

PERIODIC MATERIALS GENERATION USING TEXT-GUIDED JOINT DIFFUSION MODEL

Anonymous authors

Paper under double-blind review

ABSTRACT

Equivariant diffusion models have emerged as the prevailing approach for generating novel crystal materials due to their ability to leverage the physical symmetries of periodic material structures. However, current models do not effectively learn the joint distribution of atom types, fractional coordinates, and lattice structure of the crystal material in a cohesive end-to-end diffusion framework. Also, none of these models work under realistic setups, where users specify the desired characteristics that the generated structures must match. In this work, we introduce TGDMat, a novel text-guided diffusion model designed for 3D periodic material generation. Our approach integrates global structural knowledge through textual descriptions at each denoising step while jointly generating atom coordinates, types, and lattice structure using a periodic-E(3)-equivariant graph neural network (GNN). Through extensive experiments with popular datasets on benchmark tasks, we first demonstrate that integrating textual knowledge significantly improves the material generation capabilities of existing state-of-the-art models. Furthermore, we show that TGDMat surpasses text-guided variants of existing baseline methods by a substantial margin, highlighting the effectiveness of our joint diffusion paradigm. Additionally, incorporating textual knowledge reduces overall training and sampling computational overhead while enhancing generative performance when utilizing real-world textual prompts from experts.

1 INTRODUCTION

Screening 3D periodic structures and their atomic compositions to identify novel crystal materials with specific chemical properties remains a long-standing challenge in the materials design community. These materials have been fundamental to key innovations such as the development of batteries, solar cells, semiconductors etc. (Butler et al., 2018; Desiraju, 2002). Historically, there have been attempts to generate novel materials by conducting resource-intensive and time-consuming simulations based on Density Functional Theory (DFT) (Kohn & Sham, 1965). Recently, the equivariant diffusion models (Jiao et al., 2023; Luo et al., 2023b; Xie et al., 2021) have demonstrated great potential to generate stable 3D periodic structures of new crystal materials.

However, these models possess several inherent limitations. 1) None of these existing SOTA models learns the joint distribution of atom coordinates, types, and lattice structure of the material through an end-to-end diffusion network. Existing models like CDVAE (Xie et al., 2021) and SyMat (Luo et al., 2023b) learn lattice parameters and atom types separately using a VAE model and further use a score network to learn the conditional distribution of atom coordinates given atom types and lattice. DiffCSP (Jiao et al., 2023), on the other hand, focuses primarily on structure prediction task where it assumes atom types are given and predict the stable crystal structure (lattice and coordinates). 2) Furthermore, these models use SE(3)-equivariant GNNs as backbone denoising network, which largely relies on messages passing around the local neighborhood of the atoms. Hence they fail to incorporate global structural knowledge into the diffusion process, which can enhance the diffusion performance. 3) Finally, these models are unconditional by design. From initial noisy structures without any external constraints, they generate stable crystal structures, which are distributionally similar to structures of the training dataset. This setup may have limited utility in real-world scenarios, as it lacks a mechanism for users to specify a criteria for the material to be generated. In a realistic setup, users would want to specify certain key details about the target material, like the chemical

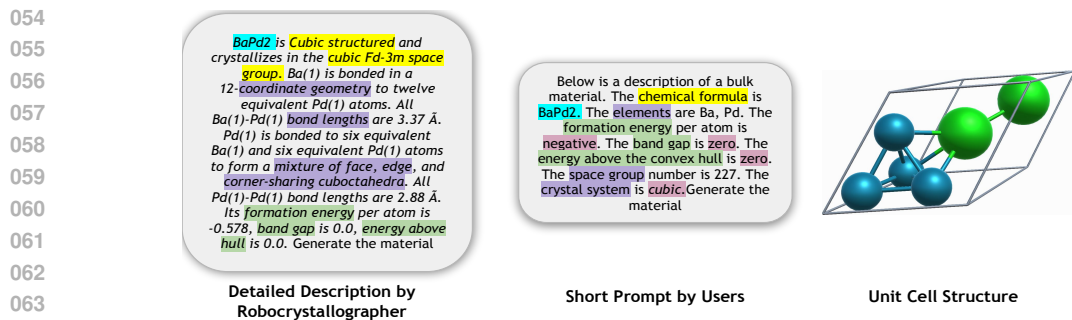


Figure 1: Detailed textual description generated by Robocrystallographer, less-detailed prompts by domain experts, and crystal unit cell structure of **BaPd₂**.

formula, space group, crystal symmetry, bond lengths, chemical properties, etc as input to the diffusion model, which the generated structure must then match.

In this paper, we propose, *TGDMat*, a novel *Text-Guided Diffusion Model for Material Generation* that mitigates the limitations mentioned above and enhances the generation capability. Though Text Guided Diffusion Models (TGDMs) produce impressively high-quality data in the form of images (Nichol et al., 2021; Ramesh et al., 2022; Rombach et al., 2022; Saharia et al., 2022), audio (Kreuk et al., 2022; Yang et al., 2022), video (Du et al., 2024), molecules (Gong et al., 2024; Luo et al., 2023a) etc, it remains largely unexplored in periodic material generation. Text-guided diffusion for new material generation has some key benefits. First, we can leverage popular tools like Robocrystallographer (Ganose & Jain, 2019) to generate a textual description of the material which provides a rich and diverse set of global structural knowledge like chemical formula, lattice constraint, space group number, crystal symmetry, chemical properties, etc. We believe this additional information is helpful for diffusion models in learning underlying crystal geometry. Second, it provides end users the flexibility to use custom prompts to guide the material generation process, ensuring that the resulting material aligns with the user’s provided description. Towards that goal, we first develop a diffusion model that jointly generates the atom coordinates, atom types, and lattice structure of crystal materials using a periodic E(3)-equivariant denoising model, satisfying periodic E(3) invariance properties of learned data distribution. Subsequently, we fuse textual information into the reverse diffusion process, which guides the denoising process in predicting material structure as specified by the textual description.

To sum up, our novel contributions in this work are as follows:

- To the best of our knowledge, we are the first to explore text-guided diffusion for material generation. Our proposed TGDMat bridges the gap between natural language understanding and material structure generation.
- Unlike prior models, TGDMat conducts joint diffusion on lattices, atom types, and coordinates, enhancing its ability to accurately capture the crystal geometry. Additionally, incorporating global structural knowledge through textual descriptions at each denoising step improves TGDMat’s ability to generate plausible materials with valid and stable structures.
- Through extensive experiments using popular datasets on benchmark tasks we show that text guidance can improve the generation capability of existing SOTA diffusion models for crystal materials. Moreover, in the generation task, TGDMat outperforms text-fusion variants of SOTA models with good margin, showcasing the effectiveness of the text guided joint diffusion paradigm.
- Fusing textual knowledge reduces the overall computational cost for both training and inference of the diffusion model. Moreover, when applied to real-world custom text prompts by experts, TGDMat demonstrates rich generative capability under general textual conditions.

	DiffCSP	TGDMat
Tasks	Only CSP Task	Both CSP and Gen Tasks
Diffusion on Atom Type	-	Discrete Diffusion (D3PM)
Model Category	Unconditional; unable to specify the criteria required by the user	Conditional; able to specify the criteria required by the user (in Text Format)
Text Guided Diffusion	No	Yes

Table 1: Key Differences between TGDMat from DiffCSP

2 PRELIMINARIES

2.1 CRYSTAL STRUCTURE REPRESENTATION

Crystal material can be modeled by a minimal *unit cell*, which gets repeated infinite times in 3D space on a regular lattice to form the periodic crystal structure. Given a material with N number of atoms in its unit cell, we can describe the unit cell by two matrices: *Atom Type Matrix* (A) and *Coordinate Matrix* (X). Atom Type Matrix $A = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_N]^T \in \mathbb{R}^{N \times k}$ denotes set of atomic type in one hot representation (k : maximum possible atom types). On the other hand, Coordinate Matrix $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T \in \mathbb{R}^{N \times 3}$ denotes atomic coordinate positions, where $\mathbf{x}_i \in \mathbb{R}^3$ corresponds to coordinates of i^{th} atom in the unit cell. Further, there is an additional *Lattice Matrix* $L = [\mathbf{l}_1, \mathbf{l}_2, \mathbf{l}_3]^T \in \mathbb{R}^{3 \times 3}$, which describes how a unit cell repeats itself in the 3D space towards $\mathbf{l}_1, \mathbf{l}_2$ and \mathbf{l}_3 direction to form the periodic 3D structure of the material. Formally, a given material can be defined as $M = (A, X, L)$ and we can represent its infinite periodic structure as $\hat{\mathbf{X}} = \{\hat{x}_i | \hat{x}_i = x_i + \sum_{j=1}^3 k_j l_j\}$; $\hat{\mathbf{A}} = \{\hat{a}_i | \hat{a}_i = a_i\}$ where $k_1, k_2, k_3, i \in \mathbb{Z}, 1 \leq i \leq N$.

2.2 INVARIANCES IN CRYSTAL STRUCTURE

The basic idea of using generative models for crystal generation is to learn the underlying data distribution of material structure $p(M)$. Since crystal materials satisfy physical symmetry properties (Dresselhaus et al., 2007; Zee, 2016), one of the major challenges here is the learned distribution must satisfy periodic E(3) invariance i.e. invariance to permutation, translation, rotation, and periodic transformations. A formal definition of these invariance properties is provided in Appendix C.

3 RELATED WORK: PERIODIC MATERIAL GENERATION

Recently, the majority of the research on material generation focuses on using popular generative models like VAEs (Kingma & Welling, 2013), GANs (Goodfellow et al., 2014) or Diffusion Models (Song & Ermon, 2019; 2020; Ho et al., 2020) to generate 3D periodic structures of materials (Hoffmann et al., 2019; Noh et al., 2019; Ren et al., 2020; Kim et al., 2020; Court et al., 2020; Long et al., 2021; Zhao et al., 2021; Xie et al., 2021; Jiao et al., 2023; Luo et al., 2023b; Zeni et al., 2023; Yang et al., 2023; Jiao et al., 2024; Miller et al., 2024). In specific, state-of-the-art models like CDVAE (Xie et al., 2021) and SyMat (Luo et al., 2023b) combine VAEs and score-based diffusion models to work directly with atomic coordinates, ensuring euclidean and periodic invariance using equivariant graph neural networks(GNNs). Moreover, DiffCSP (Jiao et al., 2023) focuses on structure prediction, jointly optimizing atom coordinates and lattice using a diffusion framework given atomic composition. We provided a comprehensive literature review of other related works in Appendix B.

Key differences between DiffCSP and TGDMat. We report key differences between DiffCSP and TGDMat in Table 1. The goal of this paper is not to introduce a new diffusion model to replace existing models like DiffCSP or CDVAE for periodic material generation. Instead, we focus on demonstrating that conditional models can outperform traditional unconditional models, such as DiffCSP. Specifically, we show that incorporating textual conditions through text-guided diffusion leads to better performance compared to using unconditional models like DiffCSP. Additionally, we enhance DiffCSP by integrating discrete diffusion over atom types in our proposed TGDMat.

162 4 METHODOLOGY

163 4.1 PROBLEM FORMULATION

164 In this work, given the textual description, we focus on generating a stable crystal structure that aligns
 165 with the provided textual description. Formally, given a dataset $\mathcal{M} = \{\mathbf{M}_i, \mathbf{T}_i\}$, containing crystal
 166 structure $\mathbf{M}_i = (\mathbf{A}_i, \mathbf{X}_i, \mathbf{L}_i)$ and its text description (\mathbf{T}_i), the goal of text guided crystal generation
 167 problem is to capture the underlying conditional data distribution $f(\mathbf{M}|\mathbf{T})$ via learning a generative
 168 model $p_\theta(\mathbf{M}|\mathbf{T})$, where θ is a set of learnable parameters. While training, we need p_θ to ensure that
 169 the learned distribution is invariant to different symmetry transformations mentioned in Section 2.2.
 170 Once trained, given a text description of a plausible material, the learned generative model can sample
 171 a valid and stable structure of the material, that is invariant to different symmetry transformations.
 172

173 4.2 TEXTUAL DATASETS

174 Leveraging textual information to guide the reverse diffusion process remains unexplored in the
 175 material design community. To the best of our knowledge, there is currently no text data available
 176 for materials in benchmark databases (mentioned in Section 5.1). Hence, we first curate the textual
 177 data of these material databases. Specifically, we propose two approaches for generating textual
 178 descriptions of materials, which are easy to follow. First, we utilize a freely available utility tool,
 179 *Robocrystallographer* (Ganose & Jain, 2019) to generate detailed textual descriptions about the
 180 periodic structure of crystal materials. These descriptions encompass local compositional details
 181 like atomic coordination, geometry, etc. as well as global structural aspects like crystal formula,
 182 mineral type, space group information, etc. Secondly, we utilized shorter and less detailed prompts
 183 that are more easily interpretable by users. We extend the prompt template proposed by (Gruver et al.,
 184 2024), which encodes minimal information about the material like its chemical formula, constituent
 185 elements, crystal system it belongs to, and its space group number. Further, we specify a few
 186 chemical properties, and instead of mentioning their actual values, we provide generic information
 187 like negative/positive formation energy, zero/nonzero band gaps, etc. Detailed information regarding
 188 the two textual datasets, including their curation process is provided in Appendix D.
 189

190 4.3 PROPOSED METHODOLOGY : TGD MAT

191 Our proposed model, TGD Mat (Fig. 2), uses an equivariant diffusion model guided by contextual
 192 representation of the textual description (\mathbf{C}_p) to generate a new crystal structure $\mathbf{M} = (\mathbf{A}, \mathbf{X}, \mathbf{L})$.
 193 Unlike prior methods (Jiao et al., 2023; Luo et al., 2023b; Xie et al., 2021), our method jointly diffuses
 194 $\mathbf{A}, \mathbf{X}, \mathbf{L}$ to learn the underlying data distribution of crystal structure $p(\mathbf{M}|\mathbf{C}_p)$. Diffusion models (Ho
 195 et al., 2020; Song & Ermon, 2019; 2020) are popular generative models that are formulated using
 196 a T steps Markov Chain. Given an input crystal material $\mathbf{M}_0 = (\mathbf{A}_0, \mathbf{X}_0, \mathbf{L}_0)$, the forward process
 197 gradually add noise to $\mathbf{A}_0, \mathbf{X}_0, \mathbf{L}_0$ independently over T steps and the reverse denoising process
 198 samples a noisy structure $\mathbf{M}_T = (\mathbf{A}_T, \mathbf{X}_T, \mathbf{L}_T)$ from a prior distribution and reconstruct back \mathbf{M}_0
 199 using some GNN model. At each t^{th} step of denoising ($0 \geq t \geq T$), the contextual representation of
 200 the crystal textual description (\mathbf{C}_p) will guide the diffusion process so that the intermediate structure
 201 \mathbf{M}_t aligns the target 3D structure constrained on textual conditions. Moreover, the learned distribution
 202 of material structure must satisfy periodic E(3) invariance. It is well studied in the literature (Xu
 203 et al., 2022) that if the prior distribution $p(x)$ is invariant to a group and the transition probabilities
 204 of a Markov chain $y \sim p(y|x)$ exhibit equivariance, the marginal distribution of y at any given time
 205 step also remains invariant to group transformations. Hence the learned distribution $p(\mathbf{M}_0)$ of the
 206 denoising model will satisfy periodic E(3) invariance if the prior distribution $p(\mathbf{M}_T)$ is invariant
 207 and the neural network used to parameterize the transition probability $q(\mathbf{M}_{t-1}|\mathbf{M}_t)$ is equivariant to
 208 permutational, translation, rotational, and periodic transformations. To satisfy that, we use periodic-
 209 E(3)-equivariant GNN model as a backbone denoising network to guide the denoising process. Next
 210 in this section, we first explain diffusion on \mathbf{M} in 4.3.1, then demonstrate the text-guided denoising
 211 network in 4.3.2 and finally training details in 4.4.
 212

213 4.3.1 JOINT EQUIVARIANT DIFFUSION ON \mathbf{M}

214 **Diffusion on Lattice (\mathbf{L}).** Since the Lattice Matrix $\mathbf{L} = [l_1, l_2, l_3]^T \in \mathbb{R}^{3 \times 3}$ is in continuous
 215 space, we leverage the idea of the Denoising Diffusion Probabilistic Model (DDPM) for diffusion

