
Unleashing the power of novel conditional generative approaches for new materials discovery

Anonymous Author(s)
Affiliation
Address
email

Abstract

1 For a very long time, computational approaches to the design of new materials have relied on an iterative process of finding a candidate material
2 and modeling its properties. AI has played a crucial role in this regard,
3 helping to accelerate the discovery and optimization of crystal properties
4 and structures through advanced computational methodologies and data-
5 driven approaches. To address the problem of new materials design and
6 fasten the process of new materials search, we have applied latest generative
7 approaches to the problem of crystal structure design, trying to solve the
8 inverse problem: by given properties generate a structure that satisfies them
9 without utilizing supercomputer powers. In our work we propose two ap-
10 proaches: 1) conditional structure modification: optimization of the stability
11 of an arbitrary atomic configuration, using the energy difference between the
12 most energetically favorable structure and all its less stable polymorphs and
13 2) conditional structure generation. We used a representation for materials
14 that includes the following information: lattice, atom coordinates, atom
15 types, chemical features, space group and formation energy of the structure.
16 The loss function was optimized to take into account the periodic boundary
17 conditions of crystal structures. We have applied Diffusion models approach,
18 Flow matching, usual Autoencoder (AE) and compared the results of the
19 models and approaches. As a metric for the study, physical pymatgen
20 matcher was employed: we compare target structure with generated one
21 using default tolerances. So far, our modifier and generator produce struc-
22 tures with needed properties with accuracy 41% and 82% respectively. To
23 prove the offered methodology efficiency, inference have been carried out,
24 resulting in several potentially new structures with formation energy below
25 the AFLOW-derived convex hulls.
26

27 1 Introduction

28 The search for novel materials with specified properties has been a cornerstone of scientific
29 exploration for decades. From the discovery of semiconductors revolutionizing electronics
30 to the development of superalloys enhancing aerospace technologies, the synthesis of new
31 materials has continually propelled technological advancements.

32 However, traditional methods for material discovery often employ exhaustive trial and error
33 experimental approaches. In turn, computational efforts, relying on density functional theory
34 (DFT)[1] approaches, usually require huge amounts of computing power. In this regard,
35 automatic descriptor generators[2], GNNs[3][4] and transferable GNN models [5] fueled
36 combination of these methods and machine learning (ML) approaches. In particular, the

37 utilization of generative machine learning models, such as Variational Autoencoder[6] and
38 GANs[7], presents a paradigm shift in how crystal structures are generated and optimized.
39 By harnessing the power of data-driven approaches, we can navigate the vast landscape of
40 possible crystal structures with unprecedented efficiency and precision.

41 Recent advancements in the field of materials discovery have yielded promising results
42 through various innovative approaches. For instance, FTCP[6] utilizes Autoencoders for
43 uncovering new materials, while CubicGAN[7] leverages GANs for the discovery of cubic
44 crystal materials. Additionally, Physics Guided Crystal Generative Model (PGCGM)[8] has
45 introduced a method for generating crystal structures based on specific space groups encoding.
46 DP-CDVAE[9] is a model, that combines VAE and diffusion approaches. MatterGen[10]
47 employed equivariant GNNs as score matching function in diffusion processes for crystal
48 structure generation.

49 One of the most discussed frameworks is GNoME[11] that has made most recent and large
50 advancements in the field of the new materials discovery employs a sophisticated pipeline
51 to discover new materials, particularly focusing on inorganic crystals. This allows for the
52 discovery of innovative materials beyond known structures.

53 After generating candidate structures through both pipelines, GNoME evaluates their stability
54 by predicting their formation energies. Based on the comparison of the obtained formation
55 energy with those of the known competing phases (i.e. stability assessment), the model selects
56 the most promising candidates for further evaluation using known theoretical frameworks.

57 The question of the completeness of chemical space arises due to two main concerns with
58 GNoME-derived stable structures. Firstly, they mostly contain three or more unique elements,
59 while ternary and quaternary structures are less explored than binary compounds. Secondly,
60 the comparison of GNoME-discovered structures to the Materials Project, which has 154,718
61 materials, is flawed since larger databases like AFLOW, NOMAD, and the Open Quantum
62 Materials Database contain millions of entries. This raises questions about the novelty of
63 the discovered materials.

64 In this study, we present an end-to-end framework for the generation of crystal structures with
65 specified properties using advanced generative AI techniques. The basis architecture of the
66 models is Autoencoder, enabling encoding and decoding structural representations. Then, the
67 most commonly used generative approaches in image generation were utilized to model prob-
68 ability distribution transformations, and to capture complex underlying structure-property
69 relationships within our dataset: Flow Matching[12], Denoising Diffusion Probabilistic
70 Models(DDPM)[13], and Denoising Diffusion Implicit Models(DDIM)[14]. Through the
71 integration of these techniques, we aim to transcend conventional limitations in materials
72 discovery, paving the way for accelerated predictions of materials with desired properties.

73 To employ model architectures often used for image/video generation, a matrix represen-
74 tation of crystal structures was developed, containing crucial information such as chemical
75 composition, atomic coordinates, symmetries (space group), and formation energies. Within
76 the approach proposed, it has become important to develop a novel metric for assessing
77 the similarity between generated structures and target configurations. This metric obviates
78 the need for computationally expensive DFT calculations, allowing for rapid validation and
79 refinement of generated structures. Furthermore, we introduce a loss function that accounts
80 for the periodic boundary conditions inherent in crystal lattices, ensuring the fidelity of the
81 generated structures.

82 Our study explores two distinct approaches for crystal structure prediction: 1) conditional
83 structure modification and 2) conditional structure generation. The former involves optimiz-
84 ing the stability of existing structures by generating more stable polymorphs, while the latter
85 entails the generation of entirely new structures based on user-defined criteria. Through
86 rigorous analysis, we demonstrate the efficacy of our approach in discovering novel materials
87 with desired properties.

88 Importantly, to validate the utility of our framework, we conducted a series of generation
89 experiments using the Vienna Ab initio simulation package(VASP)[15] as a tool for inference
90 validation. Remarkably, our method facilitated the discovery of 6 structures below the
91 corresponding convex hull. This significant outcome underscores the remarkable potential of

92 our framework in uncovering thermodynamically stable materials, thereby offering promising
93 avenues for advanced materials discovery and design.

94 2 Data, Dataset

95 2.1 Data overview

96 In this study, the AFLOW database[16] was utilized as a source of data on the structures
97 and properties of materials. AFLOW is an extensive and comprehensive database that
98 consolidates a vast array of materials-related information, offering an expansive repository
99 for crystallographic data, computed properties, and various other materials-science-related
100 datasets. AFLOW database contains more than 3.5 million structures.

101 From the extensive collection housed within AFLOW, the focus was narrowed to select only
102 polymorphs, because models are trained to distinguish composition-property and structure-
103 property relations with numerous structures of the same chemical composition. Specifically,
104 the selection process targeted polymorphic structures with 4 to 60 atoms within their unit
105 cells. This criterion aimed to encompass a diverse yet manageable subset of structures,
106 balancing complexity with computational feasibility. By filtering polymorphs based on their
107 atom count, the dataset was balanced.

108 Moreover, in order to decrease the complexity of the data, we have removed all structures
109 containing elements and space groups found in less than 1% of all structures. The entire
110 dataset consisted of more than 85000 polymorph groups including more than 2.1 million
111 structures. The minimum size of group of polymorphs was 7 samples and the maximum one
112 was 71 samples. The total number of space groups was 19 and the total number of chemical
113 species over the dataset was 55. Each structure S in the dataset is described by the following
114 features:

- 115 • Fractional coordinates of atoms in the lattice basis X_{coord} (has 60 rows with 3
116 coordinates x, y, z each) and X_{lattice} (matrix 3 by 3 constructed of 3 base vectors).
117 Overall matrix X of structure is constructed as

$$X_{64 \times 3} = \text{concatenation}(X_{\text{coord}}_{60 \times 3}, \text{padding}_{1 \times 3}, X_{\text{lattice}}_{3 \times 3}) \quad (1)$$

- Chemical elements which are presented as a one-hot matrix $elements_{ij}$ of size 64×103 (including padding), where ones are positioned at the indices corresponding to the position of a certain chemical element in the periodic table.

$$elements_{ij} = \begin{cases} 1 & \text{if } i\text{-th atom's element number from the periodic table} = j \\ 0 & \text{otherwise} \end{cases}$$

- 118 • Elemental property matrix $elementalProperties$ containing 22 chemical features
119 encoding chemical elements obtained from [8]. The properties of each element were
120 calculated using Mendeleev package[17].
- 121 • Space group spg of a structure. We use the space group encoding method presented in
122 [8], when each space group is represented by a $192 \times 4 \times 4$ matrix, which corresponds
123 to 192 possible symmetry operations.
- 124 • Structure formation energy E
- 125 • Nsites - number of atoms in a crystal lattice.

