $\mathbf{A}$ . We contrast DiffCSP against five generative methods following [9]: **FTCP** [28], **Cond-DFC-VAE** [7], **G-SchNet** [60] and its periodic variant **P-G-SchNet**, and the original version of **CDVAE** [9]. Specifically for our DiffCSP, we gather the statistics of the atom numbers from the training set, then sample the number based on the pre-computed distribution similar to Hoogeboom et al. [41], which allows DiffCSP to generate structures of variable size. Following [9], we evaluate the generation performance in terms of these metrics: **Validity**, **Coverage**, and **Property statistics**, which respectively return the validity of the predicted crystals, the similarity between the test set and the generated samples, and the property calculation regarding density, formation energy, and the number of elements. The detailed definitions of the above metrics are provided in Appendix G.

**Results** Table 4 show that our method achieves comparable validity and coverage rate with previous methods, and significantly outperforms the baselines on the similarity of property statistics, which indicates the high reliability of the generated samples.

## 6 Discussions

**Limitation 1.** Composition generation. Our model yields slightly lower compositional validity in Table 4. We provide more discussion in Appendix G, and it is promising to propose more powerful generation methods on atom types. **2.** Experimental evaluation. Further wet-lab experiments can better verify the effectiveness of the model in real applications.

**Conclusion** In this work, we present DiffCSP, a diffusion-based learning framework for crystal structure prediction, which is particularly curated with the vital symmetries existing in crystals. The diffusion is highly flexible by jointly optimizing the lattice and fractional coordinates, where

Table 3: Ablation studies on MP-20. MG: Multi-Graph edge construction [15], FT: Fourier-Transformation.

	Match rate (%) $\uparrow$	RMSE $\downarrow$
DiffCSP	<b>51.49</b>	<b>0.0631</b>
<i>w/o Joint Diffusion</i>		
$\mathbf{L} \rightarrow \mathbf{F}$	50.03	0.0921
$\mathbf{F} \rightarrow \mathbf{L}$	36.73	0.0838
<i>w/o <math>O(3)</math> Equivariance</i>		
w/o inner product	1.66	0.4002
w/ chirality	49.68	0.0637
<i>w/o Periodic Translation Invariance</i>		
w/o WN	34.09	0.2350
w/o FT	29.15	0.0926
<i>MG Edge Construction</i>		
MG w/ FT	25.85	0.1079
MG w/o FT	28.05	0.1314

<sup>3</sup>Composition-based metrics are not meaningful for Carbon-24, as all structures are composed of carbon.

Table 4: Results on ab initio generation task. The results of baseline methods are from Xie et al. [9].

Data	Method	Validity (%) $\uparrow$		Coverage (%) $\uparrow$		Property $\downarrow$		
		Struc.	Comp.	COV-R	COV-P	$d_\rho$	$d_E$	$d_{elem}$
Perov-5	FTCP [28]	0.24	54.24	0.00	0.00	10.27	156.0	0.6297
	Cond-DFC-VAE [7]	73.60	82.95	73.92	10.13	2.268	4.111	0.8373
	G-SchNet [60]	99.92	98.79	0.18	0.23	1.625	4.746	0.0368
	P-G-SchNet [60]	79.63	<b>99.13</b>	0.37	0.25	0.2755	1.388	0.4552
	CDVAE [9]	<b>100.0</b>	98.59	99.45	<b>98.46</b>	0.1258	0.0264	0.0628
	DiffCSP	<b>100.0</b>	98.85	<b>99.74</b>	98.27	<b>0.1110</b>	<b>0.0263</b>	<b>0.0128</b>
Carbon-24 <sup>3</sup>	FTCP [28]	0.08	–	0.00	0.00	5.206	19.05	–
	G-SchNet [60]	99.94	–	0.00	0.00	0.9427	1.320	–
	P-G-SchNet [60]	48.39	–	0.00	0.00	1.533	134.7	–
	CDVAE [9]	<b>100.0</b>	–	99.80	83.08	0.1407	0.2850	–
	DiffCSP	<b>100.0</b>	–	<b>99.90</b>	<b>97.27</b>	<b>0.0805</b>	<b>0.0820</b>	–
MP-20	FTCP [28]	1.55	48.37	4.72	0.09	23.71	160.9	0.7363
	G-SchNet [60]	99.65	75.96	38.33	99.57	3.034	42.09	0.6411
	P-G-SchNet [60]	77.51	76.40	41.93	99.74	4.04	2.448	0.6234
	CDVAE [9]	<b>100.0</b>	<b>86.70</b>	99.15	99.49	0.6875	0.2778	1.432
	DiffCSP	<b>100.0</b>	83.25	<b>99.71</b>	<b>99.76</b>	<b>0.3502</b>	<b>0.1247</b>	<b>0.3398</b>

the intermediate distributions are guaranteed to be invariant under necessary transformations. We demonstrate the efficacy of our approach on a wide range of crystal datasets, verifying the strong applicability of DiffCSP towards predicting high-quality crystal structures.

## 7 Acknowledgement

This work was jointly supported by the following projects: the National Natural Science Foundation of China (61925601 & 62006137); Beijing Nova Program (20230484278); the Fundamental Research Funds for the Central Universities, and the Research Funds of Renmin University of China (23XNKJ19); Alibaba Damo Research Fund; CCF-Ant Research Fund (CCF-AFSG RF20220204); National Key R&D Program of China (2022ZD0117805).

## References

- [1] Gautam R Desiraju. Cryptic crystallography. *Nature materials*, 1(2):77–79, 2002.
- [2] Keith T Butler, Daniel W Davies, Hugh Cartwright, Olexandr Isayev, and Aron Walsh. Machine learning for molecular and materials science. *Nature*, 559(7715):547–555, 2018.
- [3] Walter Kohn and Lu Jeu Sham. Self-consistent equations including exchange and correlation effects. *Physical review*, 140(4A):A1133, 1965.
- [4] Chris J Pickard and RJ Needs. Ab initio random structure searching. *Journal of Physics: Condensed Matter*, 23(5):053201, 2011.
- [5] Tomoki Yamashita, Nobuya Sato, Hiori Kino, Takashi Miyake, Koji Tsuda, and Tamio Oguchi. Crystal structure prediction accelerated by bayesian optimization. *Physical Review Materials*, 2(1):013803, 2018.
- [6] Artem R Oganov, Chris J Pickard, Qiang Zhu, and Richard J Needs. Structure prediction drives materials discovery. *Nature Reviews Materials*, 4(5):331–348, 2019.
- [7] Callum J Court, Batuhan Yildirim, Apoorv Jain, and Jacqueline M Cole. 3-d inorganic crystal structure generation and property prediction via representation learning. *Journal of chemical information and modeling*, 60(10):4518–4535, 2020.
- [8] Wenhui Yang, Edirisuriya M Dilanga Siriwardane, Rongzhi Dong, Yuxin Li, and Jianjun Hu. Crystal structure prediction of materials with high symmetry using differential evolution. *Journal of Physics: Condensed Matter*, 33(45):455902, 2021.
- [9] Tian Xie, Xiang Fu, Octavian-Eugen Ganea, Regina Barzilay, and Tommi S Jaakkola. Crystal diffusion variational autoencoder for periodic material generation. In *International Conference on Learning Representations*, 2021.

- [10] Minkai Xu, Lantao Yu, Yang Song, Chence Shi, Stefano Ermon, and Jian Tang. Geodiff: A geometric diffusion model for molecular conformation generation. In *International Conference on Learning Representations*, 2021.
- [11] Brian L. Trippe, Jason Yim, Doug Tischer, David Baker, Tamara Broderick, Regina Barzilay, and Tommi S. Jaakkola. Diffusion probabilistic modeling of protein backbones in 3d for the motif-scaffolding problem. In *The Eleventh International Conference on Learning Representations*, 2023.
- [12] Gabriele Corso, Hannes Stärk, Bowen Jing, Regina Barzilay, and Tommi S. Jaakkola. Diffdock: Diffusion steps, twists, and turns for molecular docking. In *The Eleventh International Conference on Learning Representations*, 2023.
- [13] Chence Shi, Shitong Luo, Minkai Xu, and Jian Tang. Learning gradient fields for molecular conformation generation. In *International Conference on Machine Learning*, pages 9558–9568. PMLR, 2021.
- [14] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A A Kohl, Andrew J Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Sebastian Bodenstern, David Silver, Oriol Vinyals, Andrew W Senior, Koray Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate protein structure prediction with AlphaFold. *Nature*, 596(7873):583–589, 2021. doi: 10.1038/s41586-021-03819-2.
- [15] Tian Xie and Jeffrey C. Grossman. Crystal graph convolutional neural networks for an accurate and interpretable prediction of material properties. *Phys. Rev. Lett.*, 120:145301, Apr 2018. doi: 10.1103/PhysRevLett.120.145301.
- [16] Bowen Jing, Gabriele Corso, Jeffrey Chang, Regina Barzilay, and Tommi Jaakkola. Torsional diffusion for molecular conformer generation. *arXiv preprint arXiv:2206.01729*, 2022.
- [17] Yanchao Wang, Jian Lv, Li Zhu, and Yanming Ma. Crystal structure prediction via particle-swarm optimization. *Physical Review B*, 82(9):094116, 2010.
- [18] Yunwei Zhang, Hui Wang, Yanchao Wang, Lijun Zhang, and Yanming Ma. Computer-assisted inverse design of inorganic electrides. *Physical Review X*, 7(1):011017, 2017.
- [19] TL Jacobsen, MS Jørgensen, and B Hammer. On-the-fly machine learning of atomic potential in density functional theory structure optimization. *Physical review letters*, 120(2):026102, 2018.
- [20] Evgeny V Podryabinkin, Evgeny V Tikhonov, Alexander V Shapeev, and Artem R Oganov. Accelerating crystal structure prediction by machine-learning interatomic potentials with active learning. *Physical Review B*, 99(6):064114, 2019.
- [21] Guanjian Cheng, Xin-Gao Gong, and Wan-Jian Yin. Crystal structure prediction by combining graph network and optimization algorithm. *Nature communications*, 13(1):1–8, 2022.
- [22] Jordan Hoffmann, Louis Maestrati, Yoshihide Sawada, Jian Tang, Jean Michel Sellier, and Yoshua Bengio. Data-driven approach to encoding and decoding 3-d crystal structures. *arXiv preprint arXiv:1909.00949*, 2019.
- [23] Juhwan Noh, Jaehoon Kim, Helge S Stein, Benjamin Sanchez-Lengeling, John M Gregoire, Alan Aspuru-Guzik, and Yousung Jung. Inverse design of solid-state materials via a continuous representation. *Matter*, 1(5):1370–1384, 2019.
- [24] Jianjun Hu, Wenhui Yang, and Edirisuriya M Dilanga Siriwardane. Distance matrix-based crystal structure prediction using evolutionary algorithms. *The Journal of Physical Chemistry A*, 124(51):10909–10919, 2020.
- [25] Jianjun Hu, Wenhui Yang, Rongzhi Dong, Yuxin Li, Xiang Li, Shaobo Li, and Edirisuriya MD Siriwardane. Contact map based crystal structure prediction using global optimization. *CryStEngComm*, 23(8):1765–1776, 2021.

- [26] Asma Noura, Nataliya Sokolovska, and Jean-Claude Crivello. Crystalgan: learning to discover crystallographic structures with generative adversarial networks. *arXiv preprint arXiv:1810.11203*, 2018.
- [27] Sungwon Kim, Juhwan Noh, Geun Ho Gu, Alan Aspuru-Guzik, and Yousung Jung. Generative adversarial networks for crystal structure prediction. *ACS central science*, 6(8):1412–1420, 2020.
- [28] Zekun Ren, Siyu Isaac Parker Tian, Juhwan Noh, Felipe Oviedo, Guangzong Xing, Jiali Li, Qiaohao Liang, Ruiming Zhu, Armin G. Aberle, Shijing Sun, Xiaonan Wang, Yi Liu, Qianxiao Li, Senthilnath Jayavelu, Kedar Hippalgaonkar, Yousung Jung, and Tonio Buonassisi. An invertible crystallographic representation for general inverse design of inorganic crystals with targeted properties. *Matter*, 2021. ISSN 2590-2385. doi: <https://doi.org/10.1016/j.matt.2021.11.032>.
- [29] Kristof T Schütt, Huziel E Sauceda, P-J Kindermans, Alexandre Tkatchenko, and K-R Müller. SchNet—a deep learning architecture for molecules and materials. *The Journal of Chemical Physics*, 148(24):241722, 2018.
- [30] Nathaniel Thomas, Tess Smidt, Steven Kearnes, Lusann Yang, Li Li, Kai Kohlhoff, and Patrick Riley. Tensor field networks: Rotation-and translation-equivariant neural networks for 3d point clouds. *arXiv preprint arXiv:1802.08219*, 2018.
- [31] Fabian Fuchs, Daniel E. Worrall, Volker Fischer, and Max Welling. Se(3)-transformers: 3d rotation-translation equivariant attention networks. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- [32] Victor Garcia Satorras, Emiel Hoogeboom, and Max Welling. E (n) equivariant graph neural networks. In *International Conference on Machine Learning*, pages 9323–9332. PMLR, 2021.
- [33] Philipp Thölke and Gianni De Fabritiis. Equivariant transformers for neural network based molecular potentials. In *International Conference on Learning Representations*, 2021.
- [34] Lowik Chanussot, Abhishek Das, Siddharth Goyal, Thibaut Lavril, Muhammed Shuaibi, Morgane Riviere, Kevin Tran, Javier Heras-Domingo, Caleb Ho, Weihua Hu, Aini Palizhati, Anuroop Sriram, Brandon Wood, Junwoong Yoon, Devi Parikh, C. Lawrence Zitnick, and Zachary Ulissi. Open catalyst 2020 (oc20) dataset and community challenges. *ACS Catalysis*, 2021. doi: 10.1021/acscatal.0c04525.
- [35] Richard Tran, Janice Lan, Muhammed Shuaibi, Siddharth Goyal, Brandon M Wood, Abhishek Das, Javier Heras-Domingo, Adeesh Kolluru, Ammar Rizvi, Nima Shoghi, et al. The open catalyst 2022 (oc22) dataset and challenges for oxide electrocatalysis. *arXiv preprint arXiv:2206.08917*, 2022.
- [36] Keqiang Yan, Yi Liu, Yuchao Lin, and Shuiwang Ji. Periodic graph transformers for crystal material property prediction. In *The 36th Annual Conference on Neural Information Processing Systems*, 2022.
- [37] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015.
- [38] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- [39] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021.
- [40] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.