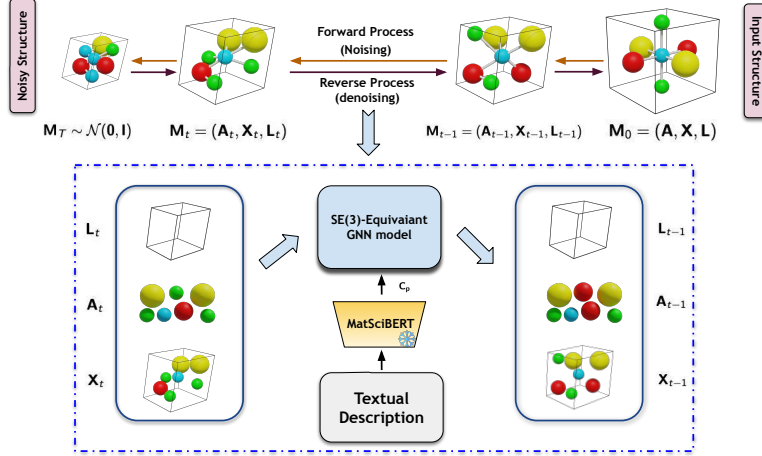


Figure 2: Model Architecture of our proposed text guided diffusion model TGDMat. At t^{th} step of reverse diffusion, given $M_t = (A_t, X_t, L_t)$, we use periodic-E(3)-equivariant GNN model guided by contextual representation of the textual prompts (C_p) to generate $M_{t-1} = (A_{t-1}, X_{t-1}, L_{t-1})$

on L . Specifically, given input lattice matrix $L_0 \sim p(L)$, at each t^{th} step, the forward diffusion process iteratively diffuses it through a transition probability $q(L_t|L_0)$ which can be derived as $q(L_t|L_0) = \mathcal{N}(L_t|\sqrt{\bar{\alpha}_t}L_0, (1 - \bar{\alpha}_t)\mathbf{I})$ where, $\bar{\alpha}_t = \prod_{k=1}^t \alpha_k$, $\alpha_t = 1 - \beta_t$ and $\{\beta_t \in (0, 1)\}_{t=1}^T$ controls the variance of diffusion step following certain noise scheduler. By reparameterization, we can rewrite $L_t = \sqrt{\bar{\alpha}_t}L_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon^L$ where, ϵ^L is noise sampled from $\mathcal{N}(0, \mathbf{I})$, added with L_0 at t^{th} step to generate L_t . After T such diffusion steps, noisy lattice matrix $L_T \sim \mathcal{N}(0, \mathbf{I})$ is generated. During reverse denoising process, given noisy $L_T \sim \mathcal{N}(0, \mathbf{I})$ we reconstruct true lattice structure L_0 through iterative denoising step via learning reverse conditional distribution, which we formulate as $p(L_{t-1}|M_t, C_p) = \mathcal{N}\{L_{t-1}|\mu^L(M_t, C_p), \beta_t \frac{(1 - \bar{\alpha}_{t-1})}{(1 - \bar{\alpha}_t)}\mathbf{I}\}$ where $\mu^L(M_t, C_p) = \frac{1}{\sqrt{\bar{\alpha}_t}}(L_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\hat{\epsilon}^L(M_t, C_p, t))$. Intuitively, $\hat{\epsilon}^L$ needs to be subtracted from L_t to generate L_{t-1} and textual representation C_p will steer this reverse diffusion process. We use a text-guided denoising network $\Phi_\theta(A_t, X_t, L_t, t, C_p)$ to model the noise term $\hat{\epsilon}^L(M_t, C_p, t)$. Following the simplified training objective proposed by (Ho et al., 2020), we train denoising model using l_2 loss between $\hat{\epsilon}^L$ and ϵ^L

$$\mathcal{L}_{lattice} = \mathbb{E}_{\epsilon^L, t \sim \mathcal{U}(1, T)} \|\epsilon^L - \hat{\epsilon}^L\|_2^2 \quad (1)$$

Diffusion on Atom Types (A). Prior studies (Jiao et al., 2023; Xie et al., 2021) consider Atom Type Matrix A as the probability distribution for k classes $\in \mathbb{R}^{N \times k}$ (continuous variable) and apply DDPM to learn the distribution. However for discrete data these models are inappropriate and produce suboptimal results (Austin et al., 2021; Campbell et al., 2022). Hence we consider A as N discrete variables belonging to k classes and leverage discrete diffusion model (D3PM) (Austin et al., 2021) for diffusion on A . In specific, with a as the one-hot representation of atom a , the transition probability for the forward process is $q(a_t|a_{t-1}) = \text{Cat}(a_t; \mathbf{p} = a_{t-1}\mathbf{Q}_t)$, where $\text{Cat}(a; \mathbf{p})$ is a categorical distribution over a with probabilities \mathbf{p} and \mathbf{Q}_t is the Markov transition matrix at time step t , defined as $[\mathbf{Q}_t]_{i,j} = q(a_t = i|a_{t-1} = j)$. Different choices of \mathbf{Q}_t and corresponding stationary distributions are proposed by (Austin et al., 2021) which provides flexibility to control the data corruption and denoising process. We adopted the absorbing state diffusion process, introducing a new absorbing state [MASK] in \mathbf{Q}_t . At each time step t , an atom either stays in its type state with probability $1 - \beta_t$ or moves to [MASK] state with probability β_t and once it moves to [MASK] state, it stays there. Hence, the stationary distribution of this diffusion process has all the mass on the [MASK] state. During denoising process, given textual representation C_p , we first sample noisy a_T and obtain a_0 through iterative denoising step via learning reverse conditional transition $p_\theta(a_{t-1}|a_t, C_p) \propto \sum_{a_0} q(a_{t-1}, a_t|a_0)p_\theta(a_0|a_t, C_p)$. We use the text-guided denoising network $\Phi_\theta(A_t, X_t, L_t, t, C_p)$ to model this denoising process, which is trained using following loss function :

$$\mathcal{L}_{type} = \mathcal{L}_{VB} + \lambda\mathcal{L}_{CE} \quad (2)$$

where \mathcal{L}_{VB} and \mathcal{L}_{CE} is the variational lower bound and cross-entropy loss respectively and λ is a hyperparameter. Details about the diffusion process and the losses \mathcal{L}_{VB} , \mathcal{L}_{CE} are in Appendix E

Diffusion on Atom Coordinates (X). We can diffuse the Coordinate Matrix $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T \in \mathbb{R}^{N \times 3}$ in two ways: either by diffusing cartesian coordinates or fractional coordinates. Prior works like CDVAE (Xie et al., 2021) and SyMat (Luo et al., 2023b) diffuse cartesian coordinates whereas DiffCSP (Jiao et al., 2023) diffuses fractional coordinates. In our setup, as we are jointly learning atom coordinates and lattice matrix, hence we follow DiffCSP and diffuse fractional coordinates. Fractional coordinates in crystal material resides in quotient space $\mathbb{R}^{N \times 3} / \mathbb{Z}^{N \times 3}$ induced by the crystal periodicity. Since the Gaussian distribution used in DDPM is unable to model the cyclical and bounded domain of X , it is not suitable to apply DDPM to model X . Hence at each step of forward diffusion, we add noise sampled from Wrapped Normal (WN) distribution (De Bortoli et al., 2022) to X and during denoising leverage Score Matching Networks (Song & Ermon, 2019; 2020) to model underlying transition probability $q(X_t|X_0) = \mathcal{N}_W(X_t|X_0, \sigma_t^2 \mathbf{I})$. In specific, at each t^{th} step of diffusion, we derive X_t as : $X_t = f_w(X_0 + \sigma_t \epsilon^X)$ where, $\epsilon^X \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, σ_t is the noise scheduler and $f_w(\cdot)$ is a truncation function. Given a fractional coordinate matrix X , truncation function $f_w(X) = (X - \lfloor X \rfloor)$ returns the fractional part of each element of X .

As argued in (Jiao et al., 2023), $q(X_t|X_0)$ is periodic translation equivariant, and approaches uniform distribution $\mathcal{U}(0, 1)$ for sufficiently large values of σ_T . Hence during the denoising process, we first sample $X_T \sim \mathcal{U}(0, 1)$ and iteratively denoise via score network for T steps to recover back the true fractional coordinates X_0 . We use the text-guided denoising network $\Phi_\theta(\mathbf{A}_t, X_t, \mathbf{L}_t, t, \mathbf{C}_p)$ to model the denoising process, which is trained using the following score-matching objective function :

$$\mathcal{L}_{coord} = \mathbb{E}_{\substack{X_t \sim q(X_t|X_0) \\ t \sim \mathcal{U}(1, T)}} \|\nabla_{X_t} \log q(X_t|X_0) - \hat{\epsilon}^X(\mathbf{M}_t, \mathbf{C}_p, t)\|_2^2 \quad (3)$$

where $\nabla_{X_t} \log q(X_t|X_0) \propto \sum_{\mathbf{K} \in \mathbb{Z}^{N \times 3}} \exp(-\frac{\|X_t - X_0 + \mathbf{K}\|_F^2}{2\sigma_t^2})$ is the score function of transitional distribution and $\hat{\epsilon}^X(\mathbf{M}_t, \mathbf{C}_p, t)$ denoising term. More Details are provided in Appendix E

4.3.2 TEXT GUIDED DENOISING NETWORK

In this subsection, we will illustrate the detailed architecture of our proposed Text Guided Denoising Network $\Phi_\theta(\mathbf{A}_t, X_t, \mathbf{L}_t, t, \mathbf{C}_p)$, which we used during denoising process to generate \mathbf{A} , X and \mathbf{L} . As mentioned in 2.2, the learned distribution of material structure $p(\mathbf{M})$ must satisfy periodic E(3) invariance. Hence we leverage a periodic-E(3)-equivariant Graph Neural Network (GNN) integrated with a pre-trained textual encoder to model the denoising process. In particular, as a text encoder, we adopt a pre-trained MatSciBERT (Gupta et al., 2022) model, which is a domain-specific language model for materials science, followed by a projection layer. MatSciBERT is effectively a pre-trained SciBERT model on a scientific text corpus of 3.17B words, which is further trained on a huge text corpus of materials science containing around 285M words. We feed textual description of material T into MatSciBERT and extract embedding of [CLS] token \mathbf{h}_{CLS} as a representation of the whole text. Further, we feed \mathbf{h}_{CLS} through a projection layer to generate the contextual textual embedding for the material $\mathbf{C}_p \in \mathbb{R}^d$, which we pass to the equivariant GNN model to guide the denoising process. Practically, as the backbone network for the denoising process, we extend CSPNet architecture (Jiao et al., 2023), originally developed for crystal structure prediction (CSP) task. CSPNet is built upon EGNN (Satorras et al., 2021), satisfying periodic E(3) invariance condition on periodic crystal structure. At the k^{th} layer message passing, the Equivariant Graph Convolutional Layer (EGCL) takes as input the set of atom embeddings $\mathbf{h}^k = [\mathbf{h}_1^k, \mathbf{h}_2^k, \dots, \mathbf{h}_N^k]$, atom coordinates $\mathbf{x}^k = [\mathbf{x}_1^k, \mathbf{x}_2^k, \dots, \mathbf{x}_N^k]$ and Lattice Matrix \mathbf{L} and outputs a transformation on \mathbf{h}^{k+1} . Formally, we can define the k^{th} layer message passing operation as:

$$\mathbf{m}_{i,j} = \rho_m\{\mathbf{h}_i^k, \mathbf{h}_j^k, \mathbf{L}^T \mathbf{L}, \psi_{FT}(\mathbf{x}_i^k - \mathbf{x}_j^k)\}; \mathbf{m}_i = \sum_{j=1}^N \mathbf{m}_{i,j}; \mathbf{h}_i^{k+1} = \mathbf{h}_i^k + \rho_h\{\mathbf{h}_i^k, \mathbf{m}_i\} \quad (4)$$

where ρ_m, ρ_h are MLPs and ψ_{FT} is a Fourier Transformation function applied on relative difference between fractional coordinates $\mathbf{x}_i^k, \mathbf{x}_j^k$. Fourier Transformation is used since it is invariant to periodic translation and extracts various frequencies of all relative fractional distances that are helpful for crystal structure modeling (Jiao et al., 2023).

We fuse textual representation \mathbf{C}_p into input atom feature \mathbf{h}_i^0 as

$$\mathbf{h}_i^0 = \rho\{f_{atom}(\mathbf{a}_i) \parallel f_{pos}(t) \parallel \mathbf{C}_p\} \quad (5)$$

Dataset	Method	Validity \uparrow		Coverage \uparrow		Property Statistics (EMD) \downarrow		
		Compositional(%)	Structural(%)	COV-R(%)	COV-P(%)	# Element	ρ	\mathcal{E}
Perov-5	CDVAE	98.29	100	99.25	98.39	0.0731	0.1462	0.0291
	CDVAE+	98.45	100	99.53	99.09	0.0609	0.1276	0.0223
	SyMat	96.83	100	99.16	98.29	0.0193	0.1991	0.2827
	SyMat+	97.88	100	99.70	98.79	0.0172	0.1755	0.2566
	DiffCSP	98.15	100	99.28	98.08	0.0132	0.1280	0.0267
	DiffCSP+	98.44	100	99.85	98.53	0.0119	0.1070	0.0241
Carbon-24	CDVAE	-	100	99.35	82.66	-	0.1539	0.2889
	CDVAE+	-	100	99.82	84.76	-	0.1377	0.2660
	SyMat	-	100	99.42	97.17	-	0.1234	3.9628
	SyMat+	-	100	99.90	97.63	-	0.1171	3.8620
	DiffCSP	-	99.9	99.49	97.27	-	0.0861	0.0876
	DiffCSP+	-	100	99.93	97.33	-	0.0763	0.0853
MP-20	CDVAE	86.30	100	99.15	99.49	1.4921	0.7085	0.3039
	CDVAE+	87.42	100	99.57	99.81	0.9720	0.6388	0.2977
	SyMat	87.96	99.9	98.30	99.37	0.5236	0.4012	0.3877
	SyMat+	88.47	99.9	99.01	99.95	0.4865	0.3879	0.3489
	DiffCSP	83.25	100	99.41	99.76	0.3411	0.3802	0.1497
	DiffCSP+	85.07	100	99.81	99.89	0.3122	0.3799	0.1355

Table 2: Summary of comparative results on *Gen* task between SOTA diffusion models (M) and their text-guided variants (M+). We highlight the best performances for each class of models in bold. The table contains ”-” values for metrics that do not apply to certain datasets.

where t is the timestamp of the diffusion model, $f_{pos}(\cdot)$ is sinusoidal positional encoding (Ho et al., 2020; Vaswani et al., 2017), $f_{atom}(\cdot)$ learned atomic embedding function and \parallel is concatenation operation. Input atom features \mathbf{h}^0 and coordinates \mathbf{x}^0 are fed through \mathcal{K} layers of EGCL to produce $\hat{\mathbf{e}}^L$, $p(\mathbf{A}_{t-1} | \mathbf{M}_t)$ and $\hat{\mathbf{e}}^X$ as follows :

$$\hat{\mathbf{e}}^L = \mathbf{L}\rho_L\left(\frac{1}{N} \sum_{i=1}^N \mathbf{h}^K\right); p(\mathbf{A}_{t-1} | \mathbf{M}_t) = \rho_A(\mathbf{h}^K); \hat{\mathbf{e}}^X = \rho_X(\mathbf{h}^K) \quad (6)$$

where ρ_L, ρ_A, ρ_X are MLPs on the final layer embeddings. Intuitively, we feed global structural knowledge about the crystal structure into the network by injecting contextual representation \mathbf{C}_p into input atom features. This added signal participates through message-passing operations in Eq. 4 and guides in denoising atom types, coordinates, and lattice parameters such that it can capture the global crystal geometry and aligned with the input stable structure specified by textual description.

4.4 TRAINING AND SAMPLING

TGDMat is trained using the following combined loss: $\mathcal{L} = \lambda_L \mathcal{L}_{lattice} + \lambda_A \mathcal{L}_{type} + \lambda_X \mathcal{L}_{coord}$ where $\mathcal{L}_{lattice}$, \mathcal{L}_{type} and \mathcal{L}_{coord} are lattice l_2 loss (Eq. 1), type loss (Eq. 2) and coordinate score matching loss (Eq. 3) respectively and $\lambda_L, \lambda_A, \lambda_X$ are hyperparameters control the relative weightage between these different loss components. During training, we freeze the MatSciBERT parameters and do not tune them further. During sampling, we use the Predictor-Corrector (Song et al., 2020) sampling mechanism to sample $\mathbf{A}_0, \mathbf{X}_0$ and \mathbf{L}_0 . Training/Sampling algorithms are provided in Appendix E.5

5 EXPERIMENTS

5.1 BENCHMARK TASKS, EVALUATION METRICS AND DATASETS

Following the prior works (Jiao et al., 2023; Xie et al., 2021), we evaluate our proposed model TGDMat on two benchmark tasks for material generation, *Random Material Generation (Gen)* and *Crystal Structure Prediction (CSP)*. In *Gen* task, the goal of the generative model is to generate novel stable materials (atom types, coordinates, and lattice). In *CSP* task, atom types of the materials are given and the goal is to predict/match the crystal structure (atom coordinates and lattice). A visual representation of both tasks is provided in Figure 5 of Appendix F.1. In TGDMat model, by design choice, we use the textual description of crystal materials during each step of the reverse diffusion process to enhance the generation capability in both tasks. Following (Jiao et al., 2023; Xie et al., 2021), for *Gen* task, we evaluate the performance using seven metrics under three broad categories:

Method	# Samples Generated/ test data	Perov-5		Carbon-24		MP-20	
		Match Rate \uparrow	RMSE \downarrow	Match Rate \uparrow	RMSE \downarrow	Match Rate \uparrow	RMSE \downarrow
CDVAE	1	45.31	0.1138	17.09	0.2969	33.90	0.1045
	20	88.51	0.0464	88.37	0.2286	66.95	0.1026
CDVAE+	1	49.25	0.1055	23.73	0.2590	41.80	0.1021
	20	89.73	0.0417	89.77	0.2053	72.56	0.0840
SyMat	1	47.32	0.1074	20.81	0.2655	33.92	0.1039
	20	90.25	0.0316	89.29	0.2184	71.03	0.0945
SyMat++	1	50.88	0.0963	28.18	0.2510	43.17	0.1016
	20	92.30	0.0201	91.65	0.1870	72.96	0.0820
DiffCSP	1	52.02	0.0760	17.54	0.2759	51.49	0.0631
	20	98.60	0.0128	88.47	0.2192	77.93	0.0492
DiffCSP++	1	90.46	0.0203	44.63	0.2266	55.15	0.0572
	20	98.59	0.0072	95.27	0.1534	82.02	0.0391

Table 3: Summary of comparative results on CSP task between SOTA diffusion models (M) and their text-guided variants (M+). We highlight the best performances for each class of models in bold.