126 2.2 Data representation. Modification task

127 The crystal pair sampling strategy involves handling a potential data leakage: possible
128 inclusion of structures from the same polymorph group but with different energies into training
129 and validation subsets. To mitigate this issue, the polymorph group formulas were initially
130 divided into distinct training and validation sets, ensuring a relatively balanced distribution
131 of chemical elements across these subsets. Subsequently, the pairs were categorized into two
132 groups: those with low-energy (lowest energy in polymorph group) targets designated as

133 $lowestEnergyPairs = (S_i, S_0) \forall i \in [1, \dots, structuresNum]$ and those with non-low-energy
 134 targets, all structures except the most optimal one, formed as $nonLowestEnergyPairs =$
 135 $(S_i, S_j) | i > j > 0$. The validation set was constructed as a subset of $lowestEnergyPairs$.
 136 The training set was dynamically constructed every epoch from $lowestEnergyPairs$ and
 137 $nonlowestEnergyPairs$, preserving equal numbers of pairs sampled and maintaining a
 138 limited count per polymorph group. This strategy ensured a robust separation between
 139 training and validation sets, thus preventing data leakage and improving model performance.

140 Each pair sample $\{S_{init}, S_{target}\} \in pairDataset$ consisted of the information about each
 141 structure (hereinafter, we will call them initial and target structures). The following data
 142 was used:

- 143 • Coordinates and lattice information of initial and target structures X_{init}, X_{target}
- 144 • Difference in formation energies between initial and target structures $E_{diff} =$
 145 $E_{target} - E_{init}$
- 146 • Space group of target structure spg_{target}
- 147 • Elements matrix $elemetsMatrix$, elemental property matrix $elementalProperties$
 148 and number of sites $numSites$, which are the same for initial and target structure
 149 because of identical chemical composition.

150 The modification task involved transforming the input structure X_{init} into the target structure
 151 X_{target} .

152 2.3 Data representation. Generation task

153 In its tern the generation task receives normal or uniform (depends on a model) noise as
 154 input from which the structure is generated, which is akin to the image generation processes
 155 in computer vision tasks.

156 For the generation task, an additional dataset was constructed. Data for the generation task
 157 is slightly simpler, while it considers only $\{S_{target}\}$. Therefore, the models can be trained on
 158 all structures available, rather than pairs. The following data is used:

- 159 • Coordinates and lattice information of target structure X_{target}
- 160 • Formation energy of target structure E_{target}
- 161 • Space group of target structure spg_{target}
- 162 • Elements matrix $elemetsMatrix$, elemental property matrix $elementalProperties$
 163 and number of sites $numSites$ of target structure.

164 3 Loss and metrics

165 3.1 Atomic coordinates

166 The atomic coordinates are represented as a 60×3 matrix, where each row corresponds
 167 to the coordinates of an atom. The L1 loss was utilized during the training of a model for
 168 predicting atomic coordinates.

169 $L_1(preds, target)_i = ||preds_i - target_i||_1 = \sum_{j=1}^3 |preds_{ij} - target_{ij}|$, where $target$ and $pred$
 170 are target and predicted atomic coordinate matrices.

171 3.2 Lattice

172 The lattice itself is represented as a 3×3 matrix, where each row signifies a directing basis
 173 vector. In this case, we have also used the L1 norm as a loss function.

174 3.3 Periodic boundary condition loss

175 This section presents an enhanced loss function, designed for the regression model (see
 176 Section 5.1), that addresses this challenge by integrating periodic boundary conditions into

177 the loss calculation, outperforming the conventional L1 loss function. In the field of ML
 178 applied to atomic structures, even slight displacement of atomic coordinates is crucial and
 179 employing appropriate loss functions that consider the periodic nature of atomic structures
 180 increases the flexibility of model predictions.

181 In the dataset representing atomic structures, it is crucial to acknowledge the presence
 182 of atoms residing at various positions within the lattice framework. Certain atoms are
 183 positioned at the vertices, edges, or faces of the lattice. According to periodic boundary
 184 conditions (PBC), identical atoms in the vicinity of vertices, edges, or faces but also exist in
 185 analogous positions across the lattice. Implementation of such an invariance within the loss
 186 function helps in effectively capturing periodic pattern of crystals, enhancing the model’s
 187 capability to learn and predict atomic structures more comprehensively.

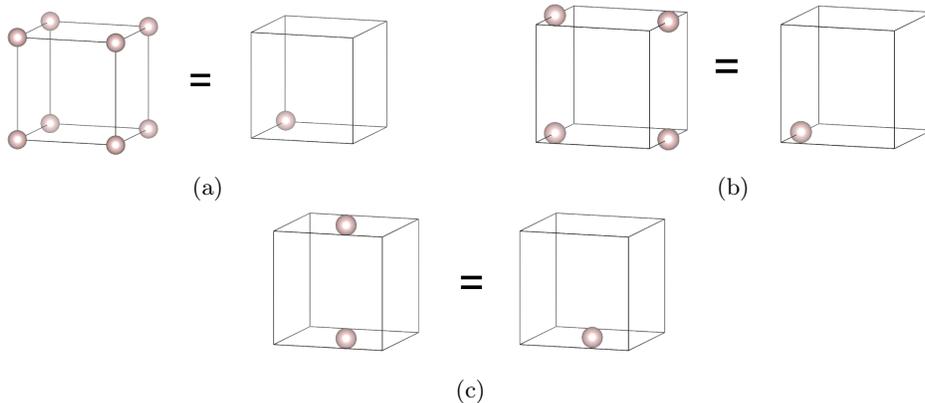


Figure 1: Illustration of atoms at a)vertices, b)edges, and c)faces of lattice under periodic boundary conditions

188 The loss function is being calculated as minimum of distances from predicted point to the
 189 target one taking into account 26 its periodic images (according to PBC) A.4.

190 The empirical validation of this enhanced loss function showcases its superiority(Figure4) in
 191 capturing discrepancies within atomic structures, thus indicating its potential as a robust
 192 tool for improving the accuracy of ML models in materials science applications.

193 3.4 Metric

194 As a metric, we have chosen an analogue of accuracy: the generated structures are compared
 195 to the target structures using a specialized matcher, yielding the proportion of structures that
 196 successfully pass the matching process. For metric calculation, we employed the Pymatgen
 197 StructureMatcher with the default set of parameters ($ltol = 0.2$, $stol = 0.3$, $angle_tol = 5$).
 198 Although this approach is less accurate than structure relaxation using ab initio calculations
 199 and comparing the structure formation energy with the energy above the hull, it enables
 200 model validation to be performed orders of magnitude faster than the traditional method.

201 4 Model

202 For experiments, a 1d UNet model (see Figure2 (b)) architecture similar to the 2d UNet model
 203 described in [18] was utilized along with 2D and 1D convolutional neural networks (CNNs)
 204 for the space group and element matrix embeddings, respectively. Based on this model, 3
 205 different training processes have been developed: ordinary regression model, Conditional
 206 Flow Matching (CFM)[19] model, and diffusion model.

207 The model was conditioned (see Figure2 (a)) on the following data: time condition (t), the
 208 same as in [18], element condition (el), formation energy difference condition (E_{diff}), and
 209 desirable space group (spg). el_{emb} , spg_{emb} and E_{diff} are concatenated into one embedding

210 C_{emb} . t is fed into the Transformer Positional Encoding Layer and transformed into an
 211 embedding T_{emb} . The two embeddings: C_{emb} and T_{emb} are then applied into one condition.