- [41] Emiel Hoogeboom, Victor Garcia Satorras, Clément Vignac, and Max Welling. Equivariant diffusion for molecule generation in 3d. In *International Conference on Machine Learning*, pages 8867–8887. PMLR, 2022.
- [42] Shitong Luo, Yufeng Su, Xingang Peng, Sheng Wang, Jian Peng, and Jianzhu Ma. Antigen-specific antibody design and optimization with diffusion-based generative models for protein structures. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [43] Valentin De Bortoli, Emile Mathieu, Michael John Hutchinson, James Thornton, Yee Whye Teh, and Arnaud Doucet. Riemannian score-based generative modelling. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [44] Chin-Wei Huang, Milad Aghajohari, Joey Bose, Prakash Panangaden, and Aaron Courville. Riemannian diffusion models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [45] Detlef WM Hofmann and Joannis Apostolakis. Crystal structure prediction by data mining. *Journal of Molecular Structure*, 647(1-3):17–39, 2003.
- [46] Ralf W Grosse-Kunstleve, Nicholas K Sauter, and Paul D Adams. Numerically stable algorithms for the computation of reduced unit cells. *Acta Crystallographica Section A: Foundations of Crystallography*, 60(1):1–6, 2004.
- [47] Thomas N Kipf and Max Welling. Variational graph auto-encoders. *arXiv preprint arXiv:1611.07308*, 2016.
- [48] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021.
- [49] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33:12438–12448, 2020.
- [50] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- [51] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [52] Chi Chen, Weike Ye, Yunxing Zuo, Chen Zheng, and Shyue Ping Ong. Graph networks as a universal machine learning framework for molecules and crystals. *Chemistry of Materials*, 31(9):3564–3572, 2019.
- [53] Niklas WA Gebauer, Michael Gastegger, Stefaan SP Hessmann, Klaus-Robert Müller, and Kristof T Schütt. Inverse design of 3d molecular structures with conditional generative neural networks. *Nature communications*, 13(1):1–11, 2022.
- [54] Ivano E Castelli, David D Landis, Kristian S Thygesen, Søren Dahl, Ib Chorkendorff, Thomas F Jaramillo, and Karsten W Jacobsen. New cubic perovskites for one-and two-photon water splitting using the computational materials repository. *Energy & Environmental Science*, 5(10):9034–9043, 2012.
- [55] Ivano E Castelli, Thomas Olsen, Soumendu Datta, David D Landis, Søren Dahl, Kristian S Thygesen, and Karsten W Jacobsen. Computational screening of perovskite metal oxides for optimal solar light capture. *Energy & Environmental Science*, 5(2):5814–5819, 2012.
- [56] Chris J. Pickard. Airss data for carbon at 10gpa and the c+n+h+o system at 1gpa, 2020.
- [57] Anubhav Jain, Shyue Ping Ong, Geoffroy Hautier, Wei Chen, William Davidson Richards, Stephen Dacek, Shreyas Cholia, Dan Gunter, David Skinner, Gerbrand Ceder, et al. Commentary: The materials project: A materials genome approach to accelerating materials innovation. *APL materials*, 1(1):011002, 2013.

- [58] Shyue Ping Ong, William Davidson Richards, Anubhav Jain, Geoffroy Hautier, Michael Kocher, Shreyas Cholia, Dan Gunter, Vincent L Chevrier, Kristin A Persson, and Gerbrand Ceder. Python materials genomics (pymatgen): A robust, open-source python library for materials analysis. *Computational Materials Science*, 68:314–319, 2013.
- [59] Colin W Glass, Artem R Oganov, and Nikolaus Hansen. Uspex—evolutionary crystal structure prediction. *Computer physics communications*, 175(11-12):713–720, 2006.
- [60] Niklas Gebauer, Michael Gastegger, and Kristof Schütt. Symmetry-adapted generation of 3d point sets for the targeted discovery of molecules. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 7566–7578. Curran Associates, Inc., 2019.
- [61] Gerhard Kurz, Igor Gilitschenski, and Uwe D Hanebeck. Efficient evaluation of the probability density function of a wrapped normal distribution. In *2014 Sensor Data Fusion: Trends, Solutions, Applications (SDF)*, pages 1–5. IEEE, 2014.
- [62] James Bergstra, Daniel Yamins, and David Cox. Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures. In *International conference on machine learning*, pages 115–123. PMLR, 2013.
- [63] Johannes Gasteiger, Shankari Giri, Johannes T. Margraf, and Stephan Günnemann. Fast and uncertainty-aware directional message passing for non-equilibrium molecules. In *Machine Learning for Molecules Workshop, NeurIPS*, 2020.
- [64] Johannes Gasteiger, Florian Becker, and Stephan Günnemann. Gemnet: Universal directional graph neural networks for molecules. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2021.
- [65] Nils ER Zimmermann and Anubhav Jain. Local structure order parameters and site fingerprints for quantification of coordination environment and crystal structure similarity. *RSC advances*, 10(10):6063–6081, 2020.
- [66] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [67] Peter E Blöchl. Projector augmented-wave method. *Physical review B*, 50(24):17953, 1994.
- [68] Georg Kresse and Jürgen Furthmüller. Efficiency of ab-initio total energy calculations for metals and semiconductors using a plane-wave basis set. *Computational materials science*, 6(1):15–50, 1996.
- [69] John P Perdew, Kieron Burke, and Matthias Ernzerhof. Generalized gradient approximation made simple. *Physical review letters*, 77(18):3865, 1996.
- [70] Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows and multinomial diffusion: Learning categorical distributions. *Advances in Neural Information Processing Systems*, 34:12454–12465, 2021.
- [71] Clement Vignac, Igor Krawczuk, Antoine Siraudin, Bohan Wang, Volkan Cevher, and Pascal Frossard. Digress: Discrete denoising diffusion for graph generation. In *The Eleventh International Conference on Learning Representations*, 2023.
- [72] Daniel W Davies, Keith T Butler, Adam J Jackson, Jonathan M Skelton, Kazuki Morita, and Aron Walsh. Smact: Semiconducting materials by analogy and chemical theory. *Journal of Open Source Software*, 4(38):1361, 2019.
- [73] Logan Ward, Ankit Agrawal, Alok Choudhary, and Christopher Wolverton. A general-purpose machine learning framework for predicting properties of inorganic materials. *npj Computational Materials*, 2(1):1–7, 2016.

- [74] Min Zhao, Fan Bao, Chongxuan Li, and Jun Zhu. EGSDE: Unpaired image-to-image translation via energy-guided stochastic differential equations. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [75] Fan Bao, Min Zhao, Zhongkai Hao, Peiyao Li, Chongxuan Li, and Jun Zhu. Equivariant energy-guided SDE for inverse molecular design. In *The Eleventh International Conference on Learning Representations*, 2023.

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Related Works</b>	<b>2</b>
<b>3</b>	<b>Preliminaries</b>	<b>3</b>
<b>4</b>	<b>The Proposed Method: DiffCSP</b>	<b>3</b>
4.1	Symmetries of Crystal Structure Distribution . . . . .	3
4.2	Joint Equivariant Diffusion . . . . .	4
4.3	The Architecture of the Denoising Model . . . . .	6
<b>5</b>	<b>Experiments</b>	<b>7</b>
5.1	Stable Structure Prediction . . . . .	7
5.2	Comparison with DFT-based Methods . . . . .	8
5.3	Ablation Studies . . . . .	8
5.4	Ab Initio Crystal Generation . . . . .	9
<b>6</b>	<b>Discussions</b>	<b>9</b>
<b>7</b>	<b>Acknowledgement</b>	<b>10</b>
<b>A</b>	<b>Theoretical Analysis</b>	<b>18</b>
A.1	Proof of Proposition 1 . . . . .	18
A.2	Proof of Proposition 2 . . . . .	19
A.3	Proof of Proposition 3 . . . . .	20
A.4	Discussion on Periodic Translation Invariance . . . . .	21
<b>B</b>	<b>Implementation Details</b>	<b>22</b>
B.1	Approximation of the Wrapped Normal Distribution . . . . .	22
B.2	Adaptation of CDVAE . . . . .	22
B.3	Algorithms for Training and Sampling . . . . .	23
B.4	Hyper-parameters and Training Details . . . . .	24
<b>C</b>	<b>Exploring the Effects of Sampling and Candidate Ranking</b>	<b>25</b>
C.1	Impact of Sampling Numbers . . . . .	25
C.2	Diversity . . . . .	25
C.3	Ranking among Multiple Candidates . . . . .	25
<b>D</b>	<b>Impact of Noise Schedulers</b>	<b>26</b>
<b>E</b>	<b>Learning Curves of Different Variants</b>	<b>26</b>
<b>F</b>	<b>Comparison with DFT-based Methods</b>	<b>26</b>

F.1	Implementation Details . . . . .	26
F.2	Results . . . . .	28
<b>G</b>	<b>Extension to More General Tasks</b>	<b>28</b>
G.1	Overview . . . . .	28
G.2	Ab Initio Generation . . . . .	28
G.3	Property Optimization . . . . .	30
<b>H</b>	<b>Computational Cost for Inference</b>	<b>32</b>
<b>I</b>	<b>Error Bars</b>	<b>32</b>
<b>J</b>	<b>More Visualizations</b>	<b>32</b>

## A Theoretical Analysis

### A.1 Proof of Proposition 1

We first introduce the following definition to describe the equivariance and invariance from the perspective of distributions.

**Definition 4.** We call a distribution  $p(x)$  is  $G$ -invariant if for any transformation  $g$  in the group  $G$ ,  $p(g \cdot x) = p(x)$ , and a conditional distribution  $p(x|c)$  is  $G$ -equivariant if  $p(g \cdot x|g \cdot c) = p(x|c)$ ,  $\forall g \in G$ .

We then provide and prove the following lemma to capture the symmetry of the generation process.

**Lemma 1** (Xu et al. [10]). Consider the generation Markov process  $p(x_0) = p(x_T) \int p(x_{0:T-1}|x_t) dx_{1:T}$ . If the prior distribution  $p(x_T)$  is  $G$ -invariant and the Markov transitions  $p(x_{t-1}|x_t)$ ,  $0 < t \leq T$  are  $G$ -equivariant, the marginal distribution  $p(x_0)$  is also  $G$ -invariant.

*Proof.* For any  $g \in G$ , we have

$$\begin{aligned} p(g \cdot x_0) &= p(g \cdot x_T) \int p(g \cdot x_{0:T-1}|g \cdot x_t) dx_{1:T} \\ &= p(g \cdot x_T) \int \prod_{t=1}^T p(g \cdot x_{t-1}|g \cdot x_t) dx_{1:T} \\ &= p(x_T) \int \prod_{t=1}^T p(g \cdot x_{t-1}|g \cdot x_t) dx_{1:T} \\ &= p(x_T) \int \prod_{t=1}^T p(x_{t-1}|x_t) dx_{1:T} \\ &= p(x_T) \int p(x_{0:T-1}|x_t) dx_{1:T} \\ &= p(x_0). \end{aligned}$$

Hence, the marginal distribution  $p(x_0)$  is  $G$ -invariant. □

The proposition Proposition 1 is rewritten and proved as follows.

**Proposition 1.** The marginal distribution  $p(\mathbf{L}_0)$  by Eq. (3) is  $O(3)$ -invariant if  $\hat{\epsilon}_{\mathbf{L}}(\mathcal{M}_t, t)$  is  $O(3)$ -equivariant, namely  $\hat{\epsilon}_{\mathbf{L}}(\mathbf{Q}\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t) = \mathbf{Q}\hat{\epsilon}_{\mathbf{L}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t)$ ,  $\forall \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$ .