Validity, Coverage, and Property Statistics, whereas for evaluating the performance of CSP task we use **Match Rate (MR)** and **RMSE**. For evaluation, we use three popular material datasets: **Perov-5** (Castelli et al., 2012a;b), **Carbon-24** (Pickard., 2020) and **MP-20** (Jain et al., 2013b). We curated textual data for these datasets with a textual description of each material. Specifically, we generate both long detailed textual descriptions and shorter prompts using approaches mentioned in 4.2. While training TGDMat, we split the datasets into the train, test, and validation sets following the convention of 60:20:20 (Xie et al., 2021). More details are in Appendix F.1 and F.2.

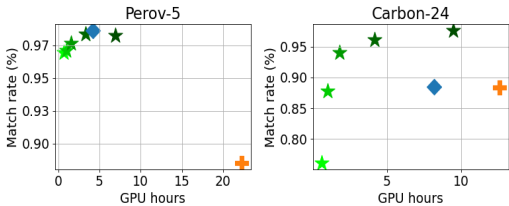
5.2 EFFICACY OF TEXT-GUIDANCE IN DIFFUSION

Setup. We begin by examining whether text guidance during reverse denoising process in diffusion model can improve the generation of stable periodic structures of crystal materials. Specifically, we compared state-of-the-art (unconditional) diffusion models with their text-guided variants across both benchmark tasks. We choose three popular state-of-the-art generative models: **CDVAE** (Xie et al., 2021), **SyMat** (Luo et al., 2023b), **DiffCSP** (Jiao et al., 2023) and build their text-guided variants named **CDVAE+**, **SyMat+**, **DiffCSP+** respectively, where we fuse the contextual representation of (long detailed) text data into denoising network of those models using our proposed algorithm as described in 4.3.2. For the CSP task, we generate k samples for each material structure in the test set using baseline models, and we determine the matching metrics (MR and RMSE) if at least one sample aligns with the ground truth structure. In the Gen task, the text-guided variants are constrained by textual prompts for generating new materials. To ensure a fair comparison regarding sample size, we generate a number of samples equal to the test data size for all baseline models, both unconditional and text-guided variants, and evaluate their performance accordingly.

Results and Discussions. We report the results of Gen and CSP task in Table 2 and 3 respectively. We observe that for both the tasks, text-guided variant of any SOTA model surpasses the vanilla variant with a good margin across three datasets. In specific, for Gen task, text-guided models consistently outperform respective vanilla models across all seven metrics. For CSP task, while prior unconditional diffusion models demonstrate improved Match Rates and lower RMSE when generating 20 samples ($k = 20$) per test material, they largely fail in both metrics when generating only one sample ($k = 1$) per test material. However, generating 20 samples per test material to match the structure is unrealistic and computationally burdensome. Using text guidance during the reverse denoising process, with just one generated sample per test material, text-guided variants outperform respective vanilla models, thereby reducing computational overhead. Moreover, even with 20 generated samples, the performance of text-guided models are better for all the benchmark datasets. Overall, in the CSP task, DiffCSP+ emerges as the top model compared to all baseline models and their text-guided variants across three benchmark datasets. These findings across different SOTA models collectively highlight the importance of text-guided diffusion: incorporating textual knowledge during reverse diffusion aids in aligning the noisy structure with the 3D geometry of stable realistic materials, enhancing both stable crystal structure prediction and the random generation task.

Dataset	Method	Validity \uparrow		Coverage \uparrow		Property Statistics (EMD) \downarrow		
		Compositional(%)	Structural(%)	COV-R(%)	COV-P(%)	# Element	ρ	\mathcal{E}
Perov-5	CDVAE+	98.45	100	99.53	99.09	0.0609	0.1276	0.0223
	SyMat+	97.88	100	99.70	98.79	0.0172	0.1755	0.2566
	DiffCSP+	98.44	100	99.85	98.53	0.0119	0.1070	0.0241
	TGDMat(Short)	98.28	100	99.71	99.24	0.0108	0.0947	0.0237
	TGDMat(Long)	98.63	100	99.87	99.52	0.0090	0.0497	0.0187
Carbon-24	CDVAE+	-	100	99.82	84.76	-	0.1377	0.2660
	SyMat+	-	100	99.90	97.63	-	0.1171	3.8620
	DiffCSP+	-	100	99.93	97.33	-	0.0763	0.0853
	TGDMat(Short)	-	100	99.81	91.77	-	0.0681	0.0865
	TGDMat(Long)	-	100	99.91	92.43	-	0.0436	0.0632
MP-20	CDVAE+	87.42	100	99.57	99.81	0.9720	0.6388	0.2977
	SyMat+	88.47	99.9	99.01	99.95	0.4865	0.3879	0.3489
	DiffCSP+	85.07	100	99.81	99.89	0.3122	0.3799	0.1355
	TGDMat(Short)	86.60	100	99.79	99.88	0.3337	0.3296	0.1189
	TGDMat(Long)	92.97	100	99.89	99.95	0.2890	0.3082	0.1154

Table 4: Summary of comparative results on *Gen* task between text-guided SOTA models and TGDMat. We highlight the best and second-best performances in bold and underlined, respectively. The table contains “-” values for metrics that do not apply to certain datasets.



(a) Match Rate vs Running time

Figure 3: (a) Match Rate vs Running time (GPU Hours) for different variants of TGDMat(Long) {50 Steps \star , 100 Steps \star , 200 Steps \star , 500 Steps \star , 1K Steps \star }, DiffCSP+ \blacklozenge and CDVAE+ \blackplus .

5.3 EFFECTIVENESS OF TGDMAT

Setup. Next, we demonstrate the effectiveness of our proposed TGDMat framework, which jointly denoises atom types, coordinates, and lattice structures of crystal materials and during the denoising process, integrates contextual representation of the textual description with backbone GNN network to guide the diffusion process. Here we only consider *Gen* task, since for *CSP* task, the atom type of the material is given, and in that setup, TGDMat jointly denoising only atom coordinates and lattice structures with textual guidance, which aligns with the DiffCSP+ framework we discussed earlier in 5.2. To compare the performance of TGDMat in *Gen* task, we choose the aforementioned text-guided variants of popular SOTA models: **CDVAE+**, **SyMat+** and **DiffCSP+**.

Results and Discussions. We report the result in Table 4. We trained TGDMat using both detailed textual descriptions and short prompts, as outlined in 4.2 and report them as **TGDMat(Long)** and **TGDMat(Short)** respectively in Table 4. We observe that both variants of TGDMat consistently enhances performance across almost all metrics across the benchmark datasets. Particularly on the Perov-5 and MP-20 datasets, TGDMat outperforms all baseline models across all metrics. In the Carbon-24 dataset, TGDMat exhibits performance improvements across all metrics except for COV-P and COV-R, where its performance is on par with state-of-the-art results. Additionally, our experiments indicate that utilizing shorter prompts results in a slight decrease in overall performance compared to the longer variant. Nonetheless, the performance remains superior or comparable to baseline models. Notably, TGDMat’s superior performance across seven metrics on MP-20 highlights its potential to generate novel materials that can be synthesized experimentally. Overall, TGDMat exhibits promising performance in *Gen* task, indicating the benefits of using text-guided joint diffusion to learn \mathbf{A} , \mathbf{X} , \mathbf{L} and generate more stable periodic structures of 3D crystal materials.

Additionally, we present visualizations of a few generated materials based on the textual descriptions in Table ?? and compare them with the ground truth material structure. The generated samples exhibit clear matches to the ground truth structure, highlighting the generation capability of the TGDMat model given information in text form. More visualization results are in Appendix F.9.

5.4 CORRECTNESS OF GENERATED MATERIALS

Setup. In this section, we investigate whether the generated material matches different features specified by the textual prompts. TGDMat has the capability to process textual prompts given by the user, enabling it to manage global attributes about crystal materials such as Formula, Space group, Crystal System, and different property values like formation energy, band-gap, etc. To ensure the fidelity of our model’s outputs concerning these specified global attributes from the text prompt, We randomly generated 1000 materials (sampled from all three Datasets) based on their respective textual descriptions(both Long and Short) and assessed the percentage of generated materials that matched the global features outlined in the text prompt. In specific, we matched the Formula, Space group, and Crystal System, and Dimensions of generated materials with the textual descriptions. Moreover, we examined whether properties such as formation energy and bandgap matched the specified criteria as per the text prompt (positive/negative, zero/nonzero).

Results and Discussions. We report the results for long prompt in Table 5 and short prompts in in Table 10. In general, using longer text, considering Perov-5 and Carbon-24 datasets, the generated material meets the specified criteria effectively. However, when dealing with the MP-20 dataset, which is more intricate due to its complex structure and composition, performance tends to decline. Additionally, when using shorter prompts, overall performance suffers across all datasets compared to longer text inputs. This is because the longer text, provided by the robocrystallographer, offers a comprehensive range of information, both global and local

	Global Features in Text Prompt	% of Matched Materials		
		Perov-5	Carbon-24	MP-20
TGDMat(Long)	Formula	97.50	98.20	70.54
	Space Group	87.00	80.79	67.88
	Crystal System	92.60	91.55	73.54
	Formation Energy	95.49	-	92.88
	Band Gap	-	98.61	96.73

Table 5: Correctness of generated materials matching conditions specified by the textual prompts.

5.5 COMPUTATIONAL COST FOR TRAINING AND SAMPLING

Integrating text knowledge during reverse diffusion for material generation offers a key advantage: it accelerates convergence towards realistic structures and reduces computational overhead. We observe, compared to other baseline models, TGDMat incurs substantially lower computation costs during training and sampling. Notably, our approach cuts down on training time, requiring only 500 epochs compared to 3K or 4K epochs for CDVAE and DiffCSP on Perov-5 and Carbon-24 respectively. Additionally, our method reduces sampling steps, making it faster to generate new structures. While CDVAE and DiffCSP need 5K and 1K steps respectively, our model only requires 500 steps. We compare the performance of CDVAE and DiffCSP with different TGDMat(Long) variants with 50, 100, 200, 500, and 1K steps and report the match rate of the predicted crystal structure vs running time (GPU hours in P100 GPU server) for Perov-5 and Carbon-24 datasets in Fig. 3(a). We notice that the inference time for CDVAE is lengthier as it necessitates 5K steps for each generation. However, for Carbon-24, TGDMat with 200 or 500 steps outperforms DiffCSP with 1K steps. Additionally, for Perov-5, TGDMat with 500 steps achieves results comparable to DiffCSP with 1K steps.

6 CONCLUSION

In this work, we explore a practical approach of generating stable crystal materials given a textual description of the material. We propose TGDMat, which jointly diffuse atom types, fractional coordinates, and lattice structure for crystal materials using a periodic-E(3)-equivariant denoising model. We further integrate textual information into the reverse diffusion process through a pre-trained transformer model, which guides the denoising process in learning the crystal 3D geometry matching the specification by textual description. Extensive experiments conducted on two benchmark generative tasks reveal that TGDMat surpasses all popular baseline models by a good margin. Furthermore, integrating textual knowledge reduces the overall computational cost for both training and inference of the diffusion model. Moreover, when applied to real-world custom text prompts by experts, TGDMat demonstrates rich generative capability under general textual conditions.

REFERENCES

- 540
541
542 Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured
543 denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing*
544 *Systems*, 34:17981–17993, 2021.
- 545
546 Keith T Butler, Daniel W Davies, Hugh Cartwright, Olexandr Isayev, and Aron Walsh. Machine
547 learning for molecular and materials science. *Nature*, 559(7715):547–555, 2018.
- 548
549 Andrew Campbell, Joe Benton, Valentin De Bortoli, Thomas Rainforth, George Deligiannidis, and
550 Arnaud Doucet. A continuous time framework for discrete denoising models. *Advances in Neural*
551 *Information Processing Systems*, 35:28266–28279, 2022.
- 552
553 Ivano E Castelli, David D Landis, Kristian S Thygesen, Søren Dahl, Ib Chorkendorff, Thomas F
554 Jaramillo, and Karsten W Jacobsen. New cubic perovskites for one-and two-photon water splitting
555 using the computational materials repository. *Energy & Environmental Science*, 5(10):9034–9043,
2012a.
- 556
557 Ivano E Castelli, Thomas Olsen, Soumendu Datta, David D Landis, Søren Dahl, Kristian S Thygesen,
558 and Karsten W Jacobsen. Computational screening of perovskite metal oxides for optimal solar
light capture. *Energy & Environmental Science*, 5(2):5814–5819, 2012b.
- 559
560 Chi Chen, Weike Ye, Yunxing Zuo, Chen Zheng, and Shyue Ping Ong. Graph networks as a universal
561 machine learning framework for molecules and crystals. *Chem. Mater.*, 31(9):3564–3572, 2019.
- 562
563 Kamal Choudhary and Brian DeCost. Atomistic line graph neural network for improved materials
564 property predictions. *npj Computational Materials*, 7(1):1–8, 2021.
- 565
566 Callum J Court, Batuhan Yildirim, Apoorv Jain, and Jacqueline M Cole. 3-d inorganic crystal
567 structure generation and property prediction via representation learning. *Journal of Chemical*
Information and Modeling, 60(10):4518–4535, 2020.
- 568
569 Kishalay Das, Bidisha Samanta, Pawan Goyal, Seung-Cheol Lee, Satadeep Bhattacharjee, and Niloy
570 Ganguly. Crysxpp: An explainable property predictor for crystalline materials. *npj Computational*
Materials, 8(1):43, 2022.
- 571
572 Kishalay Das, Pawan Goyal, Seung-Cheol Lee, Satadeep Bhattacharjee, and Niloy Ganguly. Crys-
573 mmnet: multimodal representation for crystal property prediction. In *Uncertainty in Artificial*
574 *Intelligence*, pp. 507–517. PMLR, 2023a.
- 575
576 Kishalay Das, Bidisha Samanta, Pawan Goyal, Seung-Cheol Lee, Satadeep Bhattacharjee, and Niloy
577 Ganguly. Crysnn: Distilling pre-trained knowledge to enhance property prediction for crystalline
578 materials. *arXiv preprint arXiv:2301.05852*, 2023b.
- 579
580 Daniel W Davies, Keith T Butler, Adam J Jackson, Jonathan M Skelton, Kazuki Morita, and Aron
581 Walsh. Smact: Semiconducting materials by analogy and chemical theory. *Journal of Open Source*
Software, 4(38):1361, 2019.
- 582
583 Valentin De Bortoli, Emile Mathieu, Michael Hutchinson, James Thornton, Yee Whye Teh, and
584 Arnaud Doucet. Riemannian score-based generative modelling. *Advances in Neural Information*
Processing Systems, 35:2406–2422, 2022.
- 585
586 Gautam R Desiraju. Cryptic crystallography. *Nature materials*, 1(2):77–79, 2002.
- 587
588 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
589 bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 590
591 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*
in neural information processing systems, 34:8780–8794, 2021.
- 592
593 Mildred S Dresselhaus, Gene Dresselhaus, and Ado Jorio. *Group theory: application to the physics*
of condensed matter. Springer Science & Business Media, 2007.