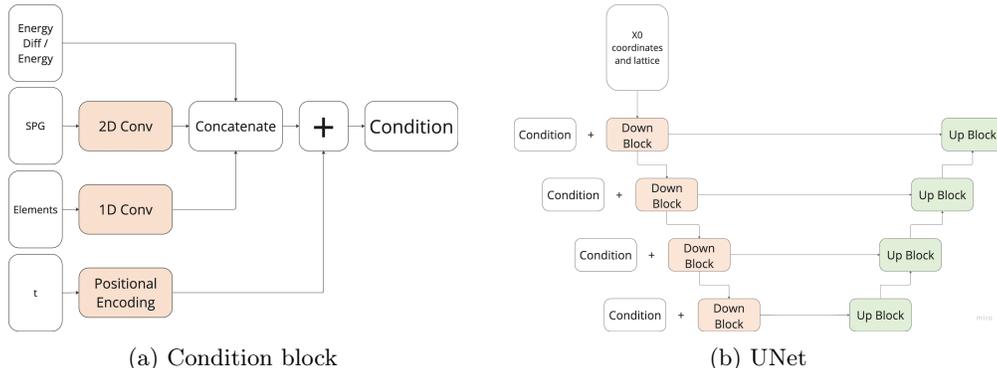


Figure 2: a)Formation of conditions using formation energy, space group, and elemental representation, and b)Schematic depiction of the model architecture

212 5 Methodology

213 In this work, two approaches are proposed: crystal structure generation and crystal structure
 214 modification. For the generation approach, crystal structures are generated from normal
 215 or uniform noise and conditioned to t , el , E , spg . Within the generation, we employed
 216 three algorithms: DDPM, DDIM, and CFM models. For the modification approach, crystal
 217 structures are generated by modifying other structures, while conditioning to el , E_{diff} , spg
 218 (and optionally t , not used in ordinary regression UNet). For the modification task, we have
 219 employed three algorithms: UNet Regression model, diffusion model, based on Palette[20]
 220 approach, and CFM model. For the generation task, we have employed four algorithms:
 221 diffusion models with DDPM and DDIM samplers, and CFM models on Uniform and Normal
 222 noise.

223 5.1 Regression model

224 During the training stage, the structure coordinates and lattice x_0 , elements features el , space
 225 group spg and E_{diff} are used as conditions. The model is trained to return x_1 structure
 226 coordinates and lattice (Algorithm 1). As for the inference process, one can see the details
 227 in the Algorithm 2

228 5.2 Conditional Flow Matching models

229 CFM is a fast method for training Continuous Normalizing Flows (CNF)[21] models without
 230 the need for simulations. It offers a training objective that enables conditional generative
 231 modeling and accelerates both training and inference.

232 The basic way of training CFM model (Algorithm 3) organized as follows: during the
 233 training stage, x_0 and x_1 are sampled from the source distribution and the target distribution
 234 respectively, then a linear interpolation x_t is calculated as $x_t = tx_1 + (1 - t)x_0$ (exponential
 235 moving average between distributions x_0 and x_1 ; t is sampled from a uniform distribution
 236 $\mathcal{U}(0, 1)$), and afterwards pass the x_t and t as inputs to our model f_θ , forcing the model to
 237 predict a velocity from the distribution x_0 to x_1 . Therefore, the loss for CFM model is
 238 the following: $L_{CFM} = E_{t,x_1,x_0}[\|f_\theta(x_t, t) - (x_1 - x_0)\|^2] = E_{t,x_1,x_0}[\|f_\theta(tx_1 + (1 - t)x_0, t) -$
 239 $(x_1 - x_0)\|^2]$

240 For the modification approach, x_0 and x_1 are both sampled from our dataset distribution
 241 according to the sampling strategy for modification mentioned in 2.2. Also, the model is
 242 conditioned to el , spg_1 , E_{diff} , besides t (see Algorithm 4)

243 For the generation approach, we tested two noise distributions for the x_0 : normal distribu-
 244 tion $\mathcal{N}(0, 1)$ and uniform noise distribution $\mathcal{U}(0, 1)$, which resulted in significantly better
 245 performance. The intuition for using uniform distribution instead of normal one was inspired
 246 by the diagram of x, y, z coordinate distribution (Figure 3). The model is also conditioned
 247 to el, spg_1, E , and t (see Algorithm 5)

248 During the sampling stage, we generate X_1 structure by the given X_0 by solving the following
 249 ordinary differential equation (ODE): $dx_t = f_\theta(x_t, t, el, spg_1, E)dt$, beginning with x_0 . In
 250 order to solve the ODE, the Euler method was employed: $x_{t+h} = x_t + hf_\theta(x_t, t, el, spg_1, E)$
 251 (Algorithm 6)

252 5.3 Diffusion models

253 In our work, we observe diffusion models. Diffusion models generate samples from a target
 254 distribution x_1 , starting from a source distribution $x_0 \sim \mathcal{N}(0, I)$.

255 During training, these models are trained to reverse a Markovian forward process, which
 256 adds noise x_0 to the data step by step. Meaning, diffusion models are trained to predict the
 257 noise added to the data samples x_1 . In order to train a model in this setup, the following loss
 258 function is used, $L_{simple} = E_{t, x_1, x_0} [||x_0 - f_\theta(\sqrt{\bar{\alpha}_t}x_1 + \sqrt{1 - \bar{\alpha}_t}x_0, t)||^2]$ where $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$
 259 and $\alpha_t = 1 - \beta_t$ (β_t is the variance by which added noise is being scheduled on each step t).

260 Our modification approach is based on Palette, which enables sample-to-sample generation
 261 (from noise $\epsilon \sim \mathcal{N}(0, 1)$) using x_0 structure coordinates and lattice, el, spg_1, E_{diff} and t as
 262 conditions for generation of x_1 using the DDPM algorithm. Sampling stage is performed by
 263 a backward diffusion process with linear scheduler (see Algorithms 7, 8).

264 For the generation approach ((Algorithm 9), x_0 is sampled from a normal distribution and
 265 el, spg_1, E, t are fed into the model as conditions. During our experiments, we tested
 266 2 approaches: DDPM(Algorithm 10) classic approach and DDIM(Algorithm 11) which
 267 results in usage of smaller number of sampling steps in order to speed up the generation
 268 process. Moreover, DDIM enables the process of generating samples from random noise to
 269 be deterministic.

270 6 Experiment Results

271 All the models presented in tables (Table 1 and Table 2) have been trained with the same
 272 hyperparameters and architectures. The metric used is described in Section 3.4. We also
 273 provide all experiment details in A.3.

Table 1: Validation metrics on generation task

DDPM	DDIM	CFM $\mathcal{N}(0, 1)$	CFM $\mathcal{U}(0, 1)$
0.8074	0.82	0.482	0.8097

Table 2: Validation metrics on modification task

Ordinary Model	Diffusion	CFM
0.4148	0.3653	0.2059

274 7 Inference

275 In order to demonstrate a potential of the proposed approaches, we have chosen a chemical
 276 composition, containing numerous variations and phases of structures composed of [W, B,
 277 Ta] with well-explored convex hull. Structures that lie on the convex hull are considered to
 278 be thermodynamically stable, and the ones above it are either metastable or unstable.

279 7.1 Inference pipeline

280 The proposed testing procedure involves generating test conditions for structures, passing
281 them to the trained generative models, pre-optimizing the generated structures to accelerate
282 the following ab initio calculations, and final relaxation and formation energy calculating using
283 VASP. Although in this work two approaches were proposed: Generation and Modification,
284 the following pipeline has only been applied to generation models, due to the fact, that
285 modification approach is based on structure-polymorphs, which leads to the necessity to
286 have at least one structure with needed composition, which is not always so. That fact
287 makes generation models much more flexible in generation structures not only with needed
288 properties, but also with needed composition. Another advantage of the generation models is
289 value of metric that is two times bigger than in modification tasks. The inference algorithm
290 is as follows:

- 291 1. Test Condition Formation:
- 292 • The chosen chemical formulas were utilized for feature extraction of *el*. Three
293 chemical compositions have been used: 1) $\text{Ta}_1\text{W}_1\text{B}_6$, 2) $\text{Ta}_1\text{W}_2\text{B}_5$ and 3)
294 $\text{Ta}_2\text{W}_1\text{B}_5$.
 - 295 • We have taken *spg* presented in the dataset as an additional condition, obtaining
296 19 space groups.
 - 297 • Finally, a set of target formation energies *E* has been formed. We have carried
298 out three experiments: 1) starting from the energy of the convex hull and
299 decreasing with a step of 0.01 eV/atom, 2) starting from the energy of the
300 convex hull and decreasing with a step of 0.1 eV/atom, and 3) starting from the
301 energy 1 eV/atom less than the energy of the convex hull and decreasing with a
302 step of 0.01 eV/atom. In total, 21 energy values were used for every inference
303 run.
 - 304 • Final inference conditions were obtained by making all possible combinations of
305 *spg* and *E* for a certain composition *el*
- 306 2. Model Inference: The conditions from the step 1 have been put to one of the trained
307 models, resulting in the generation of structures. Two models have been employed:
308 Diffusion approach and Flow matching
- 309 3. Pre-Optimization: Following the generation of all structures, each structure has been
310 pre-optimized using the PyMatGen structure relaxation method. The method used
311 m3gnet [22] model with default parameters. PyMatGen pre-optimization contributed
312 to overall speedup of further VASP relaxation.
- 313 4. Structure relaxation: Pre-optimized structures were relaxed using VASP (the rec-
314 ommended pseudopotentials, plane wave energy cutoff of 500 eV, Ediiif and Ediffg
315 convergence criteria of 10^{-5} and -10^{-2} were used).