*Proof.* Consider the transition probability in Eq. (3), we have

$$p(\mathbf{L}_{t-1}|\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}) = \mathcal{N}(\mathbf{L}_{t-1}|a_t(\mathbf{L}_t - b_t\hat{\epsilon}_{\mathbf{L}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t)), \sigma_t^2 \mathbf{I}),$$

where  $a_t = \frac{1}{\sqrt{\alpha_t}}$ ,  $b_t = \frac{\beta_t}{\sqrt{1-\alpha_t}}$ ,  $\sigma_t^2 = \beta_t \cdot \frac{1-\alpha_{t-1}}{1-\alpha_t}$  for simplicity, and  $\hat{\epsilon}_{\mathbf{L}}(\mathcal{M}_t, t)$  is completed as  $\hat{\epsilon}_{\mathbf{L}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t)$ .

As the denoising term  $\hat{\epsilon}_{\mathbf{L}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t)$  is  $O(3)$ -equivariant, we have  $\hat{\epsilon}_{\mathbf{L}}(\mathbf{Q}\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t) = \mathbf{Q}\hat{\epsilon}_{\mathbf{L}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t)$  for any orthogonal transformation  $\mathbf{Q} \in \mathbb{R}^{3 \times 3}$ ,  $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$ .

For the variable  $\mathbf{L} \sim \mathcal{N}(\bar{\mathbf{L}}, \sigma^2 \mathbf{I})$ , we have  $\mathbf{Q}\mathbf{L} \sim \mathcal{N}(\mathbf{Q}\bar{\mathbf{L}}, \mathbf{Q}(\sigma^2 \mathbf{I})\mathbf{Q}^\top) = \mathcal{N}(\mathbf{Q}\bar{\mathbf{L}}, \sigma^2 \mathbf{I})$ . That is,

$$\mathcal{N}(\mathbf{L}|\bar{\mathbf{L}}, \sigma^2 \mathbf{I}) = \mathcal{N}(\mathbf{Q}\mathbf{L}|\mathbf{Q}\bar{\mathbf{L}}, \sigma^2 \mathbf{I}). \quad (11)$$

For the transition probability  $p(\mathbf{L}_{t-1}|\mathbf{L}_t, \mathbf{F}_t, \mathbf{A})$ , we have

$$\begin{aligned}
 p(\mathbf{Q}\mathbf{L}_{t-1}|\mathbf{Q}\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}) &= \mathcal{N}(\mathbf{Q}\mathbf{L}_{t-1}|a_t(\mathbf{Q}\mathbf{L}_t - b_t\hat{\epsilon}_{\mathbf{L}}(\mathbf{Q}\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t)), \sigma_t^2\mathbf{I}) \\
 &= \mathcal{N}(\mathbf{Q}\mathbf{L}_{t-1}|a_t(\mathbf{Q}\mathbf{L}_t - b_t\mathbf{Q}\hat{\epsilon}_{\mathbf{L}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t)), \sigma_t^2\mathbf{I}) \\
 &\quad \text{(O(3)-equivariant } \hat{\epsilon}_{\mathbf{L}}) \\
 &= \mathcal{N}(\mathbf{Q}\mathbf{L}_{t-1}|\mathbf{Q}\left(a_t(\mathbf{L}_t - b_t\hat{\epsilon}_{\mathbf{L}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t))\right), \sigma_t^2\mathbf{I}) \\
 &= \mathcal{N}(\mathbf{L}_{t-1}|a_t(\mathbf{L}_t - b_t\hat{\epsilon}_{\mathbf{L}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t)), \sigma_t^2\mathbf{I}) \quad \text{(Eq. (11))} \\
 &= p(\mathbf{L}_{t-1}|\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}).
 \end{aligned}$$

As the transition is  $O(3)$ -equivariant and the prior distribution  $\mathcal{N}(0, \mathbf{I})$  is  $O(3)$ -invariant, we prove that the the marginal distribution  $p(\mathbf{L}_0)$  is  $O(3)$ -invariant based on lemma 1.  $\square$

## A.2 Proof of Proposition 2

Let  $\mathcal{N}_w(\mu, \sigma^2\mathbf{I})$  denote the wrapped normal distribution with mean  $\mu$ , variance  $\sigma^2$  and period 1. We first provide the following lemma.

**Lemma 3.** *If the denoising term  $\hat{\epsilon}_{\mathbf{F}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t)$  is periodic translation invariant, and the transition probability can be formulated as  $p(\mathbf{F}_{t-1}|\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}) = \mathcal{N}_w(\mathbf{F}_{t-1}|\mathbf{F}_t + u_t\hat{\epsilon}_{\mathbf{F}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t), v_t^2\mathbf{I})$ , where  $u_t, v_t$  are functions of  $t$ , the transition is periodic translation equivariant.*

*Proof.* As the denoising term  $\hat{\epsilon}_{\mathbf{F}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t)$  is periodic translation invariant (for short PTI), we have  $\hat{\epsilon}_{\mathbf{F}}(\mathbf{L}_t, w(\mathbf{F}_t + \mathbf{t}\mathbf{1}^\top), \mathbf{A}, t) = \hat{\epsilon}_{\mathbf{F}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t)$ , for any translation  $\mathbf{t} \in \mathbb{R}^3$ .

For the wrapping function  $w(\cdot)$ , we have

$$w(a + b) = w(w(a) + b), \forall a, b \in \mathbb{R}. \quad (12)$$

For wrapped normal distribution  $\mathcal{N}_w(\mu, \sigma^2)$  with mean  $\mu$ , variance  $\sigma^2$  and period 1, and for any  $k', k'' \in \mathbb{Z}$ , we have

$$\begin{aligned}
 \mathcal{N}_w(x + k'|\mu + k'', \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma}} \sum_{k=-\infty}^{\infty} \exp\left(-\frac{(x + k' - (\mu + k'') - k)^2}{2\sigma^2}\right) \\
 &= \frac{1}{\sqrt{2\pi\sigma}} \sum_{m=-\infty}^{\infty} \exp\left(-\frac{(x - \mu - m)^2}{2\sigma^2}\right) \quad (m = k - k' + k'') \\
 &= \mathcal{N}_w(x|\mu, \sigma^2)
 \end{aligned}$$

Let  $k' = 0, k'' = w(\mu) - \mu$ , we directly have

$$\mathcal{N}_w(x|w(\mu), \sigma^2) = \mathcal{N}_w(x|\mu, \sigma^2). \quad (13)$$

For any  $t \in \mathbb{R}$ , we have

$$\begin{aligned}
 \mathcal{N}_w(x + t|\mu + t, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma}} \sum_{k=-\infty}^{\infty} \exp\left(-\frac{(x + t - (\mu + t) - k)^2}{2\sigma^2}\right) \\
 &= \frac{1}{\sqrt{2\pi\sigma}} \sum_{k=-\infty}^{\infty} \exp\left(-\frac{(x - \mu - k)^2}{2\sigma^2}\right) \\
 &= \mathcal{N}_w(x|\mu, \sigma^2).
 \end{aligned}$$

Let  $k' = w(x + t) - (x + t), k'' = w(\mu + t) - (\mu + t)$ , we have

$$\mathcal{N}_w(w(x + t)|w(\mu + t), \sigma^2) = \mathcal{N}_w(x + t|\mu + t, \sigma^2) = \mathcal{N}_w(x|\mu, \sigma^2). \quad (14)$$

For the transition probability  $p(\mathbf{F}_{t-1}|\mathbf{L}_t, \mathbf{F}_t, \mathbf{A})$ , we have

$$\begin{aligned} & p(w(\mathbf{F}_{t-1} + \mathbf{t})|\mathbf{L}_t, w(\mathbf{F}_t + \mathbf{t}\mathbf{1}^\top), \mathbf{A}) \\ &= \mathcal{N}_w(w(\mathbf{F}_{t-1} + \mathbf{t}\mathbf{1}^\top)|w(\mathbf{F}_t + \mathbf{t}\mathbf{1}^\top) + u_t \hat{\epsilon}_{\mathbf{F}}(\mathbf{L}_t, w(\mathbf{F}_t + \mathbf{t}\mathbf{1}^\top), \mathbf{A}, t), v_t^2 \mathbf{I}) \\ &= \mathcal{N}_w(w(\mathbf{F}_{t-1} + \mathbf{t}\mathbf{1}^\top)|w(\mathbf{F}_t + \mathbf{t}\mathbf{1}^\top) + u_t \hat{\epsilon}_{\mathbf{F}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t), v_t^2 \mathbf{I}) \end{aligned} \quad (\text{PTI } \hat{\epsilon}_{\mathbf{F}})$$

$$= \mathcal{N}_w(w(\mathbf{F}_{t-1} + \mathbf{t}\mathbf{1}^\top)|w(w(\mathbf{F}_t + \mathbf{t}\mathbf{1}^\top) + u_t \hat{\epsilon}_{\mathbf{F}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t)), v_t^2 \mathbf{I}) \quad (\text{Eq. (13)})$$

$$= \mathcal{N}_w(w(\mathbf{F}_{t-1} + \mathbf{t}\mathbf{1}^\top)|w(\mathbf{F}_t + u_t \hat{\epsilon}_{\mathbf{F}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t) + \mathbf{t}), v_t^2 \mathbf{I}) \quad (\text{Eq. (12)})$$

$$= \mathcal{N}_w(\mathbf{F}_{t-1}|\mathbf{F}_t + u_t \hat{\epsilon}_{\mathbf{F}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t), v_t^2 \mathbf{I}) \quad (\text{Eq. (14)})$$

$$= p(\mathbf{F}_{t-1}|\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}).$$

□

The transition probability of the fractional coordinates during the Predictor-Corrector sampling can be formulated as

$$p(\mathbf{F}_{t-1}|\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}) = p_P(\mathbf{F}_{t-\frac{1}{2}}|\mathbf{L}_t, \mathbf{F}_t, \mathbf{A})p_C(\mathbf{F}_{t-1}|\mathbf{L}_{t-1}, \mathbf{F}_{t-\frac{1}{2}}, \mathbf{A}),$$

$$p_P(\mathbf{F}_{t-\frac{1}{2}}|\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}) = \mathcal{N}_w(\mathbf{F}_{t-\frac{1}{2}}|\mathbf{F}_t + (\sigma_t^2 - \sigma_{t-1}^2)\hat{\epsilon}_{\mathbf{F}}(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t), \frac{\sigma_{t-1}^2(\sigma_t^2 - \sigma_{t-1}^2)}{\sigma_t^2} \mathbf{I}),$$

$$p_C(\mathbf{F}_{t-1}|\mathbf{L}_{t-1}, \mathbf{F}_{t-\frac{1}{2}}, \mathbf{A}) = \mathcal{N}_w(\mathbf{F}_{t-1}|\mathbf{F}_t + \gamma \frac{\sigma_{t-1}}{\sigma_1} \hat{\epsilon}_{\mathbf{F}}(\mathbf{L}_{t-1}, \mathbf{F}_{t-\frac{1}{2}}, \mathbf{A}, t-1), 2\gamma \frac{\sigma_{t-1}}{\sigma_1} \mathbf{I}),$$

where  $p_P, p_C$  are the transitions of the predictor and corrector. According to lemma 3, both of the transitions are periodic translation equivariant. Therefore, the transition  $p(\mathbf{F}_{t-1}|\mathbf{L}_t, \mathbf{F}_t, \mathbf{A})$  is periodic translation equivariant. As the prior distribution  $\mathcal{U}(0, 1)$  is periodic translation invariant, we finally prove that the marginal distribution  $p(\mathbf{F}_0)$  is periodic translation invariant based on lemma 1.

### A.3 Proof of Proposition 3

We rewrite proposition 3 as follows.