- 594 Yilun Du, Sherry Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Josh Tenenbaum, Dale Schuurmans, and
595 Pieter Abbeel. Learning universal policies via text-guided video generation. *Advances in Neural*
596 *Information Processing Systems*, 36, 2024.
- 597 Octavian Ganea, Lagnajit Pattanaik, Connor Coley, Regina Barzilay, Klavs Jensen, William Green,
598 and Tommi Jaakkola. Geomol: Torsional geometric generation of molecular 3d conformer
599 ensembles. *Advances in Neural Information Processing Systems*, 34:13757–13769, 2021.
- 600 Alex M Ganose and Anubhav Jain. Robocrystallographer: automated crystal structure text descrip-
601 tions and analysis. *MRS Communications*, 9(3):874–881, 2019.
- 602 Haisong Gong, Qiang Liu, Shu Wu, and Liang Wang. Text-guided molecule generation with diffusion
603 language model. *arXiv preprint arXiv:2402.13040*, 2024.
- 604 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
605 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information*
606 *processing systems*, 27, 2014.
- 607 Nate Gruver, Anuroop Sriram, Andrea Madotto, Andrew Gordon Wilson, C Lawrence Zitnick, and
608 Zachary Ulissi. Fine-tuned language models generate stable inorganic materials as text. *arXiv*
609 *preprint arXiv:2402.04379*, 2024.
- 610 Tanishq Gupta, Mohd Zaki, N. M. Anoop Krishnan, and Mausam. MatSciBERT: A materials
611 domain language model for text mining and information extraction. *npj Computational Materials*,
612 8(1):102, May 2022. ISSN 2057-3960. doi: 10.1038/s41524-022-00784-w. URL <https://www.nature.com/articles/s41524-022-00784-w>.
- 613 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance, 2022.
- 614 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
615 *neural information processing systems*, 33:6840–6851, 2020.
- 616 Jordan Hoffmann, Louis Maestrati, Yoshihide Sawada, Jian Tang, Jean Michel Sellier, and Yoshua
617 Bengio. Data-driven approach to encoding and decoding 3-d crystal structures. *arXiv preprint*
618 *arXiv:1909.00949*, 2019.
- 619 Emiel Hoogetboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows
620 and multinomial diffusion: Learning categorical distributions. *Advances in Neural Information*
621 *Processing Systems*, 34:12454–12465, 2021.
- 622 Anubhav Jain, Shyue Ping Ong, Geoffroy Hautier, Wei Chen, William Davidson Richards, Stephen
623 Dacek, Shreyas Cholia, Dan Gunter, David Skinner, Gerbrand Ceder, and Kristin A. Persson.
624 Commentary: The Materials Project: A materials genome approach to accelerating materials
625 innovation. *APL Materials*, 1(1):011002, 07 2013a. ISSN 2166-532X. doi: 10.1063/1.4812323.
626 URL <https://doi.org/10.1063/1.4812323>.
- 627 Anubhav Jain, Shyue Ping Ong, Geoffroy Hautier, Wei Chen, William Davidson Richards, Stephen
628 Dacek, Shreyas Cholia, Dan Gunter, David Skinner, Gerbrand Ceder, et al. The materials project:
629 A materials genome approach to accelerating materials innovation, *apl mater.* 2013b.
- 630 Rui Jiao, Wenbing Huang, Peijia Lin, Jiaqi Han, Pin Chen, Yutong Lu, and Yang Liu. Crystal
631 structure prediction by joint equivariant diffusion. *arXiv preprint arXiv:2309.04475*, 2023.
- 632 Rui Jiao, Wenbing Huang, Yu Liu, Deli Zhao, and Yang Liu. Space group constrained crystal
633 generation. *arXiv preprint arXiv:2402.03992*, 2024.
- 634 Sungwon Kim, Juhwan Noh, Geun Ho Gu, Alan Aspuru-Guzik, and Yousung Jung. Generative
635 adversarial networks for crystal structure prediction. *ACS central science*, 6(8):1412–1420, 2020.
- 636 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*
637 *arXiv:1312.6114*, 2013.
- 638 Walter Kohn and Lu Jeu Sham. Self-consistent equations including exchange and correlation effects.
639 *Physical review*, 140(4A):A1133, 1965.

- 648 Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi
649 Parikh, Yaniv Taigman, and Yossi Adi. Audiogen: Textually guided audio generation. *arXiv*
650 *preprint arXiv:2209.15352*, 2022.
- 651 Meng Liu, Keqiang Yan, Bora Oztekin, and Shuiwang Ji. Graphebm: Molecular graph generation
652 with energy-based models. *arXiv preprint arXiv:2102.00546*, 2021.
- 653 Teng Long, Nuno M Fortunato, Ingo Opahle, Yixuan Zhang, Ilias Samathrakakis, Chen Shen, Oliver
654 Gutfleisch, and Hongbin Zhang. Constrained crystals deep convolutional generative adversarial
655 network for the inverse design of crystal structures. *npj Computational Materials*, 7(1):66, 2021.
- 656 Steph-Yves Louis, Yong Zhao, Alireza Nasiri, Xiran Wang, Yuqi Song, Fei Liu, and Jianjun Hu. Graph
657 convolutional neural networks with global attention for improved materials property prediction.
658 *Physical Chemistry Chemical Physics*, 22(32):18141–18148, 2020.
- 659 Shitong Luo, Yufeng Su, Xingang Peng, Sheng Wang, Jian Peng, and Jianzhu Ma. Antigen-specific
660 antibody design and optimization with diffusion-based generative models for protein structures.
661 *Advances in Neural Information Processing Systems*, 35:9754–9767, 2022.
- 662 Yanchen Luo, Sihang Li, Zhiyuan Liu, Jiancan Wu, Zhengyi Yang, Xiangnan He, Xi-
663 ang Wang, and Qi Tian. Text-guided diffusion model for 3d molecule generation.
664 <https://openreview.net/pdf?id=FdUloEgBSE>, 2023a.
- 665 Youzhi Luo, Chengkai Liu, and Shuiwang Ji. Towards symmetry-aware generation of periodic
666 materials. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023b. URL
667 <https://openreview.net/forum?id=Jkc74vn1aZ>.
- 668 Benjamin Kurt Miller, Ricky TQ Chen, Anuroop Sriram, and Brandon M Wood. Flowmm: Generating
669 materials with riemannian flow matching. *arXiv preprint arXiv:2406.04713*, 2024.
- 670 Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew,
671 Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with
672 text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- 673 Juhwan Noh, Jaehoon Kim, Helge S Stein, Benjamin Sanchez-Lengeling, John M Gregoire, Alan
674 Aspuru-Guzik, and Yousung Jung. Inverse design of solid-state materials via a continuous repre-
675 sentation. *Matter*, 1(5):1370–1384, 2019.
- 676 Shyue Ping Ong, William Davidson Richards, Anubhav Jain, Geoffroy Hautier, Michael Kocher,
677 Shreyas Cholia, Dan Gunter, Vincent L Chevrier, Kristin A Persson, and Gerbrand Ceder. Python
678 materials genomics (pymatgen): A robust, open-source python library for materials analysis.
679 *Computational Materials Science*, 68:314–319, 2013.
- 680 Cheol Woo Park and Chris Wolverton. Developing an improved crystal graph convolutional neural
681 network framework for accelerated materials discovery. *Physical Review Materials*, 4(6), Jun 2020.
682 ISSN 2475-9953. doi: 10.1103/physrevmaterials.4.063801. URL [http://dx.doi.org/10.](http://dx.doi.org/10.1103/PhysRevMaterials.4.063801)
683 [1103/PhysRevMaterials.4.063801](http://dx.doi.org/10.1103/PhysRevMaterials.4.063801).
- 684 Chris J. Pickard. Airss data for carbon at 10gpa and the c+n+h+o system at 1gpa. *materi-*
685 *alscloud:2020.0026/v1*, 2020.
- 686 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
687 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.
688 Learning transferable visual models from natural language supervision, 2021.
- 689 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-
690 conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- 691 Zekun Ren, Juhwan Noh, Siyu Tian, Felipe Oviedo, Guangzong Xing, Qiaohao Liang, Armin Aberle,
692 Yi Liu, Qianxiao Li, Senthilnath Jayavelu, et al. Inverse design of crystals using generalized
693 invertible crystallographic representation. *arXiv preprint arXiv:2005.07609*, 3(6):7, 2020.

- 702 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
703 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*
704 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 705
706 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar
707 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic
708 text-to-image diffusion models with deep language understanding. *Advances in neural information*
709 *processing systems*, 35:36479–36494, 2022.
- 710 Victor Garcia Satorras, Emiel Hoogeboom, and Max Welling. E (n) equivariant graph neural networks.
711 In *International conference on machine learning*, pp. 9323–9332. PMLR, 2021.
- 712
713 Jonathan Schmidt, Love Petterson, Claudio Verdozzi, Silvana Botti, and Miguel AL Marques.
714 Crystal graph attention networks for the prediction of stable materials. *Science Advances*, 7(49):
715 eabi7948, 2021.
- 716 Chence Shi, Shitong Luo, Minkai Xu, and Jian Tang. Learning gradient fields for molecular
717 conformation generation. In *International conference on machine learning*, pp. 9558–9568. PMLR,
718 2021.
- 719
720 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised
721 learning using nonequilibrium thermodynamics. In *International conference on machine learning*,
722 pp. 2256–2265. PMLR, 2015.
- 723
724 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution.
725 *Advances in neural information processing systems*, 32, 2019.
- 726
727 Yang Song and Stefano Ermon. Improved techniques for training score-based generative models.
728 *Advances in neural information processing systems*, 33:12438–12448, 2020.
- 729
730 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
731 Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint*
arXiv:2011.13456, 2020.
- 732
733 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
734 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
systems, 30, 2017.
- 735
736 Jiaxiang Wu, Tao Shen, Haidong Lan, Yatao Bian, and Junzhou Huang. Se (3)-equivariant energy-
737 based models for end-to-end protein folding. *bioRxiv*, pp. 2021–06, 2021.
- 738
739 Tian Xie and Jeffrey C Grossman. Crystal graph convolutional neural networks for an accurate and
740 interpretable prediction of material properties. *Phys. Rev. Lett.*, 120(14):145301, 2018.
- 741
742 Tian Xie, Xiang Fu, Octavian-Eugen Ganea, Regina Barzilay, and Tommi Jaakkola. Crystal diffusion
743 variational autoencoder for periodic material generation. *arXiv preprint arXiv:2110.06197*, 2021.
- 744
745 Minkai Xu, Lantao Yu, Yang Song, Chence Shi, Stefano Ermon, and Jian Tang. Geodiff: A geometric
746 diffusion model for molecular conformation generation. *arXiv preprint arXiv:2203.02923*, 2022.
- 747
748 Keqiang Yan, Yi Liu, Yuchao Lin, and Shuiwang Ji. Periodic graph transformers for crystal material
749 property prediction. In *The 36th Annual Conference on Neural Information Processing Systems*,
750 2022.
- 751
752 D Yang, J Yu, H Wang, W Wang, C Weng, Y Zou, and D Diffsound Yu. Discrete diffusion model for
753 text-to-sound generation. arxiv 2022. *arXiv preprint arXiv:2207.09983*, 2022.
- 754
755 Mengjiao Yang, KwangHwan Cho, Amil Merchant, Pieter Abbeel, Dale Schuurmans, Igor Mor-
datch, and Ekin Dogus Cubuk. Scalable diffusion for materials generation. *arXiv preprint*
arXiv:2311.09235, 2023.
- Anthony Zee. *Group theory in a nutshell for physicists*, volume 17. Princeton University Press, 2016.

756 Claudio Zeni, Robert Pinsler, Daniel Zügner, Andrew Fowler, Matthew Horton, Xiang Fu, Sasha
757 Shysheya, Jonathan Crabbé, Lixin Sun, Jake Smith, et al. Mattergen: a generative model for
758 inorganic materials design. *arXiv preprint arXiv:2312.03687*, 2023.
759
760 Yong Zhao, Mohammed Al-Fahdi, Ming Hu, Edirisuriya MD Siriwardane, Yuqi Song, Alireza Nasiri,
761 and Jianjun Hu. High-throughput discovery of novel cubic crystal materials using deep generative
762 neural networks. *Advanced Science*, 8(20):2100566, 2021.
763
764 Huaisheng Zhu, Teng Xiao, and Vasant G Honavar. 3m-diffusion: Latent multi-modal diffusion for
765 text-guided generation of molecular graphs, 2024.
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

PERIODIC MATERIALS GENERATION USING TEXT-GUIDED JOINT DIFFUSION MODEL (TECHNICAL APPENDIX)

A LIMITATIONS AND FUTURE WORK

1) One of the major limitations and scope of the future work of our proposed work is the lack of independent textual datasets for material generation tasks. In our experimental setup, our model relied on textual data extracted from existing datasets. Initially, we extracted text data from CIF files of materials in the test sets of Perov-5, Carbon-24, and MP-20, utilizing this data to evaluate our model and all baseline models. While the experimental results show promise, a more robust evaluation could have been achieved with an independent dataset containing only textual prompts. This would enable us to assess how effectively these models can generate the underlying 3D structure of materials through a text-guided diffusion process. Hence curating an independent textual dataset for material generation containing a diverse set of meta-information will be a future scope for research.

2) Given text prompts/descriptions, we generate contextual representation using a text encoder in TGDMat, where we adopted a pre-trained MatSciBERT Gupta et al. (2022) model, which is a domain-specific language model for materials science. Also, while training TGDMat, we freeze the MatSciBERT parameters and do not tune them further. Moreover, during sampling, the user must follow a specific format (Long/Short) to provide the text description of the target material. This setup limits the expressive power of the textual representation. We investigated the robustness of TGDMat with much shorter prompts to sample from pre-trained TGDMat model in F.6, but observed performance degradation across all the benchmark dataset on *Gen* task. Hence, Exploring state-of-the-art LLMs and further fine-tuning them during training may create more powerful text conditional diffusion models and provide flexibility to process text prompts of different formats. However, that might create computational overhead as it will increase the number of parameters significantly. This provides scope for further investigation and we keep it as scope for future work.

B MORE RELATED WORK

B.1 CRYSTAL REPRESENTATION LEARNING

In recent times, graph neural network (GNN) based approaches have emerged as a powerful model in learning robust representation of crystal materials, which enhance fast and accurate property prediction. CGCNN Xie & Grossman (2018) is the first proposed model, which represents a 3D crystal structure as an undirected weighted multi-edge graph and builds a graph convolution neural network directly on the graph. Following CGCNN, there are a lot of subsequent studies Chen et al. (2019); Choudhary & DeCost (2021); Das et al. (2023a); Louis et al. (2020); Park & Wolverton (2020); Schmidt et al. (2021), where authors proposed different variants of GNN architectures for effective crystal representation learning. Recently, graph transformer-based architecture Matformer Yan et al. (2022) is proposed to learn the periodic graph representation of the material, which marginally improves the performance, however, is much faster than the prior SOTA model. Moreover, scarcity of labeled data makes these models difficult to train for all the properties, and recently, some key studies Das et al. (2022; 2023b) have shown promising results to mitigate this issue using transfer learning, pre-training, and knowledge distillation respectively.

B.2 DIFFUSION MODELS

The fundamental idea of the diffusion model, as initially proposed by Sohl-Dickstein et al. (2015), is to gradually corrupt data with diffusion noise and learn a neural model to recover back data from noise. Idea of diffusion further developed in two broad categories - 1) *Score Matching Network* Song & Ermon (2019; 2020) and 2) *Denoising Diffusion Probabilistic Models (DDPM)* Ho et al. (2020). In recent times diffusion models have emerged as a powerful new family of deep generative models, achieving remarkable performance records across numerous applications such as image synthesis Dhariwal & Nichol (2021); Ramesh et al. (2022); Rombach et al. (2022), molecular conformer generation Shi et al. (2021); Xu et al. (2022), molecular graph generation Liu et al. (2021), protein folding Luo et al. (2022); Wu et al. (2021) etc.

B.3 CONDITIONAL DIFFUSION MODELS

The initial DDPM model Ho et al. (2020) demonstrated unconditional diffusion models for image generation, where the output cannot be directed towards a desired characteristic or property. In guided diffusion models, the sampling process can be steered by a prompt, which can be a textual description of the desired output, reference image, or any other type of media.

In the field of image generation by diffusion models, Ramesh et al. (2022) came up with a text-guided diffusion model called Dall-E2 which showed how textual prompt can be used to steer the sampling process. While training the model, both the image and its textual description are encoded and mapped together, and the encoding of the prompt is used to generate the image during sampling. Another way of guiding the diffusion process using a separate classifier model was shown by Dhariwal & Nichol (2021). They trained a classifier on the noised images and used the gradient of the classifier to guide the sampling process. In the classifier-free setting, Ho & Salimans (2022) trained two diffusion models, one guided and one unguided, and combined the resulting score estimated during sampling to get the desired outcome. OpenAI’s CLIP Radford et al. (2021) further improved the relevance of the generated image to the given prompt by scoring the correctness of the generated image given the textual prompt.