316 7.2 Inference results

317 To summarize, 6 experiments have been carried out for two different models and for three
318 formation energy conditionings. Every experiment includes 3*380 structures, per 380
319 structures for every single chemical composition. The results of experiments can be seen in
320 Table 3

321 As can be seen, 4 structures were obtained with formation energies significantly lower than
322 those obtained from the AFLOW-derived convex hull. Thus, it can be concluded that
323 this observation indicates the potential stability of the generated structures rather than
324 differences in the computational methods used in this work and during AFLOW generation.
325 Another four structures also have energies below the convex hull, but in the vicinity of it.
326 Thus, their potential stability should be interpreted with caution.

327 8 Data availability

328 The raw crystal dataset is downloaded from
329 <https://afloplib.org>

Table 3: Inference results. Each matrix element corresponds to either the minimal energy above the hull achieved in an experiment or the energy above the hull of structures with energies below the hull.

		Ta1W1B6, meV/atom	Ta1W2B5, meV/atom	Ta2W1B5, meV/atom
Diffusion	energy step = 0.01 energy gap = 0	10,41	3,275	13,079
	energy step = 0.1 energy gap = 0	11,835	3,83	-0,042
	energy step = 0.01 energy gap = 1	97,676	-1,409	5,539
Flow-Matching	energy step = 0.01 energy gap = 0	11,981	-0,483 -0,466 -0,387	-5,426
	energy step = 0.1 energy gap = 0	11,286	0,037	-5,497
	energy step = 0.01 energy gap = 1	9,529	-4,852	1,029

330 9 Code availability

331 The source code for training and inferencing our models can be obtained from GitHub at
 332 <https://github.com/AIRI-Institute/conditional-crystal-generation>

333 10 Conclusion

334 In this article, we have offered two approaches to generate crystal structures: conditional
 335 generation and conditional modification. The first approach is significantly more flexible
 336 as it does not require structure-polymorphs, enabling the generation of structures without
 337 restrictions on chemical composition, which can be crucial in certain scenarios. Another
 338 advantage of the first approach is the simplicity of data preprocessing; it only requires the
 339 chemical composition, space group, atom coordinates, and formation energies.

340 Our methodology has experimentally proven its effectiveness, resulting in four confident
 341 potentially new crystal structures with the following energies above the hull: $\{-1.409, -5.497,$
 342 $-5.426,$ and $-4.852\}$ meV/atom, and four uncertain candidates with energies of $\{-0.483, -0.466,$
 343 $-0.387,$ and $-0.042\}$ meV/atom. We have demonstrated that conditional generation approaches,
 344 commonly used in image generation, are also fruitful in the design of new materials.

345 Although the proposed methodology demonstrates its efficiency in generating potentially
 346 new crystal structures, it has certain limitations. Firstly, the data is represented in a matrix
 347 form, which does not account for all possible symmetries of the crystal structures. Secondly,
 348 the structures in the dataset range from 4 to 60 atoms per unit cell, with most structures
 349 containing fewer than 8 atoms per unit cell. However, to perform well on structures with a
 350 large number of atoms per unit cell, the models just should be pretrained on a dataset that
 351 includes larger structures.

352 Furthermore, despite the limited number of experiments(6) and structures generated (7182),
 353 we succeeded in identifying hypothetically new structures. We hope that our article will
 354 help to reveal the potential of generative AI in design of new materials with targeted
 355 thermodynamic properties and inspire other researchers to be part of this innovative journey
 356 in materials design. We believe that rapid and efficient generation of novel materials can lead
 357 to breakthroughs in various fields such as electronics, pharmaceuticals, and energy storage.
 358 This can accelerate technological advancements and make cutting-edge technologies more
 359 accessible and affordable.

360 References

- 361 [1] Diola Bagayoko. Understanding density functional theory (dft) and completing it in
362 practice. *AIP Advances*, 4(12), 2014.
- 363 [2] Alexander Dunn, Qi Wang, Alex Ganose, Daniel Dopp, and Anubhav Jain. Benchmark-
364 ing materials property prediction methods: the matbench test set and automatminer
365 reference algorithm. *npj Computational Materials*, 6(1):138, 2020.
- 366 [3] Roman A Eremin, Innokentiy S Humonen, Alexey A Kazakov, Vladimir D Lazarev,
367 Anatoly P Pushkarev, and Semen A Budenny. Graph neural networks for predicting
368 structural stability of cd-and zn-doped γ -cspbi3. *Computational Materials Science*,
369 232:112672, 2024.
- 370 [4] Alexey N Korovin, Innokentiy S Humonen, Artem I Samtsevich, Roman A Eremin,
371 AI Vasilev, Vladimir D Lazarev, and Semen A Budenny. Boosting heterogeneous
372 catalyst discovery by structurally constrained deep learning models. *Materials Today*
373 *Chemistry*, 30:101541, 2023.
- 374 [5] Alexandre Duval, Simon V Mathis, Chaitanya K Joshi, Victor Schmidt, Santiago
375 Miret, Fragkiskos D Malliaros, Taco Cohen, Pietro Liò, Yoshua Bengio, and Michael
376 Bronstein. A hitchhiker’s guide to geometric gnns for 3d atomic systems. *arXiv preprint*
377 *arXiv:2312.07511*, 2023.
- 378 [6] Zekun Ren, Siyu Isaac Parker Tian, Juhwan Noh, Felipe Oviedo, Guangzong Xing, Jiali
379 Li, Qiaohao Liang, Ruiming Zhu, Armin G Aberle, Shijing Sun, et al. An invertible
380 crystallographic representation for general inverse design of inorganic crystals with
381 targeted properties. *Matter*, 5(1):314–335, 2022.
- 382 [7] Yong Zhao, Mohammed Al-Fahdi, Ming Hu, Edirisuriya MD Siriwardane, Yuqi Song,
383 Alireza Nasiri, and Jianjun Hu. High-throughput discovery of novel cubic crystal
384 materials using deep generative neural networks. *Advanced Science*, 8(20):2100566,
385 2021.
- 386 [8] Yong Zhao, Edirisuriya M Dilanga Siriwardane, Zhenyao Wu, Nihang Fu, Mohammed
387 Al-Fahdi, Ming Hu, and Jianjun Hu. Physics guided deep learning for generative design
388 of crystal materials with symmetry constraints. *npj Computational Materials*, 9(1):38,
389 2023.
- 390 [9] Teerachote Pakornchote, Natthaphon Choomphon-Anomakhun, Sorjrit Arrerut,
391 Chayanon Atthapak, Sakarn Khamkao, Thiparat Chotibut, and Thiti Bovornratanaraks.
392 Diffusion probabilistic models enhance variational autoencoder for crystal structure
393 generative modeling. *Scientific Reports*, 14(1):1275, 2024.
- 394 [10] Claudio Zeni, Robert Pinsler, Daniel Zügner, Andrew Fowler, Matthew Horton, Xiang
395 Fu, Sasha Shysheya, Jonathan Crabbé, Lixin Sun, Jake Smith, et al. Mattergen: a
396 generative model for inorganic materials design. *arXiv preprint arXiv:2312.03687*, 2023.
- 397 [11] Amil Merchant, Simon Batzner, Samuel S Schoenholz, Muratahan Aykol, Gowoon
398 Cheon, and Ekin Dogus Cubuk. Scaling deep learning for materials discovery. *Nature*,
399 624(7990):80–85, 2023.
- 400 [12] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow
401 matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- 402 [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In
403 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in*
404 *Neural Information Processing Systems*, volume 33, pages 6840–6851. Curran Associates,
405 Inc., 2020.
- 406 [14] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models.
407 In *International Conference on Learning Representations*, 2021.