**Proposition 3.** *The score  $\hat{\epsilon}_{\mathbf{L}}$  by Eq. (9) is  $O(3)$ -equivariant, and the score  $\hat{\epsilon}_{\mathbf{F}}$  from Eq. (10) is periodic translation invariant. Hence, the generated distribution by DiffCSP is periodic  $E(3)$  invariant.*

*Proof.* We first prove the orthogonal invariance of the inner product term  $\mathbf{L}^\top \mathbf{L}$ . For any orthogonal transformation  $\mathbf{Q} \in \mathbb{R}^{3 \times 3}$ ,  $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$ , we have

$$(\mathbf{Q}\mathbf{L})^\top (\mathbf{Q}\mathbf{L}) = \mathbf{L}^\top \mathbf{Q}^\top \mathbf{Q}\mathbf{L} = \mathbf{L}^\top \mathbf{I}\mathbf{L} = \mathbf{L}^\top \mathbf{L}.$$

For the Fourier Transformation, consider  $k$  is even, we have

$$\begin{aligned} & \psi_{\text{FT}}(w(\mathbf{f}_j + \mathbf{t}) - w(\mathbf{f}_i + \mathbf{t}))[c, k] \\ &= \sin\left(2\pi m(w(f_{j,c} + t_c) - w(f_{i,c} + t_c))\right) \\ &= \sin\left(2\pi m(f_{j,c} - f_{i,c}) - 2\pi m\left((f_{j,c} - f_{i,c}) - (w(f_{j,c} + t_c) - w(f_{i,c} + t_c))\right)\right) \\ &= \sin(2\pi m(f_{j,c} - f_{i,c})) \\ &= \psi_{\text{FT}}(\mathbf{f}_j - \mathbf{f}_i)[c, k]. \end{aligned}$$

Similar results can be acquired as  $k$  is odd. Therefore, we have  $\psi_{\text{FT}}(w(\mathbf{f}_j + \mathbf{t}) - w(\mathbf{f}_i + \mathbf{t})) = \psi_{\text{FT}}(\mathbf{f}_j - \mathbf{f}_i)$ ,  $\forall \mathbf{t} \in \mathbb{R}^3$ , i.e., the Fourier Transformation  $\psi_{\text{FT}}$  is periodic translation invariant. According to the above, the message passing layers defined in Eq. (6)- (8) is periodic  $E(3)$  invariant. Hence, we can directly prove that the coordinate denoising term  $\hat{\epsilon}_{\mathbf{F}}$  is periodic translation invariant. Let  $\hat{\epsilon}_i(\mathbf{L}, \mathbf{F}, \mathbf{A}, t) = \varphi_L(\frac{1}{N} \sum_{i=1}^N \mathbf{h}_i^{(S)})$ . For the lattice denoising term  $\hat{\epsilon}_{\mathbf{L}} = \mathbf{L}\hat{\epsilon}_i$ , we have

$$\begin{aligned} \hat{\epsilon}_{\mathbf{L}}(\mathbf{Q}\mathbf{L}, \mathbf{F}, \mathbf{A}, t) &= \mathbf{Q}\mathbf{L}\hat{\epsilon}_i(\mathbf{Q}\mathbf{L}, \mathbf{F}, \mathbf{A}, t) \\ &= \mathbf{Q}\mathbf{L}\hat{\epsilon}_i(\mathbf{L}, \mathbf{F}, \mathbf{A}, t) \\ &= \mathbf{Q}\hat{\epsilon}_{\mathbf{L}}(\mathbf{L}, \mathbf{F}, \mathbf{A}, t), \forall \mathbf{Q} \in \mathbb{R}^{3 \times 3}, \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}. \end{aligned}$$

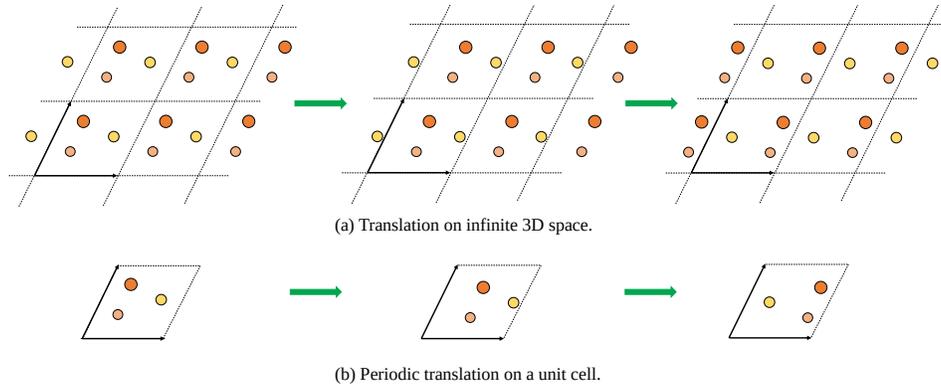


Figure 4: An example of periodic translation invariance. From the view of a unit cell, the atoms translated across the right boundary will be brought back to the left side.

Above all,  $\hat{\epsilon}_L$  is  $O(3)$ -equivariant, and  $\hat{\epsilon}_F$  is periodic translation invariant. According to proposition 1 and 2, the generated distribution by DiffCSP in Algorithm 2 is periodic  $E(3)$  invariant.  $\square$

#### A.4 Discussion on Periodic Translation Invariance

In Definition 3, we define the periodic translation invariance as a combination of translation invariance and periodicity. To see this, we illustrate an additional example in Figure 4. From a global view, when we translate all atom coordinates from left to right, the crystal structure remains unchanged, which indicates translation invariance. At the same time, from the view of a unit cell, the atom translated across the right boundary will be brought back to the left side owing to periodicity. Therefore, for convenience, we define the joint effect of translation invariance and periodicity as periodic translation invariance.

Previous works [36] have shown that shifting the periodic boundaries will not change the crystal structure. In this section, we further show that such periodic boundary shifting is equivalent to the periodic translation defined in Definition 3.

Consider two origin points  $\mathbf{p}_1, \mathbf{p}_2 \in \mathbb{R}^{3 \times 1}$  and the lattice matrix  $\mathbf{L}$ , the constructed unit cells by  $\mathbf{p}_1, \mathbf{p}_2$  can be represented as  $\mathcal{M}_1 = (\mathbf{A}_1, \mathbf{F}_1, \mathbf{L})$  and  $\mathcal{M}_2 = (\mathbf{A}_2, \mathbf{F}_2, \mathbf{L})$ , where  $\mathbf{F}_1, \mathbf{F}_2 \in \mathbb{R}^{3 \times N}$  are fractional coordinates and

$$\{(\mathbf{a}'_{1,i}, \mathbf{x}'_{1,i}) | \mathbf{a}'_{1,i} = \mathbf{a}_{1,i}, \mathbf{x}'_{1,i} = \mathbf{p}_1 + \mathbf{L}\mathbf{f}_{1,i} + \mathbf{L}\mathbf{k}, \forall \mathbf{k} \in \mathbb{Z}^{3 \times 1}\} \quad (15)$$

$$= \{(\mathbf{a}'_{1,j}, \mathbf{x}'_{1,j}) | \mathbf{a}'_{1,j} = \mathbf{a}_{1,j}, \mathbf{x}'_{1,j} = \mathbf{p}_2 + \mathbf{L}\mathbf{f}_{2,j} + \mathbf{L}\mathbf{k}, \forall \mathbf{k} \in \mathbb{Z}^{3 \times 1}\}, \quad (16)$$

which means that the unit cells formed by different origin points actually represent the same infinite crystal structures [36]. We further construct a bijection  $\mathcal{T} : \mathcal{M}_1 \rightarrow \mathcal{M}_2$  mapping each atom in  $\mathcal{M}_1$  to the corresponding atom in unit cell  $\mathcal{M}_2$ . For the pair  $(a_{1,i}, f_{1,i}) \in \mathcal{M}_1, (a_{2,j}, f_{2,j}) \in \mathcal{M}_2$ , we have  $\mathcal{T}(a_{1,i}, f_{1,i}) = (a_{2,j}, f_{2,j})$  iff  $\exists \mathbf{k}_i \in \mathbb{Z}^{3 \times 1}, s.t.$

$$\begin{cases} \mathbf{a}_{1,i} = \mathbf{a}_{2,j}, \\ \mathbf{p}_1 + \mathbf{L}\mathbf{f}_{1,i} + \mathbf{L}\mathbf{k}_i = \mathbf{p}_2 + \mathbf{L}\mathbf{f}_{2,j}. \end{cases}$$

After proper transformation, we have

$$\mathbf{f}_{2,j} = \mathbf{f}_{1,i} + \mathbf{L}^{-1}(\mathbf{p}_1 - \mathbf{p}_2) + \mathbf{k}_i. \quad (17)$$

As  $\mathbf{f}_{1,i}, \mathbf{f}_{2,i} \in [0, 1)^{3 \times 1}$ , we have

$$\begin{cases} \mathbf{k}_i = -[\mathbf{f}_{1,i} + \mathbf{L}^{-1}(\mathbf{p}_1 - \mathbf{p}_2)], \\ \mathbf{f}_{2,j} = w(\mathbf{f}_{1,i} + \mathbf{L}^{-1}(\mathbf{p}_1 - \mathbf{p}_2)), \end{cases}$$

which means shifting the periodic boundaries by changing the origin point  $\mathbf{p}_1$  into  $\mathbf{p}_2$  is equivalent to a periodic translation  $\mathbf{F}_2 = w(\mathbf{F}_1 + \mathbf{L}^{-1}(\mathbf{p}_1 - \mathbf{p}_2)\mathbf{1}^\top)$ .

## B Implementation Details

### B.1 Approximation of the Wrapped Normal Distribution

The Probability Density Function (PDF) of the wrapped normal distribution  $\mathcal{N}_w(0, \sigma_t^2)$  is

$$\mathcal{N}_w(x|0, \sigma_t^2) = \frac{1}{\sqrt{2\pi}\sigma_t} \sum_{k=-\infty}^{\infty} \exp\left(-\frac{(x-k)^2}{2\sigma_t^2}\right),$$

where  $x \in [0, 1)$ . Because the above series is convergent, it is reasonable to approximate the infinite summation to a finite truncated summation [61] as

$$f_{w,n}(x; 0, \sigma_t^2) = \frac{1}{\sqrt{2\pi}\sigma_t} \sum_{k=-n}^n \exp\left(-\frac{(x-k)^2}{2\sigma_t^2}\right).$$

And the logarithmic gradient of  $f$  can be formulated as

$$\begin{aligned} \nabla_x \log f_{w,n}(x; 0, \sigma_t^2) &= \nabla_x \log \left( \frac{1}{\sqrt{2\pi}\sigma_t} \sum_{k=-n}^n \exp\left(-\frac{(x-k)^2}{2\sigma_t^2}\right) \right) \\ &= \nabla_x \log \left( \sum_{k=-n}^n \exp\left(-\frac{(x-k)^2}{2\sigma_t^2}\right) \right) \\ &= \frac{\sum_{k=-n}^n (k-x) \exp\left(-\frac{(x-k)^2}{2\sigma_t^2}\right)}{\sigma_t^2 \sum_{k=-n}^n \exp\left(-\frac{(x-k)^2}{2\sigma_t^2}\right)} \end{aligned}$$

To estimate  $\lambda_t = \mathbb{E}_{x \sim \mathcal{N}_w(0, \sigma_t^2)}^{-1} [\|\nabla_x \log \mathcal{N}_w(x|0, \sigma_t^2)\|_2^2]$ , we first sample  $m$  points from  $\mathcal{N}_w(0, \sigma_t^2)$ , and the expectation is approximated as

$$\begin{aligned} \tilde{\lambda}_t &= \left[ \frac{1}{m} \sum_{i=1}^m \|\nabla_x \log f_{w,n}(x_i; 0, \sigma_t^2)\|_2^2 \right]^{-1} \\ &= \left[ \frac{1}{m} \sum_{i=1}^m \left\| \frac{\sum_{k=-n}^n (k-x_i) \exp\left(-\frac{(x_i-k)^2}{2\sigma_t^2}\right)}{\sigma_t^2 \sum_{k=-n}^n \exp\left(-\frac{(x_i-k)^2}{2\sigma_t^2}\right)} \right\|_2^2 \right]^{-1}. \end{aligned}$$

For implementation, we select  $n = 10$  and  $m = 10000$ .

### B.2 Adaptation of CDVAE

As illustrated in Figure 5, the original CDVAE [9] mainly consists of three parts: (1) a 3D encoder to encode the structure into the latent variable  $z_{3D}$ , (2) a property predictor to predict the lattice  $\mathbf{L}$ , the number of nodes in the unit cell  $N$ , and the proportion of each element in the composition  $c$ , (3) a 3D decoder to generate the structure from  $z_{3D}$ ,  $\mathbf{L}$ ,  $N$ ,  $c$  via the Score Matching with Langevin Dynamics (SMLD, Song and Ermon [49]) method. The training objective is composed of the loss functions on the three parts, *i.e.* the KL divergence between the encoded distribution and the standard normal distribution  $\mathcal{L}_{KL}$ , the aggregated prediction loss  $\mathcal{L}_{AGG}$  and the denoising loss on the decoder  $\mathcal{L}_{DEC}$ . Formally, we have

$$\mathcal{L}_{ORI} = \mathcal{L}_{AGG} + \mathcal{L}_{DEC} + \beta D_{KL}(\mathcal{N}(\mu_{3D}, \sigma_{3D}^2 \mathbf{I}) \| \mathcal{N}(0, \mathbf{I})).$$

We formulate  $\mathcal{L}_{KL} = \beta D_{KL}(\mathcal{N}(\mu_{3D}, \sigma_{3D}^2 \mathbf{I}) \| \mathcal{N}(0, \mathbf{I}))$  for better comparison with the adapted method.  $\beta$  is the hyper-parameter to balance the scale of the KL divergence and other loss functions.