Similar efforts have been made in the field of molecular generative models. The shortcoming of SMILES-based autoregressive models were addressed by TGM-DLM Gong et al. (2024) by utilizing diffusion models. This necessitates a two step process, text-guided generation phase, where the SMILES representation is generated from Gaussian noise with the help of a textual description, and correction phase, where necessary rectification are made for the correctness of SMILES string format. This is one of the drawbacks of the SMILES string format, which was addressed by 3M-Diffusion Zhu et al. (2024), where they have generated molecular graphs from a given textual description.

B.4 CRYSTAL MATERIAL GENERATION

In the past, there were limited efforts in creating novel periodic materials, with researchers concentrating on generating the atomic composition of periodic materials while largely neglecting the 3D structure. With the advancement of generative models, the majority of the research focuses on using popular generative models like VAEs or GANs to generate 3D periodic structures of materials, however, they either represent materials as three-dimensional voxel images Court et al. (2020); Hoffmann et al. (2019); Long et al. (2021); Noh et al. (2019) and generate images to depict material structures (atom types, coordinates, and lattices), or they directly encode material structures as embedding vectors Kim et al. (2020); Ren et al. (2020); Zhao et al. (2021). However, these models neither incorporate stability in the generated structure nor are invariant to any Euclidean and periodic transformations. In recent times equivariant diffusion models Xie et al. (2021); Luo et al. (2023b); Jiao et al. (2023); Yang et al. (2023); Jiao et al. (2024); Miller et al. (2024) have become the leading method for generating stable crystal materials, thanks to their capability to utilize the physical symmetries of periodic material structures. In specific, state-of-the-art models like CDVAE Xie et al. (2021) and SyMat Luo et al. (2023b) integrate a variational autoencoder (VAE) and powerful score-based decoder network, work directly with the atomic coordinates of the structures and uses an equivariant graph neural network to ensure euclidean and periodic invariance. However, both CDVAE and SyMat first predict the lattice parameters and atomic composition using the VAE model and subsequently update the coordinates using score based diffusion model. Moreover, given atomic composition, DiffCSP Jiao et al. (2023) jointly optimizes the atom coordinates and lattice using a diffusion framework to predict the crystal structure with high precision.

B.5 KEY DIFFERENCES BETWEEN DIFFCSP AND TGD MAT

Among the existing models, DiffCSP Jiao et al. (2023) comes close to our methodology, however, our work differs from it in multiple ways. DiffCSP primarily focuses only on the Crystal Structure Prediction (CSP) task and they didn’t explore the Crystal Generation task, whereas TGD Mat focuses on both tasks. Moreover, unlike DiffCSP, TGD Mat can leverage the informative textual descriptions

	DiffCSP	TGDMat
Tasks	Only CSP Task	Both CSP and Gen Tasks
Diffusion on Atom Type	-	Discrete Diffusion (D3PM)
Model Category	Unconditional; unable to specify the criteria required by the user	Conditional; able to specify the criteria required by the user (in Text Format)
Text Guided Diffusion	No	Yes

Table 6: Differences between TGDMat from DiffCSP

during the reverse diffusion process and can jointly learn lattices, atom types, and fractional coordinates from randomly sampled noise. This makes TGDMat more flexible and robust in Crystal Generation and Structure Prediction tasks.

We report key differences between DiffCSP and TGDMat in Table 6. The goal of this paper is not to introduce a new diffusion model to replace existing models like DiffCSP or CDVAE for periodic material generation. Instead, we focus on demonstrating that conditional models can outperform traditional unconditional models, such as DiffCSP. Specifically, we show that incorporating textual conditions through text-guided diffusion leads to better performance compared to using unconditional models like DiffCSP. Additionally, we enhance DiffCSP by integrating discrete diffusion over atom types in our proposed TGDMat framework.

C INVARIANCES IN CRYSTAL STRUCTURE

The basic idea of using generative models for crystal generation is to learn the underlying data distribution of material structure $p(\mathbf{M})$. Since crystal materials satisfy physical symmetry properties Dresselhaus et al. (2007); Zee (2016), one of the major challenges here is the learned distribution must satisfy periodic E(3) invariance i.e. invariance to permutation, translation, rotation, and periodic transformations.

- **Permutation Invariance** : If we permute the indices of constituent atoms it will not change the material. Formally, given any material $\mathbf{M} = (\mathbf{A}, \mathbf{X}, \mathbf{L})$, using any permutation matrix \mathbf{P} if we permute \mathbf{A} and \mathbf{X} as $\mathbf{P}(\mathbf{A})$ and $\mathbf{P}(\mathbf{X})$, then new material $\mathbf{M}_P = (\mathbf{P}(\mathbf{A}), \mathbf{P}(\mathbf{X}), \mathbf{L})$ will remain unchanged. Hence the underlying distribution is also the same i.e $p(\mathbf{M}) = p(\mathbf{M}_P)$.
- **Translation Invariance** : If we translate the atom coordinates by a random vector it will not change the structure of the material. Formally, given any material $\mathbf{M} = (\mathbf{A}, \mathbf{X}, \mathbf{L})$, if we translate \mathbf{X} by an arbitrary translation vector $\mathbf{u} \in \mathbb{R}^3$, new generated material $\mathbf{M}_T = (\mathbf{A}, \mathbf{X} + \mathbf{u}\mathbf{1}^T, \mathbf{L})$ will be the same as \mathbf{M} . Hence $p(\mathbf{M}) = p(\mathbf{M}_T)$ must satisfy.
- **Rotational Invariance** : If we rotate the atom coordinates and lattice matrix, the material remains unchanged. Formally, using any orthogonal rotational matrix $\mathbf{Q} \in R^{3 \times 3}$ (satisfying $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$), if we rotate \mathbf{X} and \mathbf{L} of any material \mathbf{M} and generate new $\mathbf{M}_R = (\mathbf{A}, \mathbf{Q}\mathbf{X}, \mathbf{Q}\mathbf{L})$, then actually different representations of the same material. Hence $p(\mathbf{M}) = p(\mathbf{M}_R)$ must satisfy.
- **Periodic Invariance** : Finally, since the atoms in the unit cell can periodically repeat itself infinite times along the lattice vector, there can be many choices of unit cells and coordinate matrices representing the same material. Formally, given coordinates \mathbf{X} , after applying periodic transformation using random matrix $\mathbf{K} \in R^{m \times 3}$, new coordinates $\mathbf{X}' = \mathbf{X} + \mathbf{K}\mathbf{L}$ are periodically equivalent. Hence $\mathbf{M} = (\mathbf{A}, \mathbf{X}, \mathbf{L})$ and $\mathbf{M}' = (\mathbf{A}, \mathbf{X}', \mathbf{L})$ are same material and $p(\mathbf{M}) = p(\mathbf{M}')$ must hold.

D TEXTUAL DATASET

Leveraging textual information to guide the reverse diffusion process remains unexplored in the material design community. To the best of our knowledge, there is currently no dataset available that includes textual descriptions of the materials present in standard benchmark databases (Section 5.1)

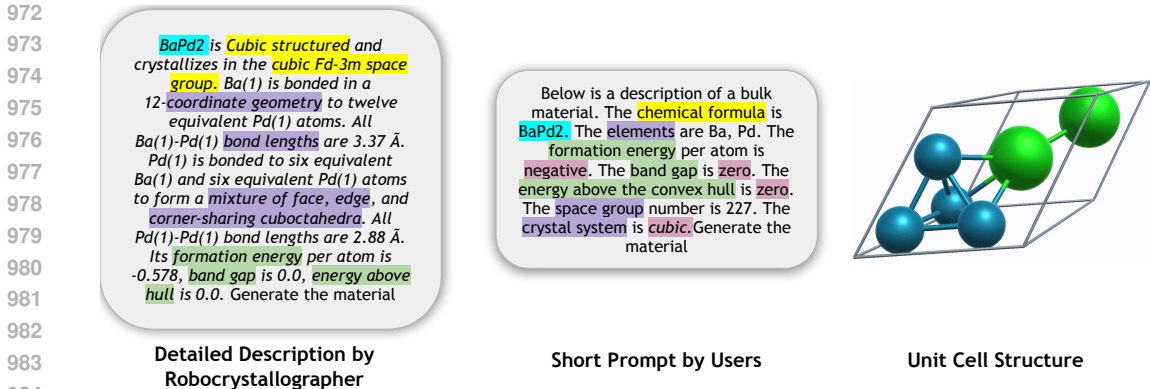


Figure 4: Detailed textual description generated by Robocrystallographer, short/less-detailed prompts by experts, and crystal unit cell structure of BaPd_2 from Material Projects dataset. Text generated by Robocrystallographer contains both local chemical compositional information related to atom/bonds (like site coordination, geometry, polyhedral connectivity, and tilt angles) and global structural knowledge (like mineral type, space group information, symmetry, and dimensionality). The shorter prompt encodes minimal information about the material like its chemical formula, constituent elements, crystal system, and few chemical properties.

used for material generation. In specific, we propose two methods for generating textual descriptions of materials. Hence, we first curate the textual dataset containing textual descriptions of these materials to train our model.

Long Detailed Textual Description: First, we utilize a freely available utility tool known as *Robocrystallographer* Ganose & Jain (2019) to generate detailed textual descriptions about the periodic structure of crystal materials encoded in Crystallographic Information Files (CIF Files). Robocrystallographer breaks down crystal structures into two main components: local compositional details such as atomic coordination, geometry, polyhedral connectivity, and tilt angles, as well as global structural aspects like crystal formula, mineral type, space group information, symmetry, and dimensionality. This information is presented in three formats: JSON for machine processing, human-readable text for easy comprehension akin to descriptions provided by humans, and machine learning format for specialized analysis. We choose the human-readable text format to compile textual datasets, which closely resemble descriptions given of the crystal structure by humans.

Short Custom Prompts: Secondly, we utilized shorter and less detailed prompts that are more easily interpretable by users. We extend the prompt template proposed by Gruver et al. (2024), which encodes minimal information about the material like its chemical formula, constituent elements, crystal system it belongs to, and its space group number. Further, we specify a few chemical properties, and instead of mentioning their actual values, we provide generic information like negative/positive formation energy, zero/nonzero band gaps, etc. We used the Pymatgen tool Ong et al. (2013) to extract this information from the Crystallographic Information Files (CIF Files) and curate the textual prompts.

An illustrative example of both these textual descriptions and the unit cell structure is provided in Figure 4.

E JOINT EQUIVARIANT DIFFUSION ON \mathbf{M}

Given an input crystal material $\mathbf{M}_0 = (\mathbf{A}_0, \mathbf{X}_0, \mathbf{L}_0)$, we define a forward diffusion process through a Markov chain over T steps to defuse \mathbf{A} , \mathbf{X} , \mathbf{L} independently as follows :

$$q(\mathbf{A}_t, \mathbf{X}_t, \mathbf{L}_t | \mathbf{A}_{t-1}, \mathbf{X}_{t-1}, \mathbf{L}_{t-1}) = q(\mathbf{A}_t | \mathbf{A}_{t-1})q(\mathbf{X}_t | \mathbf{X}_{t-1})q(\mathbf{L}_t | \mathbf{L}_{t-1}) \quad t = 1, 2, \dots, T \quad (7)$$

E.1 DIFFUSION ON LATTICE (\mathbf{L})

Lattice Matrix $\mathbf{L} = [l_1, l_2, l_3]^T \in \mathbb{R}^{3 \times 3}$ is a global feature of the material which determines the shape and symmetry of the unit cell structure. Since \mathbf{L} is in continuous space, we leverage the idea of the Denoising Diffusion Probabilistic Model (DDPM) for diffusion on \mathbf{L} . In specific, given input lattice matrix $\mathbf{L}_0 \sim p(\mathbf{L})$, the forward diffusion process iteratively diffuses it over T timesteps to a noisy lattice matrix \mathbf{L}_T through a transition probability $q(\mathbf{L}_t | \mathbf{L}_0)$ at each t^{th} step, which can be derived as follows :

$$q(\mathbf{L}_t | \mathbf{L}_0) = \mathcal{N}\left(\mathbf{L}_t | \sqrt{\bar{\alpha}_t} \mathbf{L}_0, (1 - \bar{\alpha}_t) \mathbf{I}\right) \quad (8)$$

where, $\bar{\alpha}_t = \prod_{k=1}^t \alpha_k$, $\alpha_t = 1 - \beta_t$ and $\{\beta_t \in (0, 1)\}_{t=1}^T$ controls the variance of diffusion step following certain variance scheduler. By reparameterization, we can rewrite equation 8 as:

$$\mathbf{L}_t = \sqrt{\bar{\alpha}_t} \mathbf{L}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}^L \quad (9)$$

where, $\boldsymbol{\epsilon}^L$ is a noise, sampled from $\mathcal{N}(\mathbf{0}, \mathbf{I})$, added with original input sample \mathbf{L}_0 at t^{th} step to generate \mathbf{L}_t . After T such diffusion steps, noisy lattice matrix \mathbf{L}_T is generated from prior noise distribution $\sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. In the reverse denoising process, given noisy $\mathbf{L}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ we reconstruct true lattice structure \mathbf{L}_0 through iterative denoising step via learning reverse conditional distribution, which we formulate as follows :

$$p(\mathbf{L}_{t-1} | \mathbf{M}_t, \mathbf{C}_p) = \mathcal{N}\left\{\mathbf{L}_{t-1} | \mu^L(\mathbf{M}_t, \mathbf{C}_p), \beta_t \frac{(1 - \bar{\alpha}_{t-1})}{(1 - \bar{\alpha}_t)} \mathbf{I}\right\} \quad (10)$$

where $\mu^L(\mathbf{M}_t, \mathbf{C}_p) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{L}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \hat{\boldsymbol{\epsilon}}^L(\mathbf{M}_t, \mathbf{C}_p, t)\right)$. Intuitively, $\hat{\boldsymbol{\epsilon}}^L$ is the denoising term that needs to be subtracted from \mathbf{L}_t to generate \mathbf{L}_{t-1} and textual representation \mathbf{C}_p will steer this reverse diffusion process. We use a text-guided denoising network $\Phi_\theta(\mathbf{A}_t, \mathbf{X}_t, \mathbf{L}_t, t, \mathbf{C}_p)$ to model the noise term $\hat{\boldsymbol{\epsilon}}^L(\mathbf{M}_t, \mathbf{C}_p, t)$. Following the simplified training objective proposed by Ho et al. (2020), we train the aforementioned denoising network using l_2 loss between $\hat{\boldsymbol{\epsilon}}^L$ and $\boldsymbol{\epsilon}^L$

$$\mathcal{L}_{lattice} = \mathbb{E}_{\boldsymbol{\epsilon}^L, t \sim \mathcal{U}(1, T)} \|\boldsymbol{\epsilon}^L - \hat{\boldsymbol{\epsilon}}^L\|_2^2 \quad (11)$$

E.2 DIFFUSION ON ATOM TYPES (\mathbf{A})

Prior studies Jiao et al. (2023); Xie et al. (2021) consider Atom Type Matrix \mathbf{A} as the logits/probability distribution for k classes $\in \mathbb{R}^{N \times k}$ (continuous variable in real space) and apply DDPM to learn the distribution. However for discrete data these models are inappropriate and produce suboptimal results Austin et al. (2021); Campbell et al. (2022); Hoogeboom et al. (2021). Hence we consider \mathbf{A} as N discrete variables belonging to k classes and leverage discrete denoising diffusion probabilistic model (D3PM) Austin et al. (2021) for diffusion on \mathbf{A} . In specific, denoting row vector \mathbf{a} as a one-hot representation of an atom a , we can write transition probability for forward process as:

$$q(\mathbf{a}_t | \mathbf{a}_{t-1}) = \text{Cat}(\mathbf{a}_t; \mathbf{p} = \mathbf{a}_{t-1} \mathbf{Q}_t) \quad (12)$$

where $\text{Cat}(\mathbf{a}; \mathbf{p})$ is a categorical distribution over the one-hot row vector \mathbf{a} with probabilities given by the row vector \mathbf{p} and \mathbf{Q}_t is the Markov transition matrix at time step t defined as $[\mathbf{Q}_t]_{i,j} = q(a_t = i | a_{t-1} = j)$. Different choices of \mathbf{Q}_t and corresponding stationary distributions are proposed by Austin et al. (2021) which provides flexibility to control the data corruption and denoising process. We adopted the absorbing state diffusion process, introducing a new absorbing state [MASK] in \mathbf{Q}_t . At each time step t , we can formally define the transition matrix as:

$$[\mathbf{Q}_t]_{i,j} = \begin{cases} 1, & \text{if } i = j = [\text{MASK}]. \\ 1 - \beta_t, & \text{if } i = j \neq [\text{MASK}] \\ \beta_t, & \text{if } i = j = [\text{MASK}]. \end{cases} \quad (13)$$

Intuitively, at each time step t , an atom either stays in its type state with probability $1 - \beta_t$ or moves to [MASK] state with probability β_t and once it moves to [MASK] state, it stays in that state. Hence, the stationary distribution of this diffusion process has all the mass on the [MASK] state. During reverse denoising process, given textual representation \mathbf{C}_p , we first sample noisy \mathbf{a}_T and obtain \mathbf{a}_0 through iterative denoising step via learning reverse conditional transition:

$$p_\theta(\mathbf{a}_{t-1} | \mathbf{a}_t, \mathbf{C}_p) \propto \sum_{\mathbf{a}_0} q(\mathbf{a}_{t-1}, \mathbf{a}_t | \mathbf{a}_0) p_\theta(\mathbf{a}_0 | \mathbf{a}_t, \mathbf{C}_p) \quad (14)$$

We use the text-guided denoising network $\Phi_\theta(\mathbf{A}_t, \mathbf{X}_t, \mathbf{L}_t, t, \mathbf{C}_p)$ to model this backward denoising process, which is trained using the following loss function as proposed by Austin et al. (2021) :

$$\mathcal{L}_{type} = \mathcal{L}_{VB} + \lambda \mathcal{L}_{CE} \quad (15)$$

where \mathcal{L}_{VB} is the variational lower bound loss defined as follows:

$$\mathcal{L}_{VB} = \mathbb{E}_{q(\mathbf{a}_0)} \left[\underbrace{D_{KL}\{q(\mathbf{a}_T|\mathbf{a}_0)||p(\mathbf{a}_T)\}}_{L_T} + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{a}_t|\mathbf{a}_0)} \left[\underbrace{D_{KL}\{q(\mathbf{a}_{t-1}|\mathbf{a}_t, \mathbf{a}_0)||p_\theta(\mathbf{a}_{t-1}|\mathbf{a}_t)\}}_{L_{t-1}} - \underbrace{\mathbb{E}_{q(\mathbf{a}_1|\mathbf{a}_0)}[\log p_\theta(\mathbf{a}_0|\mathbf{a}_1)]}_{L_0} \right] \right] \quad (16)$$

and \mathcal{L}_{CE} is the cross-entropy loss defined as follows:

$$\mathcal{L}_{CE} = \mathbb{E}_{q(\mathbf{a}_0)} \left[\sum_{t=2}^T \mathbb{E}_{q(\mathbf{a}_t|\mathbf{a}_0)} [\log p_\theta(\mathbf{a}_0|\mathbf{a}_t)] \right] \quad (17)$$

and λ is a hyperparameter.