- 408 [15] Guangyu Sun, Jenő Kürti, Péter Rajczy, Miklos Kertesz, Jürgen Hafner, and Georg
409 Kresse. Performance of the vienna ab initio simulation package (vasp) in chemical
410 applications. *Journal of Molecular Structure: THEOCHEM*, 624(1-3):37–45, 2003.
- 411 [16] Stefano Curtarolo, Wahyu Setyawan, Gus LW Hart, Michal Jahnatek, Roman V Chep-
412 ulskii, Richard H Taylor, Shidong Wang, Junkai Xue, Kesong Yang, Ohad Levy, et al.
413 Aflow: An automatic framework for high-throughput materials discovery. *Computational*
414 *Materials Science*, 58:218–226, 2012.
- 415 [17] Łukasz Mentel. mendeleev – a python resource for properties of chemical elements, ions
416 and isotopes.
- 417 [18] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion proba-
418 bilistic models. In *International Conference on Machine Learning*, pages 8162–8171.
419 PMLR, 2021.
- 420 [19] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le.
421 Flow matching for generative modeling. In *The Eleventh International Conference on*
422 *Learning Representations*, 2023.
- 423 [20] Chitwan Saharia, William Chan, Huiwen Chang, Chris A. Lee, Jonathan Ho, Tim
424 Salimans, David J. Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion
425 models, 2022.
- 426 [21] Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural
427 ordinary differential equations. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman,
428 N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing*
429 *Systems*, volume 31. Curran Associates, Inc., 2018.
- 430 [22] Chi Chen and Shyue Ong. A universal graph deep learning interatomic potential for
431 the periodic table. *Nature Computational Science*, 2:718–728, 11 2022.

432 A Appendix section 1

433 A.1 Distribution of atomic coordinates

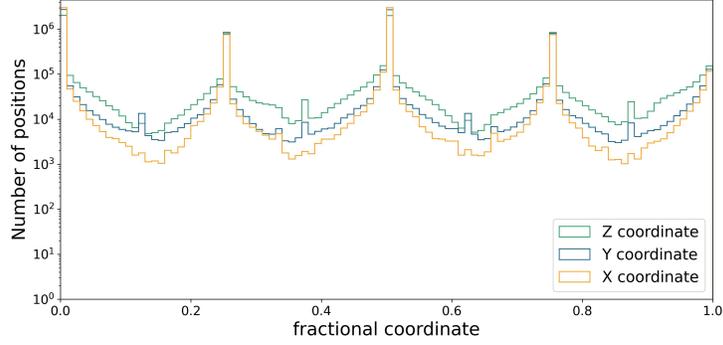


Figure 3: Distribution of the components of fractional atomic coordinates (X, Y, Z)

434 A.2 Pseudocode

Algorithm 1 Training Regression Modification Model

- 1: repeat
 - 2: $x_0 \sim q(x_0); x_1 \sim q(x_1); el \sim q(el); spg_1 \sim q(spg_1); E \sim q(E)$
 - 3: $\mathcal{L} \leftarrow \|x_1 - f_\theta(x_0, t, el, spg_1, E)\|$
 - 4: $\theta \leftarrow Update(\theta, \nabla_\theta \mathcal{L}(\theta))$
 - 5: until converge
-

Algorithm 2 Inferencing Regression Modification Model

- 1: $x_0 \sim q(x_0); el \sim q(el); spg_1 \sim q(spg_1); E \sim q(E)$
 - 2: $x_1 = f_\theta(x_0, t, el, spg_1, E)$
 - 3: return x_1
-

Algorithm 3 CFM Training

- 1: repeat
 - 2: $x_0 \sim q(x_0); x_1 \sim q(x_1)$
 - 3: $t \sim \mathcal{U}(0, 1)$
 - 4: $x_t = tx_1 + (1 - t)x_0$
 - 5: $\mathcal{L}_{CFM} \leftarrow \|f_\theta(x_t, t) - (x_1 - x_0)\|$
 - 6: $\theta \leftarrow Update(\theta, \nabla_\theta \mathcal{L}_{CFM}(\theta))$
 - 7: until converge
-

Algorithm 4 Training CFM for Modification

- 1: repeat
 - 2: $x_0 \sim q(x_0); x_1 \sim q(x_1); el \sim q(el); spg_1 \sim q(spg_1); E \sim q(E)$
 - 3: $t \sim \mathcal{U}(0, 1)$
 - 4: $x_t = tx_1 + (1 - t)x_0$
 - 5: $\mathcal{L}_{CFM} \leftarrow \|f_\theta(x_t, t, el, spg_1, E) - (x_1 - x_0)\|$
 - 6: $\theta \leftarrow Update(\theta, \nabla_\theta \mathcal{L}_{CFM}(\theta))$
 - 7: until converge
-

Algorithm 5 Training CFM for Generation

- 1: repeat
- 2: $x_0 \sim \mathcal{N}(0, 1)$ or $x_0 \sim \mathcal{U}(0, 1)$
- 3: $x_1 \sim q(x_1); el \sim q(el); spg_1 \sim q(spg_1); E \sim q(E)$
- 4: $t \sim \mathcal{U}(0, 1)$
- 5: $x_t = tx_1 + (1 - t)x_0$
- 6: $\mathcal{L}_{CFM} \leftarrow \|f_\theta(x_t, t, el, spg_1, E) - (x_1 - x_0)\|$
- 7: $\theta \leftarrow Update(\theta, \nabla_\theta \mathcal{L}_{CFM}(\theta))$
- 8: until converge

Algorithm 6 Sampling with CFM for Modification or Generation

- 1: $h = \frac{1}{T}$
- 2: $x_0 \sim q(x_0)$ or $x_0 \sim \mathcal{N}(0, 1)$ or $x_0 \sim \mathcal{U}(0, 1)$
- 3: $el \sim q(el); spg_1 \sim q(spg_1); E \sim q(E)$
- 4: for $dot = 1, \dots, T$ do
- 5: $x_{t+1} = x_t + hf_\theta(x_t, t, el, spg_1, E)$
- 6: end for
- 7: return x_1

Algorithm 7 Training DM for Modification

- 1: repeat
- 2: $x_0 \sim q(x_0); x_1 \sim q(x_1); el \sim q(el); spg_1 \sim q(spg_1); E \sim q(E)$
- 3: $t \sim \mathcal{U}(\{1, \dots, T\})$
- 4: $\epsilon \sim \mathcal{N}(0, I)$
- 5: $\mathcal{L}_D \leftarrow \|\epsilon - f_\theta(\sqrt{\alpha_t}x_1 + \sqrt{1 - \alpha_t}\epsilon, x_0, t, el, spg_1, E)\|$
- 6: $\theta \leftarrow Update(\theta, \nabla_\theta \mathcal{L}_D(\theta))$
- 7: until converge

Algorithm 8 Sampling with DM for Modification

- 1: $x_T \sim \mathcal{N}(0, I)$
- 2: for $dot = T, \dots, 1$ do
- 3: $x_0 \sim q(x_0); x_1 \sim q(x_1); el \sim q(el); spg_1 \sim q(spg_1); E \sim q(E)$
- 4: $z \sim \mathcal{N}(0, I)$ if $t > 1$ else $z = 0$
- 5: $x_{t-1} = \frac{1}{\sqrt{\alpha_t}}(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}}f_\theta(x_t, x_0, t, el, spg_1, E)) + \sqrt{1 - \alpha_t}z$
- 6: end for
- 7: return x_1