To adapt the CDVAE framework to the CSP task, we apply two main changes. Firstly, for the encoder side, to take the composition as the condition, we apply an additional 1D prior encoder to encode the composition set into a latent distribution  $\mathcal{N}(\mu_{1D}, \sigma_{1D}^2 \mathbf{I})$  and minimize the KL divergence between the 3D and 1D distribution. The training objective is modified into

$$\mathcal{L}_{ADA} = \mathcal{L}_{AGG} + \mathcal{L}_{DEC} + \beta D_{KL}(\mathcal{N}(\mu_{3D}, \sigma_{3D}^2 \mathbf{I}) \| \mathcal{N}(\mu_{1D}, \sigma_{1D}^2 \mathbf{I})).$$

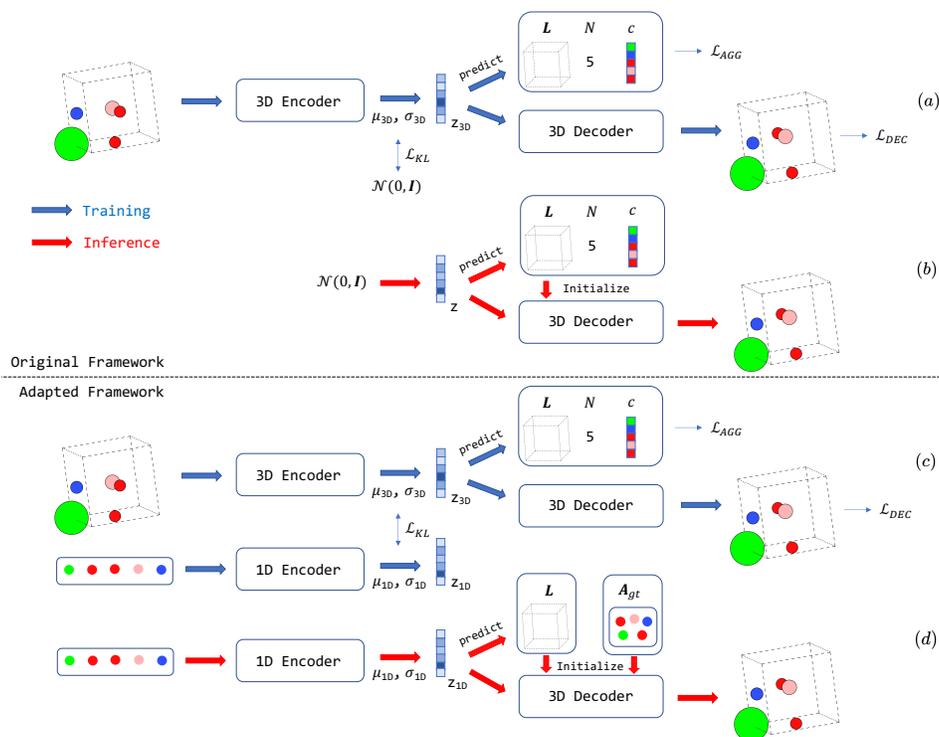


Figure 5: Overview of the original (a,b) and adapted (c,d) CDVAE. The key adaptations lie in two points. (1) We introduce an additional 1D prior encoder to fit the latent distribution of the given composition. (2) We initialize the generation procedure of the 3D decoder with the ground truth composition and keep the atom types unchanged to ensure the generated structure conforms to the given composition.

During the inference procedure, as the composition is given, the latent variable  $z_{1D}$  is sampled from  $\mathcal{N}(\mu_{1D}, \sigma_{1D}^2 \mathbf{I})$ . For implementation, we apply a Transformer [51] without positional encoding as the 1D encoder to ensure the permutation invariance. Secondly, for the generation procedure, we apply the ground truth composition for initialization and keep the atom types unchanged during the Langevin dynamics to ensure the generated structure conforms to the given composition.

### B.3 Algorithms for Training and Sampling

Algorithm 1 summarizes the forward diffusion process as well as the training of the denoising model  $\phi$ , while Algorithm 2 illustrates the backward sampling process. They can maintain the symmetries if  $\phi$  is delicately constructed. Notably, We apply the predictor-corrector sampler [50] to sample  $\mathbf{F}_0$ . In Algorithm 2, Line 7 refers to the predictor while Lines 9-10 correspond to the corrector.

---

#### Algorithm 1 Training Procedure of DiffCSP

---

- 1: **Input:** lattice matrix  $\mathbf{L}_0$ , atom types  $\mathbf{A}$ , fractional coordinates  $\mathbf{F}_0$ , denoising model  $\phi$ , and the number of sampling steps  $T$ .
  - 2: Sample  $\epsilon_L \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \epsilon_F \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $t \sim \mathcal{U}(1, T)$ .
  - 3:  $\mathbf{L}_t \leftarrow \sqrt{\alpha_t} \mathbf{L}_0 + \sqrt{1 - \alpha_t} \epsilon_L$
  - 4:  $\mathbf{F}_t \leftarrow w(\mathbf{F}_0 + \sigma_t \epsilon_F)$
  - 5:  $\hat{\epsilon}_L, \hat{\epsilon}_F \leftarrow \phi(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t)$
  - 6:  $\mathcal{L}_L \leftarrow \|\epsilon_L - \hat{\epsilon}_L\|_2^2$
  - 7:  $\mathcal{L}_F \leftarrow \lambda_t \|\nabla_{\mathbf{F}_t} \log q(\mathbf{F}_t | \mathbf{F}_0) - \hat{\epsilon}_F\|_2^2$
  - 8: Minimize  $\mathcal{L}_L + \mathcal{L}_F$
-

---

**Algorithm 2** Sampling Procedure of DiffCSP

---

1: **Input:** atom types  $\mathbf{A}$ , denoising model  $\phi$ , number of sampling steps  $T$ , step size of Langevin dynamics  $\gamma$ .  
2: Sample  $\mathbf{L}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \mathbf{F}_T \sim \mathcal{U}(0, 1)$ .  
3: **for**  $t \leftarrow T, \dots, 1$  **do**  
4: Sample  $\epsilon_L, \epsilon_F, \epsilon'_F \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$   
5:  $\hat{\epsilon}_L, \hat{\epsilon}_F \leftarrow \phi(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}, t)$ .  
6:  $\mathbf{L}_{t-1} \leftarrow \frac{1}{\sqrt{\alpha_t}}(\mathbf{L}_t - \frac{\beta_t}{\sqrt{1-\alpha_t}}\hat{\epsilon}_L) + \sqrt{\beta_t \cdot \frac{1-\alpha_{t-1}}{1-\alpha_t}}\epsilon_L$ .  
7:  $\mathbf{F}_{t-\frac{1}{2}} \leftarrow w(\mathbf{F}_t + (\sigma_t^2 - \sigma_{t-1}^2)\hat{\epsilon}_F + \frac{\sigma_{t-1}\sqrt{\sigma_t^2 - \sigma_{t-1}^2}}{\sigma_t}\epsilon_F)$   
8:  $\_, \hat{\epsilon}_F \leftarrow \phi(\mathbf{L}_{t-1}, \mathbf{F}_{t-\frac{1}{2}}, \mathbf{A}, t-1)$ .  
9:  $d_t \leftarrow \gamma\sigma_{t-1}/\sigma_1$   
10:  $\mathbf{F}_{t-1} \leftarrow w(\mathbf{F}_{t-\frac{1}{2}} + d_t\hat{\epsilon}_F + \sqrt{2d_t}\epsilon'_F)$ .  
11: **end for**  
12: **Return**  $\mathbf{L}_0, \mathbf{F}_0$ .

---

#### B.4 Hyper-parameters and Training Details

We acquire the origin datasets from CDVAE [9]<sup>4</sup> and MPTS-52 [57]<sup>5</sup>. We utilize the codebases from GN-OA [21]<sup>6</sup>, cG-SchNet [53]<sup>7</sup> and CDVAE [9]<sup>8</sup> for baseline implementations.

For the optimization methods, we apply the MEGNet [52] with 3 layers, 32 hidden states as property predictor. The model is trained for 1000 epochs with an Adam optimizer with learning rate  $1 \times 10^{-3}$ . As for the optimization algorithms, we apply RS, PSO, and BO according to Cheng et al. [21]. For RS and BO, We employ random search and TPE-based BO as implemented in Hyperopt [62]<sup>9</sup>. Specifically, we choose observation quantile  $\gamma$  as 0.25 and the number of initial random points as 200 for BO. For PSO, we used scikit-opt<sup>10</sup> and choose the momentum parameter  $\omega$  as 0.8, the cognitive as 0.5, the social parameters as 0.5 and the size of population as 20.

For P-cG-SchNet, we apply the SchNet [29] with 9 layers, 128 hidden states as the backbone model. The model is trained for 500 epochs on each dataset with an Adam optimizer with initial learning rate  $1 \times 10^{-4}$  and a Plateau scheduler with a decaying factor 0.5 and a patience of 10 epochs. We select the element proportion and the number of atoms in a unit cell as conditions for the CSP task. For CDVAE, we apply the DimeNet++ [63] with 4 layers, 256 hidden states as the encoder and the GemNet-T [64] with 3 layers, 128 hidden states as the decoder. We further apply a Transformer [51] model with 2 layers, 128 hidden states as the additional prior encoder as proposed in Appendix B.2. The model is trained for 3500, 4000, 1000, 1000 epochs for Perov-5, Carbon-24, MP-20 and MPTS-52 respectively with an Adam optimizer with initial learning rate  $1 \times 10^{-3}$  and a Plateau scheduler with a decaying factor 0.6 and a patience of 30 epochs. For our DiffCSP, we utilize the setting of 4 layer, 256 hidden states for Perov-5 and 6 layer, 512 hidden states for other datasets. The dimension of the Fourier embedding is set to  $k = 256$ . We apply the cosine scheduler with  $s = 0.008$  to control the variance of the DDPM process on  $\mathbf{L}_t$ , and an exponential scheduler with  $\sigma_1 = 0.005, \sigma_T = 0.5$  to control the noise scale of the score matching process on  $\mathbf{F}_t$ . The diffusion step is set to  $T = 1000$ . Our model is trained for 3500, 4000, 1000, 1000 epochs for Perov-5, Carbon-24, MP-20 and MPTS-52 with the same optimizer and learning rate scheduler as CDVAE. For the step size  $\gamma$  in Langevin dynamics for the structure prediction task, we apply  $\gamma = 5 \times 10^{-7}$  for Perov-5,  $1 \times 10^{-5}$  for MP-20 and MPTS-52, and for Carbon-24, we apply  $\gamma = 5 \times 10^{-6}$  to predict one sample and  $\gamma = 5 \times 10^{-7}$  for multiple samples. For the ab initio generation and optimization task on Perov-5, Carbon-24 and MP-20, we apply  $\gamma = 1 \times 10^{-6}, 1 \times 10^{-5}, 5 \times 10^{-6}$ , respectively. All models are trained on GeForce RTX 3090 GPU.

---

<sup>4</sup><https://github.com/txie-93/cdvae/tree/main/data>

<sup>5</sup><https://github.com/sparks-baird/mp-time-split>

<sup>6</sup>[http://www.comates.group/links?software=gn\\_oa](http://www.comates.group/links?software=gn_oa)

<sup>7</sup><https://github.com/atomistic-machine-learning/cG-SchNet>

<sup>8</sup><https://github.com/txie-93/cdvae>

<sup>9</sup><https://github.com/hyperopt/hyperopt>

<sup>10</sup><https://github.com/guofei9987/scikit-opt>

## C Exploring the Effects of Sampling and Candidate Ranking

### C.1 Impact of Sampling Numbers

Figure 6 illustrates the impact of sampling numbers on the match rate. The match rate of all methods increases when sampling more candidates, and DiffCSP outperforms the baselines methods under the arbitrary number of samples.

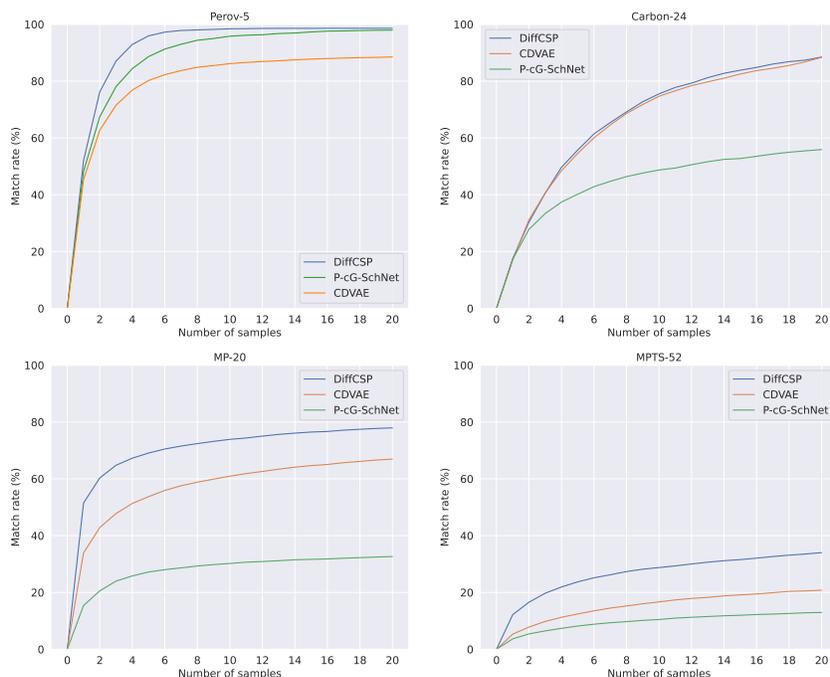


Figure 6: Comparison on different number of samples.

### C.2 Diversity

We further evaluate the diversity by yielding the CrystalNN [65] fingerprint of each generated structure, calculating the L2-distances among all pairs in the 20 samples of each composition, collating the mean and max value of the distances, and finally averaging the results from all testing candidates. We list the diversity of CDVAE and DiffCSP on Perov-5 and MP-20 in Table 5.

Table 5: Comparison on the diversity.