E.3 DIFFUSION ON ATOM COORDINATES (\mathbf{X})

Coordinate Matrix $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T \in \mathbb{R}^{N \times 3}$ denotes atomic coordinate positions, where $x_i \in \mathbb{R}^3$ corresponds to coordinates of i^{th} atom in the unit cell. We can diffuse the atom coordinates in two ways: either by diffusing cartesian coordinates or fractional coordinates. Prior works like CDVAE Xie et al. (2021) and SyMat Luo et al. (2023b) diffuse cartesian coordinates whereas DiffCSP Jiao et al. (2023) diffuse fractional coordinates. In our setup, as we are jointly learning atom coordinates and lattice matrix simultaneously, we follow the line of work by DiffCSP and diffuse fractional coordinates. Atomic fractional coordinates in crystal material lives in quotient space $\mathbb{R}^{N \times 3} / \mathbb{Z}^{N \times 3}$ induced by the crystal periodicity. Since the Gaussian distribution used in DDPM is unable to model the cyclical and bounded domain of \mathbf{X} , it is not suitable to apply DDPM to model \mathbf{X} . Hence at each step of forward diffusion, we add noise sample from Wrapped Normal (WN) distribution De Bortoli et al. (2022) to \mathbf{X} and during backward diffusion leverage Score Matching Diffusion Networks Song & Ermon (2019; 2020) to model underlying transition probability $q(\mathbf{X}_t | \mathbf{X}_0) = \mathcal{N}_W(\mathbf{X}_t | \mathbf{X}_0, \sigma_t^2 \mathbf{I})$. In specific, at each t^{th} step of diffusion, we derive \mathbf{X}_t as : $\mathbf{X}_t = f_w(\mathbf{X}_0 + \sigma_t \epsilon^{\mathbf{X}})$ where, $\epsilon^{\mathbf{X}}$ is a noise, sampled from $\mathcal{N}(\mathbf{0}, \mathbf{I})$, σ_t is the noise scale following exponential scheduler and $f_w(\cdot)$ is a truncation function. Given a fractional coordinate matrix \mathbf{X} , truncation function $f_w(\mathbf{X}) = (\mathbf{X} - \lfloor \mathbf{X} \rfloor)$ returns the fractional part of each element of \mathbf{X} .

As argued in Jiao et al. (2023), $q(\mathbf{X}_t|\mathbf{X}_0)$ is periodic translation equivariant, and approaches uniform distribution $\mathcal{U}(0, 1)$ for sufficiently large values of σ_T . Hence during the backward denoising process, we first sample $\mathbf{X}_T \sim \mathcal{U}(0, 1)$ and iteratively denoise via score network for T steps to recover back the true fractional coordinates \mathbf{X}_0 . We use the text-guided denoising network $\Phi_\theta(\mathbf{A}_t, \mathbf{X}_t, \mathbf{L}_t, t, \mathbf{C}_p)$ to model the backward diffusion process, which is trained using the following score-matching objective function :

$$\mathcal{L}_{coord} = \mathbb{E}_{\substack{\mathbf{X}_t \sim q(\mathbf{X}_t|\mathbf{X}_0) \\ t \sim \mathcal{U}(1, T)}} \|\nabla_{\mathbf{X}_t} \log q(\mathbf{X}_t|\mathbf{X}_0) - \hat{\epsilon}^{\mathbf{X}}(\mathbf{M}_t, \mathbf{C}_p, t)\|_2^2 \quad (18)$$

where $\nabla_{\mathbf{X}_t} \log q(\mathbf{X}_t|\mathbf{X}_0) \propto \sum_{\mathbf{K} \in \mathbb{Z}^{N \times 3}} \exp(-\frac{\|\mathbf{X}_t - \mathbf{X}_0 + \mathbf{K}\|_F^2}{2\sigma_t^2})$ is the score function of transitional distribution and $\hat{\epsilon}^{\mathbf{X}}(\mathbf{M}_t, \mathbf{C}_p, t)$ denoising term.

E.4 TEXT GUIDED DENOISING NETWORK

In this subsection, we will illustrate the detailed architecture of our proposed Text Guided Denoising Network $\Phi_\theta(\mathbf{A}_t, \mathbf{X}_t, \mathbf{L}_t, t, \mathbf{C}_p)$, which we used to denoise \mathbf{A} , \mathbf{X} and \mathbf{L} . As mentioned in 2.2, the learned distribution of material structure $p(\mathbf{M})$ must satisfy periodic E(3) invariance. Hence we leverage an periodic-E(3)-equivariant Graph Neural Network (GNN) integrated with a pre-trained textual encoder to model the denoising process. In particular, as a text encoder, we adopt a pre-trained MatSciBERT Gupta et al. (2022) model, which is a domain-specific language model for materials science, followed by a projection layer. MatSciBERT is effectively a pre-trained SciBERT model on a scientific text corpus of 3.17B words, which is further trained on a huge text corpus of materials science containing around 285 M words. We feed textual description of material \mathcal{T} and extract embedding of [CLS] token \mathbf{h}_{CLS} as a representation of the whole text. Further, we pass \mathbf{h}_{CLS} through a projection layer to generate the contextual textual embedding for the material $\mathbf{C}_p \in \mathbb{R}^d$, which we pass to the equivariant GNN model to guide the denoising process. Practically, as the backbone network

for the backward diffusion process, we extend CSPNet architecture Jiao et al. (2023), originally developed for crystal structure prediction (CSP) task. CSPNet is built upon EGNN Satorras et al. (2021), satisfying periodic E(3) invariance condition on periodic crystal structure. At the k^{th} layer message passing, the Equivariant Graph Convolutional Layer (EGCL) takes as input the set of atom embeddings $\mathbf{h}^k = [\mathbf{h}_1^k, \mathbf{h}_2^k, \dots, \mathbf{h}_N^k]$, atom coordinates $\mathbf{x}^k = [\mathbf{x}_1^k, \mathbf{x}_2^k, \dots, \mathbf{x}_N^k]$ and Lattice Matrix \mathbf{L} and outputs a transformation on \mathbf{h}^{k+1} . Formally, we can define the k^{th} layer message passing operation as follows :

$$\mathbf{m}_{i,j} = \rho_m \{ \mathbf{h}_i^k, \mathbf{h}_j^k, \mathbf{L}^T \mathbf{L}, \psi_{FT}(\mathbf{x}_i^k - \mathbf{x}_j^k) \}; \quad (19)$$

$$\mathbf{h}_i^{k+1} = \mathbf{h}_i^k + \rho_h \{ \mathbf{h}_i^k, \mathbf{m}_i \} \quad (20)$$

where $\mathbf{m}_i = \sum_{j=1}^N \mathbf{m}_{i,j}$, ρ_m, ρ_h are multi-layer perceptrons and ψ_{FT} is a Fourier Transformation function applied on relative difference between fractional coordinates $\mathbf{x}_i^k, \mathbf{x}_j^k$. Fourier Transformation is used since it is invariant to periodic translation and extracts various frequencies of all relative fractional distances that are helpful for crystal structure modeling.

We fuse textual representation \mathbf{C}_p into input atom feature \mathbf{h}_i^0 as

$$\mathbf{h}_i^0 = \rho \{ f_{atom}(\mathbf{a}_i) \parallel f_{pos}(t) \parallel \mathbf{C}_p \} \quad (21)$$

where t is the timestamp of the diffusion model, $f_{pos}(\cdot)$ is sinusoidal positional encoding Ho et al. (2020); Vaswani et al. (2017), $f_{atom}(\cdot)$ learned atomic embedding function and \parallel is concatenation operation. Input atom features \mathbf{h}^0 and coordinates \mathbf{x}^0 are fed through \mathcal{K} layers of EGCL to produce $\hat{\epsilon}^L, p(\mathbf{A}_{t-1} | \mathbf{M}_t)$ and $\hat{\epsilon}^X$ as follows :

$$\hat{\epsilon}^L = \mathbf{L} \rho_L \left(\frac{1}{N} \sum_N^{i=1} \mathbf{h}^{\mathcal{K}} \right); \quad (22)$$

$$p(\mathbf{A}_{t-1} | \mathbf{M}_t) = \rho_A(\mathbf{h}^{\mathcal{K}});$$

$$\hat{\epsilon}^X = \rho_X(\mathbf{h}^{\mathcal{K}})$$

where ρ_L, ρ_A, ρ_X are multi-layer perceptrons on the final layer embeddings. Intuitively, we feed global structural knowledge about the crystal structure into the network by injecting contextual representation \mathbf{C}_p into input atom features. This added signal will participate through message-passing operations in Eq. 19 and guides in denoising atom types, coordinates, and lattice parameters such that it can capture the global crystal geometry and aligned with the input stable structure specified by textual description.

Algorithm 1 Training Algorithm

- 1: **Input:** Atom type Matrix \mathbf{A}_0 (One hot Vector Representation), Coordinate Matrix \mathbf{X}_0 , Lattice matrix \mathbf{L}_0 , Markov Transition Matrix $[\mathbf{Q}_t]_{t=1}^T$, Textual Representation \mathbf{C}_p , Number of diffusion step T and hyperparameters $\lambda_A, \lambda_X, \lambda_L$.
 - 2: **repeat**
 - 3: Sample $t \sim \mathcal{U}(0, T)$
 - 4: Sample Noise $\epsilon^X, \epsilon^L \sim N(0, I)$
 - 5: $\mathbf{L}_t = \sqrt{\bar{\alpha}_t} \mathbf{L}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon^L$
 - 6: $\mathbf{X}_t = f_w(\mathbf{X}_0 + \sigma_t \epsilon^x)$
 - 7: $\mathbf{A}_t = \text{Cat}(\mathbf{A}_t; \mathbf{p} = \mathbf{A}_{t-1} \mathbf{Q}_t)$
 - 8: $\hat{\epsilon}^L, \hat{\epsilon}^X, \mathbf{A}'_t \leftarrow \Phi_\theta(\mathbf{A}_t, \mathbf{X}_t, \mathbf{L}_t, t, \mathbf{C}_p)$
 - 9: $\mathcal{L}_{lattice} = \|\epsilon^L - \hat{\epsilon}^L\|_2^2$
 - 10: $\mathcal{L}_{coord} = \|\nabla_{\mathbf{X}_t} \log q(\mathbf{X}_t | \mathbf{X}_0) - \hat{\epsilon}^X\|_2^2$
 - 11: $\mathcal{L}_{type} = \mathcal{L}_{VB} + \lambda \mathcal{L}_{CE}$
 - 12: Minimize $\mathcal{L} = \lambda_L \mathcal{L}_{lattice} + \lambda_A \mathcal{L}_{type} + \lambda_X \mathcal{L}_{coord}$ and update parameters of Φ_θ
 - 13: **until** converged
-

E.5 TRAINING AND SAMPLING

TGDMat is trained using the following combined loss:

$$\mathcal{L} = \lambda_L \mathcal{L}_{lattice} + \lambda_A \mathcal{L}_{type} + \lambda_X \mathcal{L}_{coord} \quad (23)$$

Algorithm 2 Sampling Algorithm

```

1188 1: Sample  $\mathbf{L}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \mathbf{X}_T \sim \mathcal{U}(0, 1)$ 
1189 2: Randomly sample each atom type between 0 to 99 (Max possible atom type) and form  $\mathbf{A}_T$ 
1190 3:  $\mathbf{C}_p \leftarrow$  Textual Representation
1191 4: for  $t \leftarrow T$  to 1 do
1192 5:    $\epsilon^{\mathbf{A}}, \epsilon^{\mathbf{X}}, \epsilon^{\mathbf{L}} \sim N(0, I) / * Sample * /$ 
1193 6:    $\hat{\mathbf{A}}, \hat{\epsilon}^{\mathbf{X}}, \hat{\epsilon}^{\mathbf{L}} \leftarrow \Phi_{\theta}(\mathbf{A}_t, \mathbf{X}_t, \mathbf{L}_t, t, \mathbf{C}_p)$ 
1194 7:    $\mathbf{L}_{t-1} \leftarrow \frac{1}{\sqrt{\alpha_t}}(\mathbf{L}_t - \frac{\beta_t}{\sqrt{1-\alpha_t}}\hat{\epsilon}^{\mathbf{L}}) + \sqrt{\beta_t \frac{1-\bar{\alpha}_{t-1}}{1-\alpha_t}}\epsilon^{\mathbf{L}}$ 
1195 8:    $\mathbf{A}_{t-1} \leftarrow \text{Softmax}(\hat{\mathbf{A}} + \sigma_t \epsilon^{\mathbf{A}})$ 
1196 9:    $\mathbf{X}_{t-\frac{1}{2}} \leftarrow w(\mathbf{X}_t + (\sigma_t^2 - \sigma_{t-1}^2)\hat{\epsilon}^{\mathbf{X}} + \frac{\sigma_{t-1}\sqrt{\sigma_t^2 - \sigma_{t-1}^2}}{\sigma_t}\epsilon^{\mathbf{X}})$ 
1197 10:    $\hat{\epsilon}^{\mathbf{X}} \leftarrow \Phi_{\theta}(\mathbf{A}_t, \mathbf{X}_{t-\frac{1}{2}}, \mathbf{L}_{t-1}, t, \mathbf{C}_p)$ 
1198 11:    $\eta_t \leftarrow step\_size * \frac{\sigma_{t-1}}{\sigma_t}$ 
1199 12:    $\mathbf{X}_{t-1} \leftarrow w(\mathbf{X}_{t-\frac{1}{2}} + \eta_t \hat{\epsilon}^{\mathbf{X}} + \sqrt{2\eta_t}\epsilon^{\mathbf{X}})$ 
1200 13: end for

```

Model	# parameters	Model size
CDVAE	4,920,414	18.771 MB
SyMat	3,385,601	12.915 MB
DiffCSP	12,294,656	46.923 MB
TGDMat	12,432,228	47.448 MB

Table 7: Model size comparison of Baselines and TGDMat

where $\mathcal{L}_{lattice}$, \mathcal{L}_{type} and \mathcal{L}_{coord} are lattice l_2 loss (Eq. 11), type cross-entropy loss (Eq. 15) and coordinate score matching loss (Eq. 18) respectively and λ_L , λ_A , λ_X are hyperparameters control the relative weightage between these different loss components. During training, we freeze the MatSciBERT parameters and do not tune it further. During sampling, we use the Predictor-Corrector sampling mechanism to sample \mathbf{A}_0 , \mathbf{X}_0 and \mathbf{L}_0 . Next we explain algorithms for training and sampling.