Algorithm 9 Training DM for Generation

- 1: repeat
- 2: $x_1 \sim q(x_1); el \sim q(el); spg_1 \sim q(spg_1); E \sim q(E)$
- 3: $t \sim \mathcal{U}(\{1, \dots, T\})$
- 4: $\epsilon \sim \mathcal{N}(0, I)$
- 5: $\mathcal{L}_D \leftarrow \|\epsilon - f_\theta(\sqrt{\alpha_t}x_1 + \sqrt{1 - \alpha_t}\epsilon, t, el, spg_1, E)\|$
- 6: $\theta \leftarrow Update(\theta, \nabla_\theta \mathcal{L}_D(\theta))$
- 7: until converge

Algorithm 10 DDPM Sampling

- 1: $x_T \sim \mathcal{N}(0, I)$
- 2: for $dot = T, \dots, 1$ do
- 3: $x_1 \sim q(x_1); el \sim q(el); spg_1 \sim q(spg_1); E \sim q(E)$
- 4: $z \sim \mathcal{N}(0, I)$ if $t > 1$ else $z = 0$
- 5: $x_{t-1} = \frac{1}{\sqrt{\alpha_t}}(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}}f_\theta(x_t, t, el, spg_1, E)) + \sqrt{1 - \alpha_t}z$
- 6: end for
- 7: return x_1

Algorithm 11 DDIM Sampling

```
1:  $x_T \sim \mathcal{N}(0, I)$ 
2: for dot =  $T, \dots, 1$  with step C do
3:    $x_0 \sim q(x_0); x_1 \sim q(x_1); el \sim q(el); spg_1 \sim q(spg_1); E \sim q(E)$ 
4:    $z \sim \mathcal{N}(0, I)$  if  $t > 1$  else  $z = 0$ 
5:    $x_\theta = f_\theta(x_t, t, el, spg_1, E)$ 
6:    $x_{t-1} = \sqrt{\alpha_{t-1}} \left( \frac{x_t - \sqrt{1 - \alpha_t} x_\theta}{\sqrt{\alpha_t}} \right) + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} x_\theta + \sigma_t z$ 
7:
8: end for
9: return  $x_1$ 
```

435 **A.3 Experiment Details**

436 All the experiments use the same hyperparameters for the model:

- 437 • num_res_blocks = 7
- 438 • attention_resolution = (1, 2, 4, 8)
- 439 • model_channels = 128

440 In all the experiments models are trained with the same training parameters:

- 441 • optimizer = Adam
 - 442 - betas = (0.9, 0.999)
 - 443 - eps = 1e-08
 - 444 - weight_decay = 0
- 445 • batch_size = 256
- 446 • epochs = 400
- 447 • learning_rate = 1e-4
- 448 • lr_warmup_steps = 500
- 449 • random_state = 42

450 An important note, that all our experiments have been conducted in mixed precision in fp16.

451 Generation task: Diffusion Model (DDPM):

- 452 • num_train_timesteps = 1000 (diffusion process discretization)
- 453 • beta_start = 0.0001
- 454 • beta_end = 0.02
- 455 • num_inference_steps = 100
- 456 • beta_schedule = "squaredcos_cap_v2" (cosine)

457 Diffusion Model (DDIM):

- 458 • num_train_timesteps = 1000 (diffusion process discretization)
- 459 • beta_start = 0.0001
- 460 • beta_end = 0.02
- 461 • num_inference_steps = 100
- 462 • beta_schedule = "squaredcos_cap_v2" (cosine)

463 Flow Matching $x_0 \sim \mathcal{N}(0, 1)$:

- 464 • num_inference_steps = 100

465 Flow Matching $x_0 \sim \mathcal{U}(0, 1)$:

466 • num_inference_steps = 100

467 Modification task:

468 Regression UNet:

469 • num_inference_steps = 1

470 Diffusion Model:

471 • num_train_timesteps = 1000 (diffusion process discretization)

472 • beta_start = 0.0001

473 • beta_end = 0.02

474 • num_inference_steps = 100

475 • beta_schedule = "squaredcos_cap_v2" (cosine)

476 Flow Matching:

477 • num_inference_steps = 100

478 A.4 PBC Loss details

479 The PBC loss function operates through several steps:

1. Vertices evaluation: If the target coordinate of the atom is lattice vertex (all 3 coordinates x, y, z are equal to 1 or 0), then loss between prediction point $preds_i$ and target point $target_i$ is being calculated using following formula:

$$L_{vertex}(preds_i, target_i) = \min_{v \in vertices} ||preds_i - v||,$$

480 where *vertices* is a set of 8 possible positions according to PBC ($\{0, 0, 0\}, \{0, 0, 1\},$
481 $\dots, \{1, 1, 1\}$).

2. Edges evaluation: If the target coordinate of the atom is located on lattice edge (two coordinates are equal to 1 or 0 and one is not). For example, a lattice edge atom at point $\{0, 1, 0.3\}$ has identical atoms at points $\{0, 0, 0.3\}, \{1, 0, 0.3\}, \{1, 1, 0.3\}$. As we can see, in this example z -coordinate is fixed but x and y are exchangeable. Therefore, if the target point is represented as $\{x, y, z\}$, we can use the following formula:

$$L_{edge}(preds_i, target_i) = \min_{e \in edgePoints} ||preds_i - e||,$$

488 where *edgePoints* is a set of 4 possible positions according to PBC.

- 489 • Case of fixed point x : $edgePoints = \{\{x, 0, 0\}, \{x, 0, 1\}, \{x, 1, 0\}, \{x, 1, 1\}\}$
- 490 • Case of fixed point y : $edgePoints = \{\{0, y, 0\}, \{0, y, 1\}, \{1, y, 0\}, \{1, y, 1\}\}$
- 491 • Case of fixed point z : $edgePoints = \{\{0, 0, z\}, \{0, 1, z\}, \{1, 0, z\}, \{1, 1, z\}\}$

3. Sides evaluation: If the target coordinate of the atom is located on lattice side (one coordinate is equal to 1 or 0 and two are not). For example, lattice side atom at point $\{0, 0.5, 0.3\}$ has identical atom at point $\{1, 0.5, 0.3\}$. In this example y and z coordinates are fixed but x is exchangeable. Therefore, if the target point is represented as $\{x, y, z\}$, we can use the following formula:

$$L_{size}(preds_i, target_i) = \min_{s \in sidePoints} ||preds_i - s||,$$

497 where *sidePoints* is a set of 2 possible positions according to PBC.

- 498 • Case of exchangeable point x : $sidePoints = \{\{0, y, z\}, \{1, y, z\}\}$
- 499 • Case of exchangeable point y : $sidePoints = \{\{x, 0, z\}, \{x, 1, z\}\}$
- 500 • Case of exchangeable point z : $sidePoints = \{\{x, y, 0\}, \{x, y, 1\}\}$



Figure 4: Example of using PBC-aware loss. The depicted structures ($\text{Mo}_2\text{Nb}_2\text{Ta}_2\text{W}_2$) are visually different, but in fact they are exact the same. It is confirmed by insignificant value of PBC-aware loss

501 4. Points, which don't belong to the groups above, are processed using the default loss
 502 function.

Since the $\min(x_1, x_2, \dots, x_n)$ function is undifferentiable at multiple points ($x_i = x_j \forall i \neq j$), it makes a loss function to have a more complicated surface. Therefore, we used a norm function with order $k \rightarrow -\infty$ which is differentiable at all points as a replacement.

$$\min_{diff}(x_1, x_2, \dots, x_n) = \left(\sum_{i=1}^n |x_i|^k \right)^{\frac{1}{k}} \quad k \rightarrow -\infty$$

503 Therefore, overall PBC-aware loss for a structure is represented as:

$$\begin{aligned} L_{PBC}(preds, target) = & \sum_{i=1}^n \mathbb{I}(target_i \text{ is vertex point}) L_{vertex}(preds_i, target_i) \\ & + \mathbb{I}(target_i \text{ is edge point}) L_{edge}(preds_i, target_i) \\ & + \mathbb{I}(target_i \text{ is side point}) L_{side}(preds_i, target_i) \\ & + \mathbb{I}(target_i \text{ is usual point}) L_2(preds_i, target_i) \end{aligned}$$

504 As the count of atoms varies across different structures, the L_{PBC} metric tends to yield
 505 higher values for structures featuring a larger number of atoms. Thus, it is important to
 506 normalize the loss function with the number of atoms in the structure if it would be used in
 507 batches with structures with different number of atoms. Therefore, a PBC-aware loss for a
 508 batch of structures is formulated as:

$$\begin{aligned} & L_{batchPBC}(batchPreds, batchTargets) = \\ & \sum_{i=1}^{batchSize} \frac{1}{numSites_i} L_{PBC}(batchPreds_i, batchTargets_i) \end{aligned}$$