	Perov-5		MP-20	
	Mean	Max	Mean	Max
CDVAE	0.3249	0.7316	<b>0.3129</b>	<b>0.6979</b>
DiffCSP	<b>0.3860</b>	<b>0.8911</b>	0.2030	0.5292

### C.3 Ranking among Multiple Candidates

In § 5.1, we match each candidate with the ground truth structure to pick up the best sample. However, in real CSP scenarios, the ground truth structure is not available, necessitating a confidence model to rank the generated candidates. To address this, we develop three types of confidence models to score each sample for ranking from different perspectives: **Energy Predictor (EP)**. Since lower formation energy typically leads to more stable structures, we directly train a predictor using energy labels and apply the negative of the predicted

Table 6: Results on different confidence models. *Oracle* means applying the negative RMSD against the ground truth as the ranking score.

	Match rate (%) $\uparrow$	RMSE $\downarrow$
EP	51.96	0.0589
MD (d=0.1)	60.30	0.0357
MD (d=0.3)	60.13	0.0382
MD (d=0.5)	59.20	0.0469
CS	58.81	0.0443
Oracle	77.93	0.0492

energy as the confidence score. **Match Discriminator (MD)**. Inspired by Diffdock [12], we first generate five samples for each composition in the training/validation set using DiffCSP and calculate their RMSDs with the ground truth. We then train a binary classifier to predict whether the sample matches the ground truth with an RMSD below a threshold  $d$ . The predicted probability serves as the score. **Contrastive Scorer (CS)**. Drawing inspiration from CLIP [66], we train a contrastive model between a 1D and 3D model to align the corresponding compositions and structures. The inner product of the 1D and 3D models is used as the score. We select the Top-1 result among the 20 candidates ranked by each confidence model on the MP-20 dataset, as shown in Table 6. The results indicate that MD and CS perform relatively better, but there remains a gap between the heuristic ranking models and the oracle ranker. Designing powerful ranking models is an essential problem, which we leave for future research.

## D Impact of Noise Schedulers

We explore the noise schedulers from three perspectives and list the results in Table 7 and 8. **1.** For lattices, we originally use the cosine scheduler with  $\beta = 0.008$ , and we change it into the linear and sigmoid schedulers with  $\beta_1 = 0.0001$ ,  $\beta_T = 0.02$ . We find that the linear scheduler yields comparable results, while the sigmoid scheduler hinders the performance. **2.** For fractional coordinates, we use the exponential scheduler with  $\sigma_{min} = 0.005$ ,  $\sigma_{max} = 0.5$ , and we change the value of  $\sigma_{max}$  into 0.1 and 1.0. Results show that the small- $\sigma_{max}$  variant performs obviously worse, as only sufficiently large  $\sigma_{max}$  could approximate the prior uniform distribution. We visualize the PDF curves in Figure 7 for better understanding. **3.** For atom types, we conduct similar experiments as lattices, and the results indicate that the original cosine scheduler performs better. In conclusion, we suggest applying the proposed noise schedulers.

Table 7: CSP results on different schedulers. DiffCSP utilizes cosine scheduler on  $L$  and  $\sigma_{max} = 0.5$  on  $F$ .

	Match rate (%) $\uparrow$	RMSE $\downarrow$
DiffCSP	51.49	0.0631
<i>Schedulers of <math>L</math></i>		
Linear	50.06	0.0590
Sigmoid	45.24	0.0664
<i>Schedulers of <math>F</math></i>		
$\sigma_{max} = 0.1$	32.56	0.0913
$\sigma_{max} = 1.0$	47.89	0.0675

Table 8: Ab initio generation results on different type schedulers. DiffCSP utilizes cosine scheduler on  $A$ .

	Validity		Coverage	
	Struct.(%)	Comp.(%)	Recall	Precision
DiffCSP	100.00	83.25	99.71	99.76
Linear	99.70	79.78	98.29	99.48
Sigmoid	99.88	81.59	99.33	99.55

## E Learning Curves of Different Variants

We plot the curves of training and validation loss of different variants proposed in § 5.3 in Figure 8 with the following observations. **1.** The multi-graph methods struggle with higher training and validation loss, as the edges constructed under different disturbed lattices vary significantly, complicating the training procedure. **2.** The Fourier transformation, expanding the relative coordinates and maintaining the periodic translation invariance, helps the model converge faster at the beginning of the training procedure. **3.** The variant utilizing the fully connected graph without the Fourier transformation (named “DiffCSP w/o FT” in Figure 8) encounters obvious overfitting as the periodic translation invariance is violated, highlighting the necessity to leverage the desired invariance into the model.

## F Comparison with DFT-based Methods

### F.1 Implementation Details

We utilize USPEX [59], a DFT-based software equipped with the evolutionary algorithm to search for stable structures. We use 20 populations in each generation, and end up with 20 generations for each compound. We set 60% of the lowest-enthalpy structures allowed to produce the next generation through heredity (50%), lattice mutation (10%), and atomic permutation (20%). Moreover, two

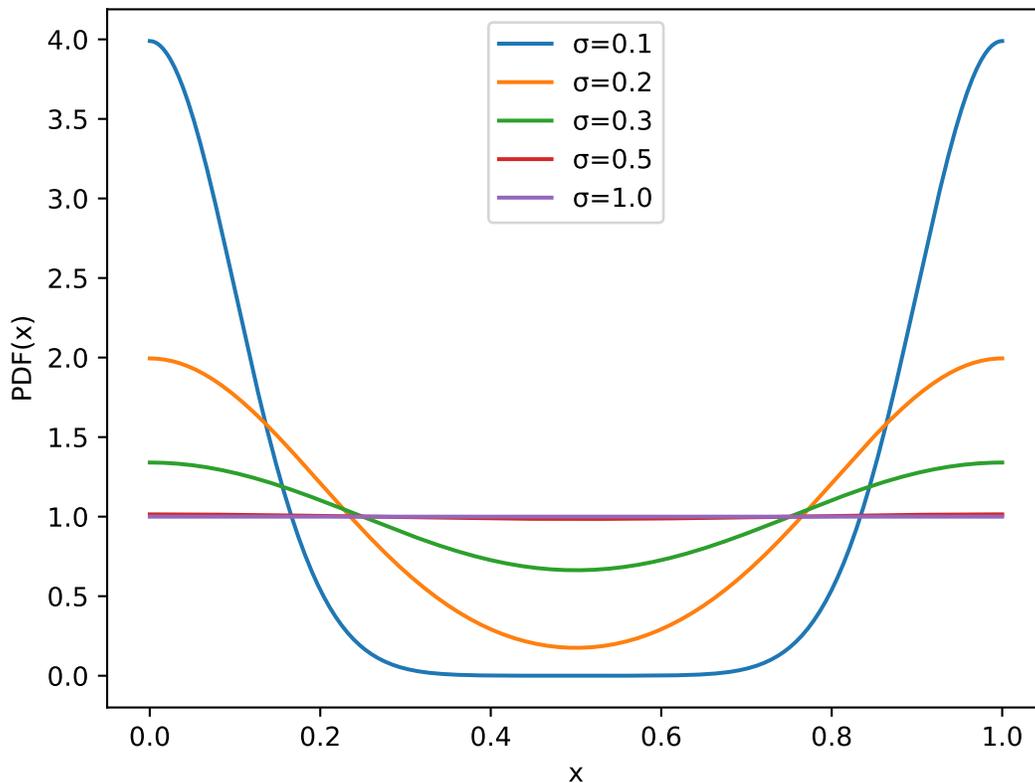


Figure 7: PDF curves of the wrapped normal distribution with periodic as 1 and different noise scales. It can be find that  $\sigma = 0.5$  is practically large enough to approximate the prior uniform distribution.

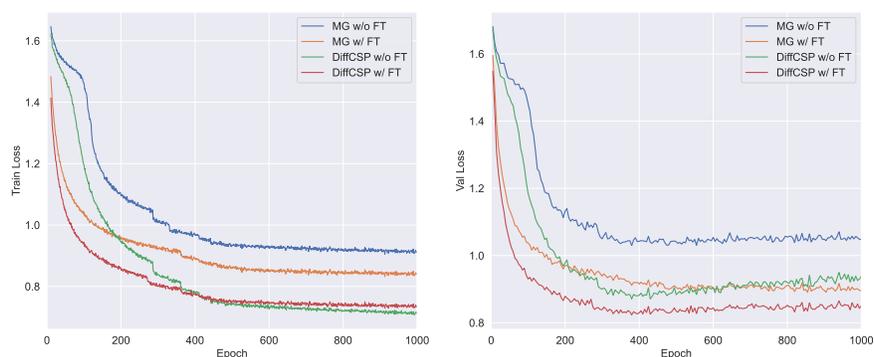


Figure 8: Learning curves of different variants proposed in § 5.3. MG and FT denote multi-graph edge construction and Fourier transformation, respectively.

lowest-enthalpy structures are allowed to survive into the next generation. The structural relaxations are calculated by the frozen-core all-electron projector augmented wave (PAW) method [67] as implemented in the Vienna ab initio simulation package (VASP) [68]. Each sample's calculation is performed using one node with 48 cores (Intel(R) Xeon(R) CPU E5-2692 v2 @ 2.20GHz), while the populations within the same generation are concurrently computed across 20 nodes in parallel. The exchange-correlation energy is treated within the generalized gradient approximation (GGA), using the Perdew-Burke-Ernzerhof (PBE) function [69].

## F.2 Results

We select 10 binary and 5 ternary compounds in MP-20 testing set and compare our model with USPEX. For our method, we sample 20 candidates for each compound following the setting in Table 1. For USPEX, we apply 20 generations, 20 populations for each compound, and select the best sample in each generation, leading to 20 candidates as well. We list the minimum RMSD of each compound in Table 9, and additionally summarize the match rate over the 15 compounds, the averaged RMSD over the matched structures, and the averaged inference time to generate 20 candidates for each compound in Table 10. The results show that DiffCSP correctly predicts more structures with higher match rate and significantly lower time cost.

Table 9: The minimum RMSD of 20 candidates of USPEX and DiffCSP, "N/A" means none of the candidates match with the ground truth.

Binary	Co <sub>2</sub> Sb <sub>2</sub>	Sr <sub>2</sub> O <sub>4</sub>	AlAg <sub>4</sub>	YMg <sub>3</sub>	Cr <sub>4</sub> Si <sub>4</sub>
USPEX	0.0008	0.0121	N/A	0.0000	N/A
DiffCSP	0.0005	0.0133	0.0229	0.0003	0.0066
Binary	Sn <sub>4</sub> Pd <sub>4</sub>	Ag <sub>6</sub> O <sub>2</sub>	Co <sub>4</sub> B <sub>2</sub>	Ba <sub>2</sub> Cd <sub>6</sub>	Bi <sub>2</sub> F <sub>8</sub>
USPEX	0.0184	0.0079	0.0052	N/A	N/A
DiffCSP	0.0264	N/A	N/A	0.0028	N/A
Ternary	KZnF <sub>3</sub>	Cr <sub>3</sub> CuO <sub>8</sub>	Bi <sub>4</sub> S <sub>4</sub> Cl <sub>4</sub>	Si <sub>2</sub> (CN <sub>2</sub> ) <sub>4</sub>	Hg <sub>2</sub> S <sub>2</sub> O <sub>8</sub>
USPEX	0.0123	N/A	N/A	N/A	0.0705
DiffCSP	0.0006	0.0482	0.0203	N/A	0.0473

Table 10: Overall results over the 15 selected compounds.

	Match Rate (%) $\uparrow$	Avg. RMSD $\downarrow$	Avg. Time $\downarrow$
USPEX	53.33	<b>0.0159</b>	12.5h
DiffCSP	<b>73.33</b>	0.0172	<b>10s</b>

## G Extension to More General Tasks

### G.1 Overview

Our method mainly focuses on CSP, aiming at predicting the stable structures from the fixed composition  $\mathbf{A}$ . We first enable the generation on atom types and extend to ab initio generation task in § G.2, and then adopt the energy guidance for property optimization in § G.3. Figure 9 illustrated the differences and connections of CSP, ab initio generation, and property optimization. Besides, CDVAE [9] proposes a reconstruction task specific for VAE-based models, which first encodes the ground truth structure, and require the model to recover the input structure. As our method follows a diffusion-based framework instead of VAE, it is unsuitable to conduct the reconstruction task of our method.

### G.2 Ab Initio Generation

Apart from  $\mathbf{L}$  and  $\mathbf{F}$ , the ab initio generation task additionally requires the generation on  $\mathbf{A}$ . As the atom types can be viewed as either  $N$  discrete variables in  $h$  classes, or the one-hot representation  $\mathbf{A} \in \mathbb{R}^{h \times N}$  in continuous space. We apply two lines of diffusion processes for the type generation, which are detailedly shown as follows.