F EXPERIMENTS

F.1 EXPERIMENTAL SETUP

Benchmark Tasks. We evaluate our proposed model TGDMat on two different categories of tasks for material generation, *Random Material Generation (Gen)* and *Crystal Structure Prediction (CSP)*. In *Gen* task, the goal of the generative model is to generate novel stable materials (atom types, fractional coordinates, and lattice structure). In *CSP* task, atom types of the materials are given and the goal is to predict/match the crystal structure (atom coordinates and lattice). In TGDMat model, by design choice, we use the textual description of crystal materials during each step of the reverse diffusion process to enhance the generation capability in both tasks. A pictorial illustration of both tasks is provided at 5

Dataset. Following Xie et al. (2021) we evaluate our model on three baseline datasets: **Perov-5**, **Carbon-24** and **MP-20**. **Perov-5** Castelli et al. (2012a;b) dataset consists of 18,928 perovskite materials, each with 5 atoms in a cell. They generally can be denoted by \mathbf{ABX}_3 indicating the three different types of atoms usually observed in such materials. **Carbon-24** Pickard. (2020) dataset has 10,153 materials with 6 to 24 atoms of carbon in the crystal lattice. Finally, **MP-20** Jain et al. (2013b) dataset has 45,231 materials curated from the Materials Project library Jain et al. (2013a), where each material has at most 20 atoms in the lattice. Crystals from **Perov-5** dataset share the same structure but differ in composition, whereas Crystals from **Carbon-24** share the same composition but differ in structure. Crystals from **MP-20** differs in both structure and composition. We curated textual data for these datasets with a textual description of each material. Specifically, we

1242
 1243
 1244
 1245
 1246
 1247
 1248
 1249
 1250
 1251
 1252
 1253
 1254
 1255
 1256
 1257
 1258
 1259
 1260
 1261
 1262
 1263
 1264
 1265
 1266
 1267
 1268
 1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282
 1283
 1284
 1285
 1286
 1287
 1288
 1289
 1290
 1291
 1292
 1293
 1294
 1295

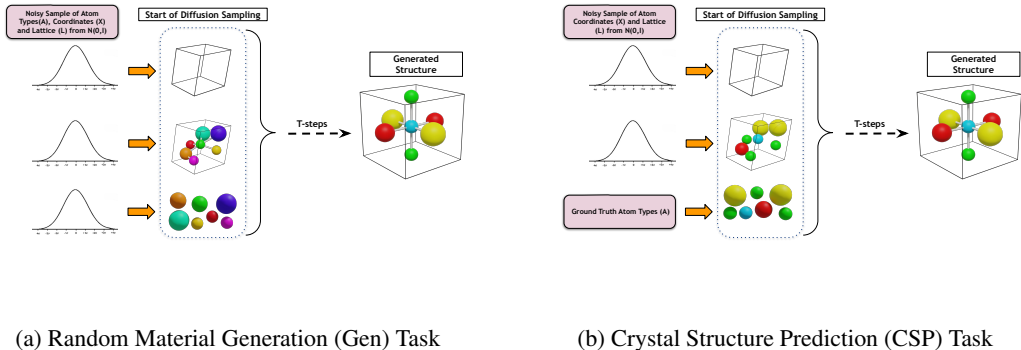


Figure 5

generate both long detailed textual descriptions and shorter prompts using approaches mentioned in Appendix D.

The structures in all three datasets are derived from quantum mechanical simulations and are all at local energy minima. Most materials in **Perov-5** and **Carbon-24** are hypothetical, whereas **MP-20** represents a realistic dataset that includes many experimentally known inorganic materials, each with a maximum of 20 atoms in the unit cell, most of which are globally stable. A model that performs well on MP-20 could potentially generate novel materials that can be synthesized experimentally. While training TGDMat, we split the datasets into the train, test, and validation sets following the convention of 60:20:20 as done by Xie et al. (2021).

Hyper-Parameters Details. In our TGDMat model, we adopted 4 layers CSPNet as message passing layer with hidden dimension set as 512. Further, we use pre-trained MatSciBERT Gupta et al. (2022) followed by a two-layer projection layer (projection dimension 64) as the text encoder module. We keep the dimension of time embedding at each diffusion timestep as 64. We train it for 500 epochs using the same optimizer, and learning rate scheduler as DiffCSP and keep the batch size as 512. We perform all the experiments in the Tesla P100-PCIE-16GB GPU server.

F.2 EVALUATION METRICS

Random Material Generation (Gen) Task. Following CDVAE Xie et al. (2021), we evaluate the performance of TGDMat and baseline models on generating novel material structure using seven metrics under three broad categories: **Validity**, **Coverage**, and **Property Statistics**. Under **Validity**, following the prior line of work Court et al. (2020); Xie et al. (2021), we measure structural and compositional validity, representing the percentages of generated crystals with valid periodic structures and atom types, respectively. A structure is valid as long as the shortest distance between any pair of atoms is larger than 0.5 \AA whereas the composition is valid if the overall charge is neutral as computed by SMOG Davies et al. (2019). In **Coverage**, we consider two coverage metrics, COV-R (Recall) and COV-P (Precision). COV-R measures the percentage of the test set materials being correctly predicted, whereas COV-P measures the percentage of generated materials that cover at least one of the test set materials. (More detailed discussions can be found in Xie et al. (2021) and Ganea et al. (2021)). Finally, we evaluate the similarity between the generated materials and those in the test set using various **Property Statistics**, where we compute the earth mover’s distance (EMD) between the distributions in element number (# Elem), density (ρ , unit g/cm^3), and formation energy (\mathcal{E} , unit eV/atom) predicted by a GNN model.

Crystal Structure Prediction (CSP) Task. We evaluate the performance of TGDMat and baseline models on stable structure prediction using standard metrics proposed by the prior works Jiao et al. (2023); Xie et al. (2021), by matching the generated structure and the input ground truth structure in the test set. In Specific, for each material structure in the test set, we generate k samples given the textual description and then identify the matching if at least one of the samples matches the ground truth structure. We calculate the **Match Rate** and **RMSE** metrics using the StructureMatcher

class in Pymatgen, which identifies the best match between two structures while accounting for all material invariances. Match rate indicates the percentage of the matched structures over the test set satisfying thresholds $\text{stol}=0.5$, $\text{angle_tol}=10$, $\text{ltol}=0.3$. RMSE is computed 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 baselines and TGDMat, we evaluate using $k = 1$ and $k = 20$.

F.3 COMPLETE AND DETAILED RESULTS

In this subsection, we provide full comprehensive results on both Gen and CSP tasks across three benchmark datasets and evaluate the performance of all the baseline models, their text-guided variants (both short and long), and our proposed TGDMat(Long) & TGDMat(Short). We report the CSP and Gen task results in Table 8 and 9 respectively.

Following are the Insights or Observations:

- For both tasks, across all the datasets, text guidance outperforms the vanilla diffusion models in almost all metrics.
- Our experiments suggest that using shorter prompts text-guided models outperforms the vanilla baseline models. However, performance is even superior when using text-guided diffusion using longer prompts.
- For the CSP task, using text guidance during the reverse denoising process, with just one generated sample per test material, text-guided variants outperform respective vanilla models, thereby reducing computational overhead.
- Our proposed TGDMat (Long) stands out as the leading model when compared to all baseline models and their text-guided variants across three benchmark datasets. In specific, for Gen Task, TGDMat (Long) outperforms the closest baseline DiffCSP+ (Long) because we leveraged discrete diffusion on atom types, which is more powerful in learning discrete variables like atom types.
- Finally, results indicate that utilizing shorter prompts TGDMat (Short) results in a slight decrease in overall performance compared to the longer variant TGDMat (Long). Nonetheless, the performance remains superior or comparable to baseline models (vanilla and text-guided variants).

F.4 CORRECTNESS OF GENERATED MATERIALS

Setup. In this section, we investigate whether the generated material matches different features specified by the textual prompts. TGDMat has the capability to process textual prompts given by the user, enabling it to manage global attributes about crystal materials such as Formula, Space group, Crystal System, and different property values like formation energy, band-gap, etc. To ensure the fidelity of our model’s outputs concerning these specified global attributes from the text prompt, We randomly generated 1000 materials (sampled from all three Datasets) based on their respective textual descriptions(both Long and Short) and assessed the percentage of generated materials that matched the global features outlined in the text prompt. In specific, we matched the Formula, Space group, and Crystal System, and Dimensions of generated materials with the textual descriptions. Moreover, we examined whether properties such as formation energy and bandgap matched the specified criteria as per the text prompt (positive/negative, zero/nonzero).

Results and Discussions. We report the results in Table 10. In general, using longer text, considering Perov-5 and Carbon-24 datasets, the generated material meets the specified criteria effectively. However, when dealing with the MP-20 dataset, which is more intricate due to its complex structure and composition, performance tends to decline. Additionally, when using shorter prompts, overall performance suffers across all datasets compared to longer text inputs. This is because the longer text, provided by the robocrystallographer, offers a comprehensive range of information, both global and local, thereby enhancing the generation capabilities of TGDMat.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

Method	# samples	Perov-5		Carbon-24		MP-20	
		Match	RMSE	Match	RMSE	Match	RMSE
CDVAE	1	45.31	0.1138	17.09	0.2969	33.9	0.1045
	20	88.51	0.0464	88.37	0.2286	66.95	0.1026
CDVAE+(short)	1	48.97	0.1063	22.65	0.264	40.33	0.1037
	20	89.54	0.0423	89.61	0.2188	70.22	0.0876
CDVAE+(long)	1	49.25	0.1055	23.73	0.259	41.8	0.1021
	20	89.73	0.0417	89.77	0.2053	72.56	0.084
SyMat	1	47.32	0.1074	20.81	0.2655	33.92	0.1039
	20	90.25	0.0316	89.29	0.2184	71.03	0.0945
SyMat+(short)	1	49.39	0.0985	23.71	0.2567	40.84	0.1027
	20	92.1	0.0255	90.86	0.2069	71.31	0.0875
SyMat+(long)	1	50.88	0.0963	28.18	0.251	43.17	0.1016
	20	92.3	0.0201	91.65	0.187	72.96	0.082
DiffCSP	1	52.02	0.076	17.54	0.2759	51.49	0.0631
	20	98.6	0.0128	88.47	0.2192	77.93	0.0492
DiffCSP+(short)	1	56.54	0.0583	24.13	0.2424	52.22	0.0597
	20	98.25	0.0137	88.28	0.2252	80.97	0.0443
DiffCSP+(long)	1	90.46	0.0203	44.63	0.2266	55.15	0.0572
	20	98.59	0.0072	95.27	0.1534	82.02	0.0391
TGDMat (short)	1	56.54	0.0583	24.13	0.2424	52.22	0.0597
	20	98.25	0.0137	88.28	0.2252	80.97	0.0443
TGDMat (long)	1	90.46	0.0203	44.63	0.2266	55.15	0.0572
	20	98.59	0.0072	95.27	0.1534	82.02	0.0391

Table 8: Summary of the Complete and Detailed Results on the CSP Task.

Dataset	Method	Validity		Coverage		Property		
		Comp	Struct	Cov-R	Cov-P	# element	Density	form_energy
Perov 5	CDVAE	98.29	100	99.25	98.39	0.0731	0.1462	0.0291
	CDVAE+ (Short)	98.17	100	99.4	99.01	0.0706	0.1395	0.0246
	CDVAE+ (Long)	98.45	100	99.53	99.09	0.0609	0.1276	0.0223
	SyMat	96.83	100	99.16	98.29	0.0193	0.1991	0.2827
	SyMat+ (Short)	96.94	100	99.22	98.4	0.0192	0.1827	0.2633
	SyMat+ (Long)	97.88	100	99.71	98.79	0.0172	0.1755	0.2566
	DiffCSP	98.15	100	99.28	98.08	0.0132	0.1281	0.0267
	DiffCSP+ (Short)	98.21	100	99.61	98.39	0.0123	0.1193	0.0266
	DiffCSP+ (Long)	98.44	100	99.85	98.53	0.0119	0.1071	0.0241
	TGDMat(Short)	98.28	100	99.71	99.24	0.0108	0.0947	0.0237
	TGDMat(Long)	98.63	100	99.87	99.52	0.009	0.0497	0.0187
	Carbon 24	CDVAE	-	100	99.35	82.66	-	0.1539
CDVAE+ (Short)		-	100	99.34	82.96	-	0.1398	0.2804
CDVAE+ (Long)		-	100	99.82	84.76	-	0.1377	0.266
SyMat		-	100	99.42	97.17	-	0.1234	3.9628
SyMat+ (Short)		-	100	99.52	97.2	-	0.1206	3.7422
SyMat+ (Long)		-	100	99.9	97.63	-	0.1171	3.862
DiffCSP		-	99.9	99.49	97.27	-	0.0861	0.0876
DiffCSP+ (Short)		-	100	99.61	97.29	-	0.0811	0.087
DiffCSP+ (Long)		-	100	99.93	97.33	-	0.0763	0.0853
TGDMat(Short)		-	100	99.81	91.77	-	0.0681	0.0865
TGDMat(Long)		-	100	99.91	92.43	-	0.0436	0.0632
MP 20		CDVAE	86.3	100	99.15	99.49	1.4921	0.7085
	CDVAE+ (Short)	87.05	100	99.36	99.6	0.993	0.642	0.297
	CDVAE+ (Long)	87.42	100	99.57	99.81	0.972	0.6388	0.2977
	SyMat	87.96	99.9	98.3	99.37	0.5236	0.4012	0.3877
	SyMat+ (Short)	88.08	99.9	98.59	99.47	0.5031	0.3917	0.3622
	SyMat+ (Long)	88.47	99.9	99.01	99.95	0.4865	0.3879	0.3489
	DiffCSP	83.25	100	99.41	99.76	0.3411	0.3802	0.1497
	DiffCSP+ (Short)	84.57	100	99.52	99.85	0.331	0.38	0.1379
	DiffCSP+ (Long)	85.07	100	99.81	99.89	0.3122	0.3799	0.1355
	TGDMat(Short)	86.6	100	99.79	99.88	0.3337	0.3296	0.1189
	TGDMat(Long)	92.97	100	99.89	99.95	0.289	0.3082	0.1154

Table 9: Summary of the Complete and Detailed Results on the Gen Task.

Method	Global Features in Text Prompt	% of Matched Materials		
		Perov-5	Carbon-24	MP-20
TGDMat(Long)	Formula	97.50	98.20	70.54
	Space Group	87.00	80.79	67.88
	Crystal System	92.60	91.55	73.54
	Formation Energy	95.49	-	92.88
	Band Gap	-	98.61	96.73
TGDMat(Short)	Formula	90.70	92.56	65.22
	Space Group	86.51	80.50	58.77
	Crystal System	83.19	81.64	72.77
	Formation Energy	90.33	-	91.00
	Band Gap	-	95.90	93.33

Table 10: Summary of results on % of generated materials matching different global features specified by the textual prompts.

F.5 CHOICE OF TEXT ENCODER

Further, we investigate the expressiveness of textual representation during the reverse diffusion process. In particular, we are interested in understanding whether there are any benefits we are gaining from using a domain-specific pre-trained text encoder MatSciBERT. We conduct an ablation study where we substitute MatSciBERT with pre-trained BERTDevlin et al. (2018) model (which is domain agnostic) as text encoder in TGDMat and evaluate the performance on both tasks. The results presented in Table 11 demonstrate that MatSciBERT surpasses BERTDevlin et al. (2018) in performance for both tasks. This highlights the richer expressiveness of contextual representation achieved through the use of a domain-specific pre-trained language model.

Table 11: Ablation study results on different choices of Text Encoders.

Text Encoder	Perov-5		Carbon-24		MP-20	
	MR \uparrow	RMSE \downarrow	MR \uparrow	RMSE \downarrow	MR \uparrow	RMSE \downarrow
BERT	96.64	0.0109	72.21	0.2679	79.53	0.057
MatSciBERT	98.63	0.0072	95.27	0.1534	82.02	0.039
	Comp \uparrow	Struct \uparrow	Comp \uparrow	Struct \uparrow	Comp \uparrow	Struct \uparrow
BERT	97.44	99.97	-	100	84.73	98.37
MatSciBERT	98.63	100	-	100	92.97	100

F.6 PERFORMANCE ON MORE SHORTER PROMPTS

In this section, we explore the generalizability and robustness of our model by examining potential variability in text description lengths. The goal of this paper is, given the text prompt, to generate specific material, not any generic or class of materials. Hence some minimum essential information about the crystal, like formula, space group, crystal system, property value, etc must be given as input to the pre-trained model. However, to investigate the robustness of our proposed TGDMat model with more custom and shorter prompts, we did an experiment where we evaluated TGDMat (trained with full text) with even shorter custom prompts with very little information as follows:

- **Specifying only Formula:** *"The chemical formula is GaSiSO₂. The elements are Ga, Si, S, O. Generate the material."*
- **Specifying only Space Group Info:** *"The spacegroup number is 1. Generate the material."*
- **Specifying only Property Info:** *"The formation energy per atom is positive. Generate the material."*

We report the results in table 12. We observe that though TGDMat can handle more custom prompts, but it affects the quality of generated materials. Hence we conclude some minimum essential information about the crystal must be given as input to TGDMat to generate high-quality crystal materials.