509 Compute resources

510 For our computational needs in model training and inference, we deployed a total of three
 511 GPU servers with the following configurations:

512 Server 1:

- 513 • GPU: NVIDIA A100/80G
- 514 • CPU: 8vCPU of Intel(R) Xeon(R) Gold 6248R @ 3.00 GHz
- 515 • RAM: 64Gb

516 Server 2:

- 517 • GPU: NVIDIA V100 (32GB)
- 518 • CPU: 8vCPU of Intel(R) Xeon(R) Gold 6278C @ 2.60 GHz
- 519 • RAM: 64Gb

520 Server 3:

- 521 • GPU: NVIDIA V100 (32GB)
- 522 • CPU: 8vCPU of Intel(R) Xeon(R) Gold 6278C @ 2.60 GHz
- 523 • RAM: 64Gb

524 Every model training time consumed up to 2 weeks employing computing power of one GPU
525 server.

526 For the ab-initio calculations implemented in VASP, we deployed a total of 5 identical CPU
527 servers with the following configurations:

- 528 • CPU: 64vCPU of Intel(R) Xeon(R) Gold 6278C CPU @ 2.60GHz
- 529 • RAM: 256Gb

530 Structure relaxation with VASP for all six experiments mentioned in Table 3 took more than
531 180 thousand CPU hours. Computing power of all CPU servers was employed.

532 NeurIPS Paper Checklist

533 The checklist is designed to encourage best practices for responsible machine learning research,
534 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do
535 not remove the checklist: The papers not including the checklist will be desk rejected. The
536 checklist should follow the references and follow the (optional) supplemental material. The
537 checklist does NOT count towards the page limit.

538 Please read the checklist guidelines carefully for information on how to answer these questions.
539 For each question in the checklist:

- 540 • You should answer [Yes] , [No] , or [NA] .
- 541 • [NA] means either that the question is Not Applicable for that particular paper or
542 the relevant information is Not Available.
- 543 • Please provide a short (1–2 sentence) justification right after your answer (even for
544 NA).

545 The checklist answers are an integral part of your paper submission. They are visible to the
546 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also
547 include it (after eventual revisions) with the final version of your paper, and its final version
548 will be published with the paper.

549 The reviewers of your paper will be asked to use the checklist as one of the factors in their
550 evaluation. While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to
551 answer "[No]" provided a proper justification is given (e.g., "error bars are not reported
552 because it would be too computationally expensive" or "we were unable to find the license for
553 the dataset we used"). In general, answering "[No]" or "[NA]" is not grounds for rejection.
554 While the questions are phrased in a binary way, we acknowledge that the true answer
555 is often more nuanced, so please just use your best judgment and write a justification to
556 elaborate. All supporting evidence can appear either in the main paper or the supplemental
557 material, provided in appendix. If you answer [Yes] to a question, in the justification please
558 point to the section(s) where related material for the question can be found.

559 IMPORTANT, please:

- 560 • Delete this instruction block, but keep the section heading "NeurIPS paper checklist",
- 561 • Keep the checklist subsection headings, questions/answers and guidelines below.
- 562 • Do not modify the questions and only use the provided macros for your answers.

563 1. Claims

564 Question: Do the main claims made in the abstract and introduction accurately
565 reflect the paper's contributions and scope?

566 Answer: [Yes] ,

567 Justification: The main claims made in the abstract and introduction do accurately
568 reflect the paper's contributions and scope. Every aspects mentioned in the abstract
569 and introduction are further revealed in the main paper.

570 Guidelines:

- 571 • The answer NA means that the abstract and introduction do not include the
572 claims made in the paper.
- 573 • The abstract and/or introduction should clearly state the claims made, including
574 the contributions made in the paper and important assumptions and limitations.
575 A No or NA answer to this question will not be perceived well by the reviewers.
- 576 • The claims made should match theoretical and experimental results, and reflect
577 how much the results can be expected to generalize to other settings.
- 578 • It is fine to include aspirational goals as motivation as long as it is clear that
579 these goals are not attained by the paper.

580 2. Limitations

581 Question: Does the paper discuss the limitations of the work performed by the
582 authors?

583 Answer: [Yes]

584 Justification: All the limitations are discussed in the conclusion section. More
585 specifically, limitations are 1) in number of atoms per unit cell of crystal structure
586 that model can work with and 2) symmetries that our structure representation is
587 able to encode. For example, crystals have infinite periodic structure. There is no
588 possible way to

589 Guidelines:

- 590 • The answer NA means that the paper has no limitation while the answer No
591 means that the paper has limitations, but those are not discussed in the paper.
- 592 • The authors are encouraged to create a separate "Limitations" section in their
593 paper.
- 594 • The paper should point out any strong assumptions and how robust the results
595 are to violations of these assumptions (e.g., independence assumptions, noiseless
596 settings, model well-specification, asymptotic approximations only holding
597 locally). The authors should reflect on how these assumptions might be violated
598 in practice and what the implications would be.
- 599 • The authors should reflect on the scope of the claims made, e.g., if the approach
600 was only tested on a few datasets or with a few runs. In general, empirical
601 results often depend on implicit assumptions, which should be articulated.
- 602 • The authors should reflect on the factors that influence the performance of the
603 approach. For example, a facial recognition algorithm may perform poorly when
604 image resolution is low or images are taken in low lighting. Or a speech-to-text
605 system might not be used reliably to provide closed captions for online lectures
606 because it fails to handle technical jargon.
- 607 • The authors should discuss the computational efficiency of the proposed algo-
608 rithms and how they scale with dataset size.
- 609 • If applicable, the authors should discuss possible limitations of their approach
610 to address problems of privacy and fairness.
- 611 • While the authors might fear that complete honesty about limitations might
612 be used by reviewers as grounds for rejection, a worse outcome might be that
613 reviewers discover limitations that aren't acknowledged in the paper. The
614 authors should use their best judgment and recognize that individual actions in
615 favor of transparency play an important role in developing norms that preserve
616 the integrity of the community. Reviewers will be specifically instructed to not
617 penalize honesty concerning limitations.

618 3. Theory Assumptions and Proofs

619 Question: For each theoretical result, does the paper provide the full set of assump-
620 tions and a complete (and correct) proof?

621 Answer: [NA] .

622 Justification: Our work does not include significant theoretical results due to the
623 fact that it is mainly focused on experiments.

624 Guidelines:

- 625 • The answer NA means that the paper does not include theoretical results.
- 626 • All the theorems, formulas, and proofs in the paper should be numbered and
627 cross-referenced.
- 628 • All assumptions should be clearly stated or referenced in the statement of any
629 theorems.
- 630 • The proofs can either appear in the main paper or the supplemental material,
631 but if they appear in the supplemental material, the authors are encouraged to
632 provide a short proof sketch to provide intuition.
- 633 • Inversely, any informal proof provided in the core of the paper should be
634 complemented by formal proofs provided in appendix or supplemental material.

635 • Theorems and Lemmas that the proof relies upon should be properly referenced.

636 4. Experimental Result Reproducibility

637 Question: Does the paper fully disclose all the information needed to reproduce
638 the main experimental results of the paper to the extent that it affects the main
639 claims and/or conclusions of the paper (regardless of whether the code and data are
640 provided or not)?

641 Answer: [Yes]

642 Justification: The paper fully discloses all the information needed to reproduce the
643 main experimental results. The information includes: 1) all hyperparameters for
644 the models, including the random seed they were performed on 2) crystal structures
645 used as a data, 3) VASP settings.