**Multinomial Diffusion for Discrete Generation** Regarding the atom types as discrete features, we apply the multinomial diffusion with the following forward diffusion process [70, 42, 71],

$$q(\mathbf{A}_t | \mathbf{A}_0) = \mathcal{C}(\mathbf{A}_t | \bar{\alpha}_t \mathbf{A}_0 + \frac{(1 - \bar{\alpha}_t)}{h} \mathbf{1}_h \mathbf{1}_N^\top), \quad (18)$$

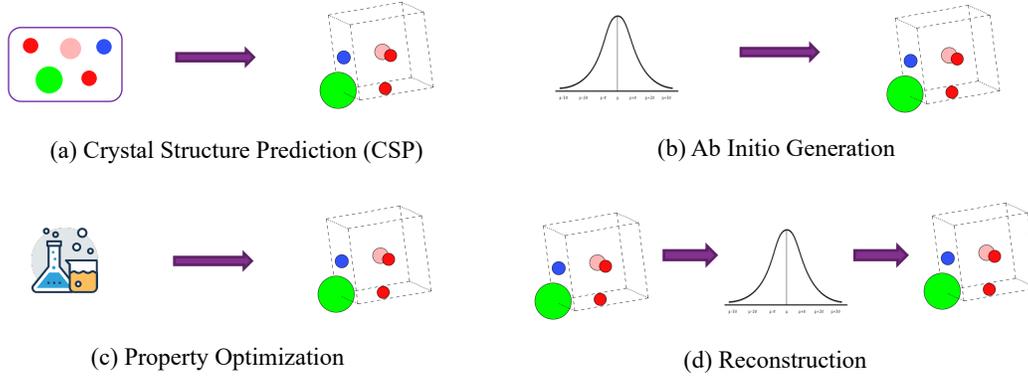


Figure 9: Different tasks for crystal generation. Our approach mainly focuses on the CSP task (a), and is capable to extend into the ab initio generation task (b) by further generating the composition, and the property optimization task (c) via introducing the guidance on the target property. The reconstruction task (d) in Xie et al. [9] is specific for VAE, which is unnecessary for our diffusion-based method.

where  $\mathbf{1}_h \in \mathbb{R}^{h \times 1}$ ,  $\mathbf{1}_N \in \mathbb{R}^{N \times 1}$  are vectors with all elements setting to one,  $\mathbf{A}_0$  is the one-hot representation of the origin composition, and the function  $\mathcal{C}(\cdot)$  samples the multinomial distribution with the conditional probability and returns the one-hot representation of  $\mathbf{A}_t$ .

The corresponding backward generation process is defined as

$$p(\mathbf{A}_{t-1}|\mathcal{M}_t) = \mathcal{C}(\mathbf{A}_{t-1}|\tilde{\theta}_t / \sum_{k=1}^h \tilde{\theta}_{t,k}), \quad (19)$$

where

$$\tilde{\theta}_t = \left( \alpha_t \mathbf{A}_t + \frac{(1 - \alpha_t)}{h} \mathbf{1}_h \mathbf{1}_N^\top \right) \odot \left( \bar{\alpha}_t \hat{\epsilon}_A(\mathcal{M}_t, t) + \frac{(1 - \bar{\alpha}_t)}{h} \mathbf{1}_h \mathbf{1}_N^\top \right), \quad (20)$$

and  $\hat{\epsilon}_A \in \mathbb{R}^{h \times N}$  is predicted by the denoising model. We further find that specific for  $t = 1$ ,  $\mathbf{A}_0 = \text{argmax}(\hat{\epsilon}_A(\mathcal{M}_1, 1))$  works better.

The training objective for multinomial diffusion is

$$\mathcal{L}_{\mathbf{A}, \text{discrete}} = \mathbb{E}_{\mathbf{A}_t \sim q(\mathbf{A}_t|\mathbf{A}_0), t \sim \mathcal{U}(1, T)} \left[ \text{KL}(q(\mathbf{A}_{t-1}|\mathbf{A}_t) \| p(\mathbf{A}_{t-1}|\mathbf{A}_t)) \right]. \quad (21)$$

**One-hot Diffusion for Continuous Generation** Another approach is to simply consider the composition  $\mathbf{A}$  as a continuous variable in real space  $\mathbb{R}^{h \times N}$ , which enables the application of standard DDPM-based method [41]. Similar to Eq. (2)-(4), the forward diffusion process is defined as

$$q(\mathbf{A}_t|\mathbf{A}_0) = \mathcal{N}\left(\mathbf{L}_t | \sqrt{\bar{\alpha}_t} \mathbf{A}_0, (1 - \bar{\alpha}_t) \mathbf{I}\right). \quad (22)$$

And the backward generation process is defined as

$$p(\mathbf{A}_{t-1}|\mathcal{M}_t) = \mathcal{N}(\mathbf{A}_{t-1} | \mu_{\mathbf{A}}(\mathcal{M}_t), \sigma_{\mathbf{A}}^2(\mathcal{M}_t) \mathbf{I}), \quad (23)$$

where  $\mu_{\mathbf{A}}(\mathcal{M}_t) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{A}_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \hat{\epsilon}_A(\mathcal{M}_t, t) \right)$ ,  $\sigma_{\mathbf{A}}^2(\mathcal{M}_t) = \beta_t \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}$ . The denoising term  $\hat{\epsilon}_A(\mathcal{M}_t, t) \in \mathbb{R}^{h \times N}$  is predicted by the model.

The training objective for one-hot diffusion is

$$\mathcal{L}_{\mathbf{A}, \text{continuous}} = \mathbb{E}_{\epsilon_{\mathbf{A}} \sim \mathcal{N}(0, \mathbf{I}), t \sim \mathcal{U}(1, T)} \left[ \|\epsilon_{\mathbf{A}} - \hat{\epsilon}_A(\mathcal{M}_t, t)\|_2^2 \right]. \quad (24)$$

The entire objective for training the joint diffusion model on  $\mathbf{L}$ ,  $\mathbf{F}$ ,  $\mathbf{A}$  is combined as

$$\mathcal{L}_{\mathcal{M}} = \lambda_{\mathbf{L}} \mathcal{L}_{\mathbf{L}} + \lambda_{\mathbf{F}} \mathcal{L}_{\mathbf{F}} + \lambda_{\mathbf{A}} \mathcal{L}_{\mathbf{A}}. \quad (25)$$

We select  $\lambda_{\mathbf{L}} = \lambda_{\mathbf{F}} = 1, \lambda_{\mathbf{A}} = 0$  for the CSP task as  $\mathbf{A}$  is fixed during the generation process, and  $\lambda_{\mathbf{L}} = \lambda_{\mathbf{F}} = 1, \lambda_{\mathbf{A}} = 20$  for the ab initio generation task to balance the scale of each loss component. Specifically, we do not optimize  $\mathcal{L}_{\mathbf{A}}$  on the Carbon-24 dataset, as all atoms in this dataset are carbon.

**Sample Structures with Arbitrary Numbers of Atoms** As the number of atoms in a unit cell (*i.e.*  $N$ ) is unchanged during the generation process, we first sample  $N$  according to the distribution of  $N$  in the training set, which is similar to Hooeboom et al. [41]. The sampled distribution  $p(\mathcal{M})$  can be more concisely described as  $p(\mathcal{M}, N) = p(N)p(\mathcal{M}|N)$ . The former term  $p(N)$  is sampled from pre-computed data distribution, and the latter conditional distribution  $p(\mathcal{M}|N)$  is modeled by DiffCSP.

**Evaluation Metrics** The results are evaluated from three perspectives. **Validity:** We consider both the structural validity and the compositional validity. The structural valid rate is calculated as the percentage of the generated structures with all pairwise distances larger than  $0.5\text{\AA}$ , and the generated composition is valid if the entire charge is neutral as determined by SMACT [72]. **Coverage:** It measures the structural and compositional similarity between the testing set  $\mathcal{S}_t$  and the generated samples  $\mathcal{S}_g$ . Specifically, letting  $d_S(\mathcal{M}_1, \mathcal{M}_2), d_C(\mathcal{M}_1, \mathcal{M}_2)$  denote the  $L_2$  distances of the CrystalNN structural fingerprints [65] and the normalized Magpie compositional fingerprints [73], the COverage Recall (COV-R) is determined as  $\text{COV-R} = \frac{1}{|\mathcal{S}_t|} |\{\mathcal{M}_i | \mathcal{M}_i \in \mathcal{S}_t, \exists \mathcal{M}_j \in \mathcal{S}_g, d_S(\mathcal{M}_i, \mathcal{M}_j) < \delta_S, d_C(\mathcal{M}_i, \mathcal{M}_j) < \delta_C\}|$  where  $\delta_S, \delta_C$  are pre-defined thresholds. The COverage Precision (COV-P) is defined similarly by swapping  $\mathcal{S}_g, \mathcal{S}_t$ . **Property statistics:** We calculate three kinds of Wasserstein distances between the generated and testing structures, in terms of density, formation energy, and the number of elements [9], denoted as  $d_\rho, d_E$  and  $d_{\text{elem}}$ , individually. The validity and coverage metrics are calculated on 10,000 generated samples, and the property metrics are evaluated on a subset with 1,000 samples passing the validity test.

**Results** We denote the abovementioned discrete and continuous generation methods as DiffCSP-D, DiffCSP-C, respectively. The results of the two variants and the strongest baseline CDVAE on MP-20 are provided in Table 11. We observe that DiffCSP-D yields slightly lower validity and coverage rates than DiffCSP-C. Moreover, DiffCSP-D tends to generate structures with more types of elements, which is far from the data distribution. Hence, we select DiffCSP-C for the other experiments in Table 4 (abbreviated as DiffCSP). We further find that DiffCSP-C supports the property optimization task in the next section. Besides, both variants have lower composition validity than CDVAE, implying that more powerful methods for composition generation are required. As our paper mainly focuses on the CSP task, we leave this for future studies.

### G.3 Property Optimization

On top of DiffCSP-C, we further equip our method with energy guidance [74, 75] for property optimization. Specifically we train a time-dependent property predictor  $E(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}_t, t)$  with the same message passing blocks as Eq. (6)-(8). And the final prediction is acquired by the final layer as

$$E = \varphi_E \left( \frac{1}{N} \sum_{i=1}^N \mathbf{h}_i^{(S)} \right). \quad (26)$$

And the gradients *w.r.t.*  $\mathbf{L}, \mathbf{F}, \mathbf{A}$  additionally guide the generation process. As the inner product term  $\mathbf{L}^\top \mathbf{L}$  is  $O(3)$  invariant, and the Fourier transformation term  $\phi_{\text{FT}}(\mathbf{f}_j - \mathbf{f}_i)$  is periodic translation invariant, the predicted energy  $E$  is periodic  $E(3)$  invariant. That is,

$$E(\mathbf{Q}\mathbf{L}_t, w(\mathbf{F}_t + \mathbf{t}\mathbf{1}^\top, \mathbf{A}_t, t) = E(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}_t, t). \quad (27)$$

Table 11: Results on MP-20 ab initio generation task.

Data	Method	Validity (%) $\uparrow$		Coverage (%) $\uparrow$		Property $\downarrow$		
		Struc.	Comp.	COV-R	COV-P	$d_\rho$	$d_E$	$d_{\text{elem}}$
MP-20	CDVAE [9]	<b>100.0</b>	<b>86.70</b>	99.15	99.49	0.6875	0.2778	1.432
	DiffCSP-D	99.70	83.11	99.68	99.53	<b>0.1730</b>	0.1366	0.9959
	DiffCSP-C	<b>100.0</b>	83.25	<b>99.71</b>	<b>99.76</b>	0.3502	<b>0.1247</b>	<b>0.3398</b>

**Algorithm 3** Energy-Guided Sampling Procedure of DiffCSP(-C)

---

1: **Input:** denoising model  $\phi$ , energy predictor  $E$ , step size of Langevin dynamics  $\gamma$ , guidance magnitude  $s$ , input structure before optimization  $\mathcal{M} = (\mathbf{L}, \mathbf{F}, \mathbf{A})$ , number of sampling steps  $T'$ , maximum number of sampling steps  $T$ .

2: **if**  $T'=T$  **then**

3:   Sample  $\mathbf{L}_{T'} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \mathbf{F}_{T'} \sim \mathcal{U}(0, 1), \mathbf{A}_{T'} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .

4: **else**

5:   Sample  $\mathbf{L}_{T'} \sim q(\mathbf{L}_{T'}|\mathbf{L}), \mathbf{F}_{T'} \sim q(\mathbf{F}_{T'}|\mathbf{F}), \mathbf{A}_{T'} \sim q(\mathbf{A}_{T'}|\mathbf{A})$  according to Eq. (2),(5) and (22).

6: **end if**

7: **for**  $t \leftarrow T', \dots, 1$  **do**

8:   Sample  $\epsilon_L, \epsilon_F, \epsilon_A, \epsilon'_F \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

9:    $\hat{\epsilon}_L, \hat{\epsilon}_F, \hat{\epsilon}_A \leftarrow \phi(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}_t, t)$ .

10:   Acquire  $\nabla_{\mathbf{L}} E, \nabla_{\mathbf{F}} E, \nabla_{\mathbf{A}} E$  from  $E(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}_t, t)$ .

11:    $\mathbf{L}_{t-1} \leftarrow \frac{1}{\sqrt{\alpha_t}}(\mathbf{L}_t - \frac{\beta_t}{\sqrt{1-\alpha_t}}\hat{\epsilon}_L) - s\beta_t \cdot \frac{1-\alpha_{t-1}}{1-\alpha_t}\nabla_{\mathbf{L}} E + \sqrt{\beta_t \cdot \frac{1-\alpha_{t-1}}{1-\alpha_t}}\epsilon_L$ .