Text Information	Perov-5		Carbon-24		MP-20	
	Comp(%) \uparrow	Struct (%) \uparrow	Comp(%) \uparrow	Struct (%) \uparrow	Comp(%) \uparrow	Struct (%) \uparrow
Only Formula	97.06	99.19	-	98.76	86.16	96.01
Only Space Group	85.91	98.97	-	95.39	84.22	96.88
Only Property	96.62	98.53	-	94.21	86.53	91.73
Full Text	98.28	100	-	100	86.60	100

Table 12: Summary of results on generated materials using more custom/shorter Prompt.

F.7 UTILITY OF TEXT-GUIDANCE THAN FEATURE VECTORS-GUIDANCE

In this subsection, we conducted an additional experiment, where we fed all relevant conditional information e.g. Formula, Space Group, Crystal Symmetry, Bond Length and Property Values as feature vectors to guide the diffusion model. We report the results for Gen Task in Table 13. Across three datasets, we did not observe performance improvements, which proved text-guidance is superior than using value guidance specifying a single or a handful of target properties or features as feature vectors.

Dataset	Method	Validity		Coverage			Property	
		Comp	Struct	Cov-R	Cov-P	# Element	Density	form_energy
Perov-5	TGDMat (Features)	96.55	98.73	99.18	97.06	0.0149	0.12	0.029
	TGDMat (Text Emb)	98.63	100	99.87	99.52	0.009	0.049	0.018
Carbon-24	TGDMat (Features)	-	99.56	99.32	92.17	-	0.104	0.087
	TGDMat (Text Emb)	-	100	99.91	92.43	-	0.043	0.063
MP-20	TGDMat (Features)	83.85	99.61	98.97	98.3	0.377	0.375	0.126
	TGDMat (Text Emb)	92.97	100	99.89	99.95	0.289	0.308	0.115

Table 13: Comparison of using features vs using text embedding of features

F.8 ABLATION STUDY FOR JOINT LEARNING OF CRYSTAL GEOMETRY.

In this subsection, we conducted an ablation study where we use three diffusion models to learn A,X,L separately. While sampling we sample A,X,L separately and merge them together. We fuse the textual representation in the same way in all three diffusion models. We present the results in following table and compare with TDGMat in Table 14. We observe significant performance degradation in all metrics across datasets if we learn A,X,L separately.

Dataset	Method	Validity		Coverage			Property	
		Comp	Struct	Cov-R	Cov-P	# Element	Density	form_energy
Perov-5	TGDMat (separately)	90.1	85.43	85.77	83.51	0.341	0.591	0.376
	TGDMat (Jontly)	98.63	100	73.57	99.52	0.009	0.049	0.018
Carbon-24	TGDMat (separately)	-	75.64	80.95	-	-	0.435	0.584
	TGDMat (Jontly)	-	100	99.91	92.43	-	0.043	0.063
MP-20	TGDMat (separately)	73.18	77.01	82.99	72.41	0.861	0.887	0.634
	TGDMat (Jontly)	92.97	100	99.89	99.95	0.289	0.308	0.115

Table 14: Comparison of learning features independently vs learning features jointly

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

F.9 MORE VISUALIZATION ON **PEROV-5**, **CARBON-24** AND **MP-20**

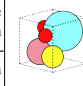
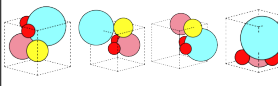
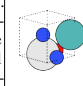
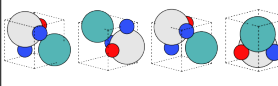
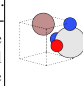
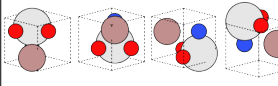
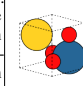
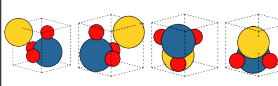
Detailed Description	Short Prompt	Ground truth	Generated Samples
<p>YCoSO₂ crystallizes in the orthorhombic Pmm2 space group. Y is bonded in a distorted square co-planar geometry to two equivalent S, two equivalent O, and two equivalent O atoms. Both Y-S bond lengths are 2.74 Å. Both Y-O bond lengths are 2.24 Å. There is one shorter (2.09 Å) and one longer (2.39 Å) Y-O bond length. Co is bonded in a distorted see-saw-like geometry to two equivalent S and two equivalent O atoms. Both Co-S bond lengths are 2.31 Å. Both Co-O bond lengths are 2.40 Å. S is bonded in a 6-coordinate geometry to two equivalent Y, two equivalent Co, and two equivalent O atoms. Both S-O bond lengths are 2.30 Å. There are two inequivalent O sites. In the first O site, O is bonded in a rectangular see-saw-like geometry to two equivalent Y and two equivalent Co atoms. In the second O site, O is bonded in a distorted square co-planar geometry to two equivalent Y and two equivalent S atoms.</p>	<p>Below is a description of a bulk material. The chemical formula is YCoSO₂. The elements are Y, Co, S, O. The formation energy per atom is positive. The spacegroup number is 24. The crystal system is orthorhombic. Generate the material:</p>		
<p>ScMoN₂O is (Cubic) Perovskite-derived structured and crystallizes in the tetragonal P4mm space group. Sc is bonded to four equivalent N and two equivalent O atoms to form ScN₄O₂ octahedra that share corners with six equivalent ScN₄O₂ octahedra and faces with eight equivalent MoN₈O₄ cuboctahedra. The corner-sharing octahedral tilt angles range from 0-1°. All Sc-N bond lengths are 2.00 Å. There is one shorter (2.00 Å) and one longer (2.01 Å) Sc-O bond length. Mo is bonded to eight equivalent N and four equivalent O atoms to form MoN₈O₄ cuboctahedra that share corners with twelve equivalent MoN₈O₄ cuboctahedra, faces with six equivalent MoN₈O₄ cuboctahedra, and faces with eight equivalent ScN₄O₂ octahedra. There are four shorter (2.83 Å) and four longer (2.84 Å) Mo-N bond lengths. All Mo-O bond lengths are 2.83 Å. N is bonded in a linear geometry to two equivalent Sc and four equivalent Mo atoms. O is bonded in a linear geometry to two equivalent Sc and four equivalent Mo atoms. The formation energy per atom is 1.8931.</p>	<p>Below is a description of a bulk material. The chemical formula is ScMoN₂O. The elements are Sc, Mo, N, O. The formation energy per atom is positive. The spacegroup number is 98. The crystal system is tetragonal. Generate the material:</p>		
<p>ScNO₂Ga is alpha Rhenium trioxide-derived structured and crystallizes in the orthorhombic Pmm2 space group. The structure consists of one Ga cluster inside a ScNO₂ framework. In the Ga cluster, Ga is bonded in a 1-coordinate geometry to atoms. In the ScNO₂ framework, Sc is bonded to two equivalent N, two equivalent O, and two equivalent O atoms to form corner-sharing ScN₂O₄ octahedra. The corner-sharing octahedral tilt angles range from 0-1°. Both Sc-N bond lengths are 2.09 Å. Both Sc-O bond lengths are 2.09 Å. Both Sc-O bond lengths are 2.09 Å. N is bonded in a linear geometry to two equivalent Sc atoms. There are two inequivalent O sites. In the first O site, O is bonded in a linear geometry to two equivalent Sc atoms. In the second O site, O is bonded in a linear geometry to two equivalent Sc atoms. The formation energy per atom is 1.4796.</p>	<p>Below is a description of a bulk material. The chemical formula is ScGaNO₂. The elements are Sc, Ga, N, O. The formation energy per atom is positive. The spacegroup number is 24. The crystal system is orthorhombic. Generate the material:</p>		
<p>OsAuO₃ is (Cubic) Perovskite structured and crystallizes in the cubic Pm-3m space group. Os is bonded to six equivalent O atoms to form OsO₆ octahedra that share corners with six equivalent OsO₆ octahedra and faces with eight equivalent AuO₁₂ cuboctahedra. The corner-sharing octahedra are not tilted. All Os-O bond lengths are 1.97 Å. Au is bonded to twelve equivalent O atoms to form distorted AuO₁₂ cuboctahedra that share corners with twelve equivalent AuO₁₂ cuboctahedra, faces with six equivalent AuO₁₂ cuboctahedra, and faces with eight equivalent OsO₆ octahedra. All Au-O bond lengths are 2.79 Å. O is bonded in a linear geometry to two equivalent Os and four equivalent Au atoms. The formation energy per atom is 1.4248.</p>	<p>Below is a description of a bulk material. The chemical formula is OsAuO₃. The elements are Os, Au, O. The formation energy per atom is 1.4248. The spacegroup number is 220. The crystal system is cubic. Generate the material.</p>		

Table 15: Visualization of the generated structures given textual description for **Perov-5** dataset

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

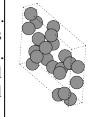
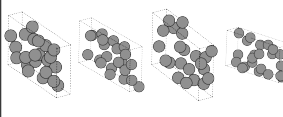
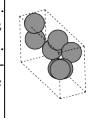
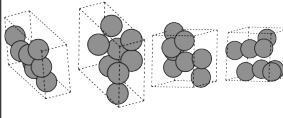
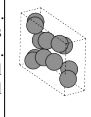
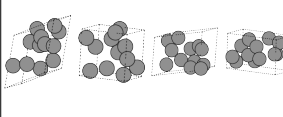
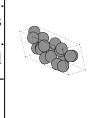
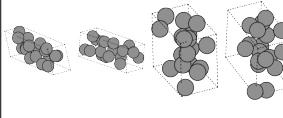
Detailed Description	Short Prompt	Ground truth	Generated Samples
C crystallizes in the triclinic P1 space group. There are twenty-two inequivalent C sites. In the first C site, C(1) is bonded to one C(18), one C(5), and two equivalent C(9) atoms to form corner-sharing CC4 tetrahedra. ... two equivalent C(12) atoms to form a mixture of distorted corner and edge-sharing CC4 trigonal pyramids. The energy per atom is -154.1336.	Below is a description of a bulk material. The chemical formula is C. The elements are C. The energy per atom is negative. The spacegroup number is 0. The crystal system is triclinic. Generate the material		
C crystallizes in the orthorhombic Cmc21 space group. There are two inequivalent C sites. In the first C site, C(1) is bonded to one C(2) and three equivalent C(1) atoms to form a mixture of corner and edge-sharing CC4 trigonal pyramids. The C(1)-C(2) bond length is 1.49 Å. There are two shorter (1.51 Å) and one longer (1.56 Å) C(1)-C(1) bond length. In the second C site, C(2) is bonded to one C(1) and three equivalent C(2) atoms to form corner-sharing CC4 tetrahedra. There are two shorter (1.54 Å) and one longer (1.56 Å) C(2)-C(2) bond length. The energy per atom is -154.2425.	Below is a description of a bulk material. The chemical formula is C. The elements are C. The energy per atom is negative. The spacegroup number is 62. The crystal system is orthorhombic. Generate the material.		
C crystallizes in the triclinic P-1 space group. There are six inequivalent C sites. In the first C site, C(1) is bonded to one C(3), one C(5), and two equivalent C(4) atoms to form corner-sharing CC4 tetrahedra. ... In the sixth C site, C(6) is bonded to one C(2), one C(4), and two equivalent C(3) atoms to form distorted corner-sharing CC4 tetrahedra. The energy per atom is -154.1338.	Below is a description of a bulk material. The chemical formula is C. The elements are C. The energy per atom is negative. The spacegroup number is 1. The crystal system is triclinic. Generate the material		
C is a Theoretical Carbon Structure-like structure and crystallizes in the triclinic P-1 space group. There are nine inequivalent C sites. In the first C site, C(1) is bonded to one C(5), one C(6), one C(7), and one C(8) atom to form a mixture of corner and edge-sharing CC4 tetrahedra. The C(1)-C(5) bond length is 1.51 Å. The C(1)-C(6) bond length is 1.56 Å. The C(1)-C(7) bond length is 1.54 Å. ... In the ninth C site, C(9) is bonded to one C(4), one C(6), one C(7), and one C(8) atom to form a mixture of corner and edge-sharing CC4 tetrahedra. The energy per atom is -154.2197.	Below is a description of a bulk material. The chemical formula is C. The elements are C. The energy per atom is -154.2197. The spacegroup number is 1. The crystal system is triclinic. Generate the material.		

Table 16: Visualization of the generated structures given textual description for **Carbon-24** dataset

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

Detailed Description	Short Prompt	Ground truth	Generated Samples
Eu2PbCl4 is Caswellsilverite-like structured and crystallizes in the tetragonal P4/mmm space group. There are two inequivalent Eu sites. In the first Eu site, Eu(1) is bonded to two equivalent P(1) and four equivalent Cl(1) atoms to form EuP2Cl4 octahedra that share corners with six equivalent Eu(1)P2Cl4 octahedra, edges with four equivalent Eu(1)P2Cl4 octahedra, and edges with eight equivalent Eu(2)P4Cl2 octahedra. ... The corner-sharing octahedra are not tilted. The formation energy per atom is -1.7615. The band gap is zero. The energy above the convex hull is zero.	Below is a description of a bulk material. The chemical formula is Eu2PbCl4. The elements are Eu, P, and Cl. The formation energy per atom is -1.7615. The band gap is 0.0. The energy above the convex hull is 0.0. The spacegroup number is 122. The crystal system is tetragonal. Generate the material.		
MgNdHg2 is Heusler structured and crystallizes in the cubic Fm-3m space group. Mg(1) is bonded in a body-centered cubic geometry to eight equivalent Hg(1) atoms. All Mg(1)-Hg(1) bond lengths are 3.18 Å. Nd(1) is bonded in a body-centered cubic geometry to eight equivalent Hg(1) atoms. All Nd(1)-Hg(1) bond lengths are 3.18 Å. Hg(1) is bonded in a body-centered cubic geometry to four equivalent Mg(1) and four equivalent Nd(1) atoms. The formation energy per atom is -0.4708. The band gap is 0.0. The energy above the convex hull is 0.0. The spacegroup number is 224.	Below is a description of a bulk material. The chemical formula is NdMgHg2. The elements are Nd, Mg, and Hg. The formation energy per atom is -0.4708. The band gap is 0.0. The energy above the convex hull is 0.0. The spacegroup number is 224. The crystal system is cubic. Generate the material.		
MgNdTi crystallizes in the hexagonal P-62m space group. Mg(1) is bonded in a 4-coordinate geometry to two equivalent Ti(1) and two equivalent Ti(2) atoms. Both Mg(1)-Ti(1) bond lengths are 3.01 Å. Both Mg(1)-Ti(2) bond lengths are 3.03 Å. Nd(1) is bonded in a 5-coordinate geometry to one Ti(2) and four equivalent Ti(1) atoms. The Nd(1)-Ti(2) bond length is 3.31 Å. All Nd(1)-Ti(1) bond lengths are 3.32 Å. There are two inequivalent Ti sites. In the first Ti site, Ti(2) is bonded in a distorted q6 geometry to six equivalent Mg(1) and three equivalent Nd(1) atoms. In the second Ti site, Ti(1) is bonded in a 9-coordinate geometry to three equivalent Mg(1) and six equivalent Nd(1) atoms. The formation energy per atom is -0.355. The band gap is 0.0. The energy above the convex hull is 0.0. The spacegroup number is 188.	Below is a description of a bulk material. The chemical formula is NdMgTi. The elements are Nd, Mg, and Ti. The formation energy per atom is -0.355. The band gap is 0.0. The energy above the convex hull is 0.0. The spacegroup number is 188. The crystal system is hexagonal. Generate the material.		
LaNi2Ge2 crystallizes in the tetragonal I4/mmm space group. La(1) is bonded in a 16-coordinate geometry to eight equivalent Ni(1) and eight equivalent Ge(1) atoms. All La(1)-Ni(1) bond lengths are 3.25 Å. All La(1)-Ge(1) bond lengths are 3.26 Å. Ni(1) is bonded in a 4-coordinate geometry to four equivalent La(1) and four equivalent Ge(1) atoms. All Ni(1)-Ge(1) bond lengths are 2.39 Å. Ge(1) is bonded in a 9-coordinate geometry to four equivalent La(1), four equivalent Ni(1), and one Ge(1) atom. The Ge(1)-Ge(1) bond length is 2.66 Å. The formation energy per atom is -0.691. The band gap is 0.0. The energy above the convex hull is 0.0.	Below is a description of a bulk material. The chemical formula is La(NiGe)2. The elements are La, Ni, and Ge. The formation energy per atom is -0.691. The band gap is 0.0. The energy above the convex hull is 0.0. The spacegroup number is 138. The crystal system is tetragonal. Generate the material.		

Table 17: Visualization of the generated structures given textual description for **MP-20** dataset