646 Guidelines:

- 647 • The answer NA means that the paper does not include experiments.
- 648 • If the paper includes experiments, a No answer to this question will not be
649 perceived well by the reviewers: Making the paper reproducible is important,
650 regardless of whether the code and data are provided or not.
- 651 • If the contribution is a dataset and/or model, the authors should describe the
652 steps taken to make their results reproducible or verifiable.
- 653 • Depending on the contribution, reproducibility can be accomplished in various
654 ways. For example, if the contribution is a novel architecture, describing the
655 architecture fully might suffice, or if the contribution is a specific model and
656 empirical evaluation, it may be necessary to either make it possible for others to
657 replicate the model with the same dataset, or provide access to the model. In
658 general, releasing code and data is often one good way to accomplish this, but
659 reproducibility can also be provided via detailed instructions for how to replicate
660 the results, access to a hosted model (e.g., in the case of a large language model),
661 releasing of a model checkpoint, or other means that are appropriate to the
662 research performed.
- 663 • While NeurIPS does not require releasing code, the conference does require all
664 submissions to provide some reasonable avenue for reproducibility, which may
665 depend on the nature of the contribution. For example
 - 666 (a) If the contribution is primarily a new algorithm, the paper should make it
667 clear how to reproduce that algorithm.
 - 668 (b) If the contribution is primarily a new model architecture, the paper should
669 describe the architecture clearly and fully.
 - 670 (c) If the contribution is a new model (e.g., a large language model), then there
671 should either be a way to access this model for reproducing the results or a
672 way to reproduce the model (e.g., with an open-source dataset or instructions
673 for how to construct the dataset).
 - 674 (d) We recognize that reproducibility may be tricky in some cases, in which
675 case authors are welcome to describe the particular way they provide for
676 reproducibility. In the case of closed-source models, it may be that access to
677 the model is limited in some way (e.g., to registered users), but it should be
678 possible for other researchers to have some path to reproducing or verifying
679 the results.

680 5. Open access to data and code

681 Question: Does the paper provide open access to the data and code, with sufficient
682 instructions to faithfully reproduce the main experimental results, as described in
683 supplemental material?

684 Answer: [Yes]

685 Justification: The raw crystal dataset can be downloaded from <https://afloplib.org>.
686 The source code for training and inferencing our models can be obtained from
687 GitHub at <https://github.com/AIRI-Institute/conditional-crystal-generation>. Both
688 data and code are also referenced in the main paper's Section 9 and Section 9 .

689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [\[Yes\]](#)

Justification: All the training and testing details and all the hyperparameters for the experiments are mentioned in the paper.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [\[No\]](#)

Justification: Firstly, calculating statistical significance for all our experiments is computationally expensive. Secondly, experimental results (crystal structures obtained from model inference) have been approved by further ab-initio calculations implemented in VASP.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)

- 742 • The assumptions made should be given (e.g., Normally distributed errors).
- 743 • It should be clear whether the error bar is the standard deviation or the standard
- 744 error of the mean.
- 745 • It is OK to report 1-sigma error bars, but one should state it. The authors
- 746 should preferably report a 2-sigma error bar than state that they have a 96%
- 747 CI, if the hypothesis of Normality of errors is not verified.
- 748 • For asymmetric distributions, the authors should be careful not to show in
- 749 tables or figures symmetric error bars that would yield results that are out of
- 750 range (e.g. negative error rates).
- 751 • If error bars are reported in tables or plots, The authors should explain in the
- 752 text how they were calculated and reference the corresponding figures or tables
- 753 in the text.

754 8. Experiments Compute Resources

755 Question: For each experiment, does the paper provide sufficient information on the
 756 computer resources (type of compute workers, memory, time of execution) needed
 757 to reproduce the experiments?

758 Answer: [\[Yes\]](#)

759 Justification: In the paper, we state all the compute powers that we have employed
 760 for AI model training and ab-initio calculations. One can observe it in Section A.4.

761 Guidelines:

- 762 • The answer NA means that the paper does not include experiments.
- 763 • The paper should indicate the type of compute workers CPU or GPU, internal
- 764 cluster, or cloud provider, including relevant memory and storage.
- 765 • The paper should provide the amount of compute required for each of the
- 766 individual experimental runs as well as estimate the total compute.
- 767 • The paper should disclose whether the full research project required more
- 768 compute than the experiments reported in the paper (e.g., preliminary or failed
- 769 experiments that didn't make it into the paper).

770 9. Code Of Ethics

771 Question: Does the research conducted in the paper conform, in every respect, with
 772 the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

773 Answer: [\[Yes\]](#)

774 Justification: Our research, working process, data and code correspond with NeurIPS
 775 Code of Ethics in <https://neurips.cc/public/EthicsGuidelines>

776 Guidelines:

- 777 • The answer NA means that the authors have not reviewed the NeurIPS Code
- 778 of Ethics.
- 779 • If the authors answer No, they should explain the special circumstances that
- 780 require a deviation from the Code of Ethics.
- 781 • The authors should make sure to preserve anonymity (e.g., if there is a special
- 782 consideration due to laws or regulations in their jurisdiction).

783 10. Broader Impacts

784 Question: Does the paper discuss both potential positive societal impacts and
 785 negative societal impacts of the work performed?

786 Answer: [\[Yes\]](#)

787 Justification: The broad societal impacts are discussed int the conclusion section.

788 Guidelines:

- 789 • The answer NA means that there is no societal impact of the work performed.
- 790 • If the authors answer NA or No, they should explain why their work has no
- 791 societal impact or why the paper does not address societal impact.

- 792
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

813 11. Safeguards

814 Question: Does the paper describe safeguards that have been put in place for
815 responsible release of data or models that have a high risk for misuse (e.g., pretrained
816 language models, image generators, or scraped datasets)?

817 Answer: [NA]

818 Justification: The paper poses no such risks.

819 Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

830 12. Licenses for existing assets

831 Question: Are the creators or original owners of assets (e.g., code, data, models),
832 used in the paper, properly credited and are the license and terms of use explicitly
833 mentioned and properly respected?

834 Answer: [Yes]

835 Justification: The models utilized in this research are appropriately cited within
836 the text, including references to the AFLOW database. We have no commercial
837 interests related to the use of these models and the crystal structure database. The
838 code used in this study was entirely developed by our team.

839 Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.

- 845 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 846 • For scraped data from a particular source (e.g., website), the copyright and
- 847 terms of service of that source should be provided.
- 848 • If assets are released, the license, copyright information, and terms of use in
- 849 the package should be provided. For popular datasets, [paperswithcode.com/](https://paperswithcode.com/datasets)
- 850 [datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help
- 851 determine the license of a dataset.
- 852 • For existing datasets that are re-packaged, both the original license and the
- 853 license of the derived asset (if it has changed) should be provided.
- 854 • If this information is not available online, the authors are encouraged to reach
- 855 out to the asset's creators.

856 13. New Assets

857 Question: Are new assets introduced in the paper well documented and is the

858 documentation provided alongside the assets?

859 Answer: [\[Yes\]](#)

860 Justification: The primary assets of our research are: 1) the methodology, and 2) the

861 code for training models and performing inference. Both are thoroughly documented.

862 Guidelines:

- 863 • The answer NA means that the paper does not release new assets.
- 864 • Researchers should communicate the details of the dataset/code/model as part
- 865 of their submissions via structured templates. This includes details about
- 866 training, license, limitations, etc.
- 867 • The paper should discuss whether and how consent was obtained from people
- 868 whose asset is used.
- 869 • At submission time, remember to anonymize your assets (if applicable). You
- 870 can either create an anonymized URL or include an anonymized zip file.

871 14. Crowdsourcing and Research with Human Subjects

872 Question: For crowdsourcing experiments and research with human subjects, does

873 the paper include the full text of instructions given to participants and screenshots,

874 if applicable, as well as details about compensation (if any)?

875 Answer: [\[NA\]](#)

876 Justification: The paper does not involve crowdsourcing nor research with human

877 subjects.

878 Guidelines:

- 879 • The answer NA means that the paper does not involve crowdsourcing nor
- 880 research with human subjects.
- 881 • Including this information in the supplemental material is fine, but if the main
- 882 contribution of the paper involves human subjects, then as much detail as
- 883 possible should be included in the main paper.
- 884 • According to the NeurIPS Code of Ethics, workers involved in data collection,
- 885 curation, or other labor should be paid at least the minimum wage in the
- 886 country of the data collector.

887 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human

888 Subjects

889 Question: Does the paper describe potential risks incurred by study participants,

890 whether such risks were disclosed to the subjects, and whether Institutional Review

891 Board (IRB) approvals (or an equivalent approval/review based on the requirements

892 of your country or institution) were obtained?

893 Answer: [\[NA\]](#)

894 Justification: The paper does not involve crowdsourcing nor research with human

895 subjects

896 Guidelines:

897
898
899
900
901
902
903
904
905
906

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.