12:    $\mathbf{A}_{t-1} \leftarrow \frac{1}{\sqrt{\alpha_t}}(\mathbf{A}_t - \frac{\beta_t}{\sqrt{1-\alpha_t}}\hat{\epsilon}_A) - s\beta_t \cdot \frac{1-\alpha_{t-1}}{1-\alpha_t}\nabla_{\mathbf{A}} E + \sqrt{\beta_t \cdot \frac{1-\alpha_{t-1}}{1-\alpha_t}}\epsilon_A$ .

13:    $\mathbf{F}_{t-\frac{1}{2}} \leftarrow w(\mathbf{F}_t + (\sigma_t^2 - \sigma_{t-1}^2)\hat{\epsilon}_F - s\frac{\sigma_{t-1}^2(\sigma_t^2 - \sigma_{t-1}^2)}{\sigma_t^2}\nabla_{\mathbf{F}} E + \frac{\sigma_{t-1}\sqrt{\sigma_t^2 - \sigma_{t-1}^2}}{\sigma_t}\epsilon_F)$

14:    $-, \hat{\epsilon}_{F,-} \leftarrow \phi(\mathbf{L}_{t-1}, \mathbf{F}_{t-\frac{1}{2}}, \mathbf{A}_{t-1}, t-1)$ .

15:    $d_t \leftarrow \gamma\sigma_{t-1}/\sigma_1$

16:    $\mathbf{F}_{t-1} \leftarrow w(\mathbf{F}_{t-\frac{1}{2}} + d_t\hat{\epsilon}_{F,-} + \sqrt{2d_t}\epsilon'_F)$ .

17: **end for**

18: **Return**  $\mathbf{L}_0, \mathbf{F}_0, \text{argmax}(\mathbf{A}_0)$ .

---

Table 12: Results on property optimization task. The results of baselines are from Xie et al. [9].

Method	Perov-5			Carbon-24			MP-20		
	SR5	SR10	SR15	SR5	SR10	SR15	SR5	SR10	SR15
FTCP	0.06	0.11	0.16	0.00	0.00	0.00	0.02	0.04	0.05
Cond-DFC-VAE	0.55	0.64	0.69	-	-	-	-	-	-
CDVAE	0.52	0.65	0.79	0.00	0.06	0.06	0.78	0.86	0.90
DiffCSP	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>0.50</b>	<b>0.69</b>	<b>0.69</b>	<b>0.82</b>	<b>0.98</b>	<b>1.00</b>

Taking gradient to both sides *w.r.t.*  $\mathbf{L}, \mathbf{F}, \mathbf{A}$ , respectively, we have

$$\mathbf{Q}^\top \nabla_{\mathbf{L}'_t} E(\mathbf{L}'_t, w(\mathbf{F}_t + t\mathbf{1}^\top, \mathbf{A}_t, t)|_{\mathbf{L}'_t = \mathbf{Q}\mathbf{L}_t}) = \nabla_{\mathbf{L}_t} E(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}_t, t), \quad (28)$$

$$\nabla_{\mathbf{F}'_t} E(\mathbf{Q}\mathbf{L}_t, \mathbf{F}'_t, \mathbf{A}_t, t)|_{\mathbf{F}'_t = w(\mathbf{F}_t + t\mathbf{1}^\top)} = \nabla_{\mathbf{F}_t} E(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}_t, t), \quad (29)$$

$$\nabla_{\mathbf{A}_t} E(\mathbf{Q}\mathbf{L}_t, w(\mathbf{F}_t + t\mathbf{1}^\top, \mathbf{A}_t, t)) = \nabla_{\mathbf{A}_t} E(\mathbf{L}_t, \mathbf{F}_t, \mathbf{A}_t, t), \quad (30)$$

which implies that the gradient to  $\mathbf{L}_t$  is O(3) equivariant, and the gradient to  $\mathbf{F}_t$  and  $\mathbf{A}_t$  is periodic E(3) invariant. Such symmetries maintain that the introduction of energy guidance does not violate the periodic E(3) invariance of the marginal distribution.

The detailed algorithm for energy-guided sampling is provided in Algorithm 3. We find that  $s = 50$  works well on the three datasets. We evaluate the performance of the energy-guided model with the same metrics as Xie et al. [9]. We sample 100 structures from the testing set for optimization. For each structure, we apply  $T = 1,000$  and  $T' = 100, 200, \dots, 1,000$ , leading to 10 optimized structures. We use the same independent property predictor as in Xie et al. [9] to select the best one from the 10 structures. We calculate the success rate (SR) as the percentage of the 100 optimized structures achieving 5, 10, and 15 percentiles of the property distribution. We select the formation energy per atom as the target property, and provide the results on Perov-5, Carbon-24 and MP-20 in Table 12, showcasing the notable superiority of DiffCSP over the baseline methods.

Aside from the Carbon-24 dataset, the composition is flexible in the above experiments. We also attempt to follow the CSP setting and optimize the crystal structures for lower energy based on the fixed composition. We visualize eight cases in Figure 10 and calculate the energies of the structures before and after optimization by VASP [68]. Results show that our method is capable to search novel structures with lower energies compared to existing ones.

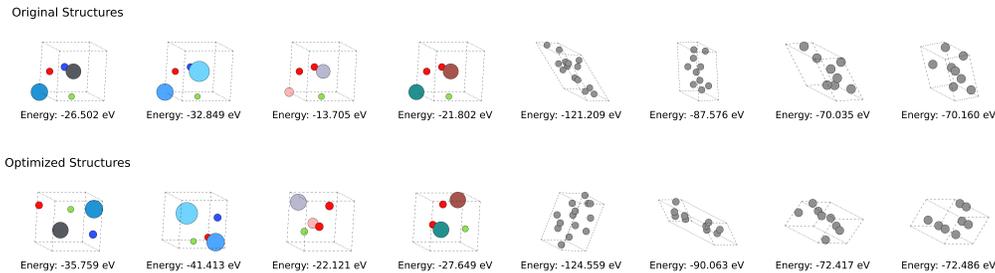


Figure 10: Visualization of 8 pairs of structures before and after optimization.

Table 13: GPU hours for yielding 20 candidates over the testing set.

	Perov-5	Carbon-24	MP-20	MPTS-52
P-cG-SchNet	2.1	3.0	10.0	22.5
CDVAE	21.9	9.8	95.2	178.0
DiffCSP	1.5	3.2	18.4	34.0

Table 14: Results on MP-20 with different inference steps.

	Steps	Match rate (%) $\uparrow$	RMSD $\downarrow$
DiffCSP	1,000	51.49	0.0631
	5,000	52.95	0.0541
CDVAE	1,000	30.71	0.1288
	5,000	33.90	0.1045

## H Computational Cost for Inference

We provide the GPU hours (GeForce RTX 3090) for different generative methods to predict 20 candidates on the 4 datasets. Table 13 demonstrates the diffusion-based models (CDVAE and our DiffCSP) are slower than P-cG-SchNet. Yet, the computation overhead of our method is still acceptable given its clearly better performance than P-cG-SchNet. Additionally, our DiffCSP is much faster than CDVAE across all datasets mainly due to the fewer generation steps. CDVAE requires 5,000 steps for each generation, whereas our approach only requires 1,000 steps. We further compare the performance of CDVAE and DiffCSP with 1,000 and 5,000 generation steps on MP-20 in Table 14. Our findings indicate that both models exhibit improved performance with an increased number of steps. Notably, DiffCSP with 1,000 steps outperforms CDVAE with 5,000 steps.

## I Error Bars

We provide a single run to generate 20 candidates in Table 1. We apply two more inferences of each generative method on Perov-5 and MP-20. Table 15 shows similar results as Table 1.

Table 15: Results on Perov-5 and MP-20 with error bars.

	# of samples	Perov-5		MP-20	
		Match rate (%) $\uparrow$	RMSE $\downarrow$	Match rate (%) $\uparrow$	RMSE $\downarrow$
P-cG-Schnet [53]	1	47.34 $\pm$ 0.63	0.4170 $\pm$ 0.0006	15.59 $\pm$ 0.41	0.3747 $\pm$ 0.0020
	20	97.92 $\pm$ 0.02	0.3464 $\pm$ 0.0004	32.70 $\pm$ 0.12	0.3020 $\pm$ 0.0002
CDVAE [9]	1	45.31 $\pm$ 0.49	0.1123 $\pm$ 0.0026	33.93 $\pm$ 0.15	0.1069 $\pm$ 0.0018
	20	88.20 $\pm$ 0.26	0.0473 $\pm$ 0.0007	67.20 $\pm$ 0.23	0.1012 $\pm$ 0.0016
DiffCSP	1	52.35 $\pm$ 0.26	0.0778 $\pm$ 0.0030	51.89 $\pm$ 0.30	0.0611 $\pm$ 0.0015
	20	<b>98.58<math>\pm</math>0.02</b>	<b>0.0129<math>\pm</math>0.0003</b>	<b>77.85<math>\pm</math>0.23</b>	<b>0.0493<math>\pm</math>0.0011</b>

## J More Visualizations

In this section, we first present additional visualizations of the predicted structures from DiffCSP and other generative methods in Figure 11. In line with Figure 3, our DiffCSP provides more accurate predictions compared with the baseline methods. Figure 12 illustrates 16 generated structures on Perov-5, Carbon-24 and MP-20. The visualization shows the capability of DiffCSP to generate diverse

structures. We further visualize the generation process of 5 structures from MP-20 in Figure 13. We find that the generated structure  $\mathcal{M}_0$  is periodically translated from the ground truth structure, indicating that the marginal distribution  $p(\mathcal{M}_0)$  follows the desired periodic translation invariance. We provide the detailed generation process in the Supplementary Materials.

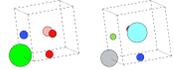
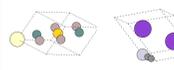
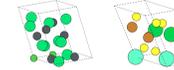
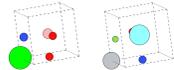
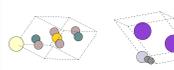
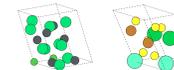
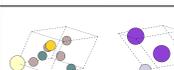
	Perov-5	Carbon-24	MP-20	MPTS-52
Ground Truth				
DiffCSP				
CDVAE				
P-cG-SchNet				

Figure 11: Additional visualizations of the predicted structures from different methods. We translate the same atom to the origin for better visualization and comparison.

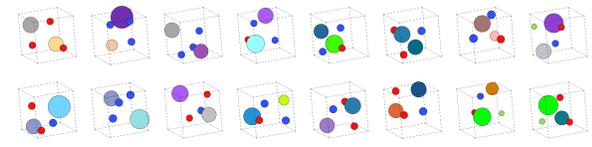
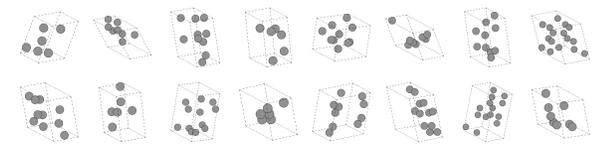
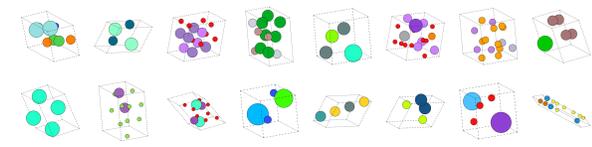
Perov-5	
Carbon-24	
MP-20	

Figure 12: Visualization of the generated structures by DiffCSP.

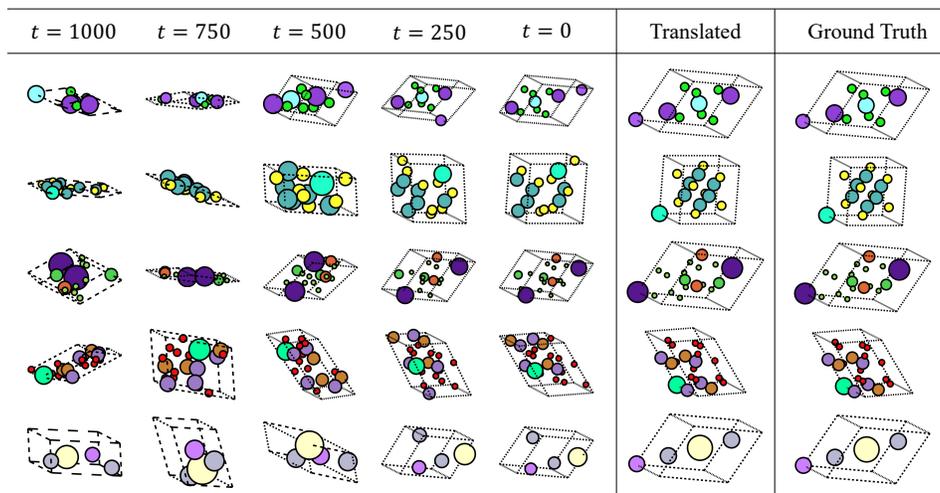


Figure 13: Visualization of the generation process on MP-20. The column “Translated” means translating the same atom in the generated structure  $\mathcal{M}_0$  to the origin as the ground truth for better comparison.