

---

# Navigating Order-(Dis)Order Family Trees via Group-Subgroup Transitions

---

Anonymous Authors<sup>1</sup>

## Abstract

Novelty in materials discovery is conventionally assessed against databases of ordered crystal structures, yet a predicted ordered structure may simply be a particular ordering of a known disordered phase, rendering its apparent novelty misleading. We refer to such cases as ordered children of disordered parents. To address this limitation, we introduce order-(dis)order family trees, a symmetry-based framework that organizes ordered and disordered structures through group-subgroup relations and enables novelty to be explicitly evaluated. Using a high-throughput family-matching procedure, we identify candidate disordered parents and symmetry-related ordered relatives for a given ordered structure. We test our framework using A-Lab results, where several targeted ordered structures are correctly identified as ordered children of known disordered parents. Extending the benchmark across experimental databases (ICSD), simulation datasets (MP-20, Alex-MP-20, GNoME), and crystal generative models reveals that many seemingly novel ordered structures are, in fact, better understood as members of experimentally known order-(dis)order family trees, establishing family trees as a key requirement for achieving genuine novelty in data-driven materials discovery.

## 1. Introduction

Materials discovery systems are pushing the field toward a closed loop between computational prediction and experimental synthesis. Agentic workflows now produce millions of candidate compounds at a rate that far outpaces experimental characterization, as they scale from proposing structures to directing robotic synthesis (Szymanski et al., 2023; Merchant et al., 2023). Consequently, the central question

is no longer simply whether we can predict stable structures beyond the training distribution, but whether the structures we predict are genuinely new to synthesis (Cheetham & Seshadri, 2024; Leeman et al., 2024). In current practice, novelty is typically defined by absence from reference databases of ordered crystal structures (Xie et al., 2021; Zeni et al., 2025; Betala et al., 2025). Defining novelty in these systems has become a critical concern: a system that rediscovers known phases dressed as new compounds wastes synthesis cycles and corrupts its own reward.

Yet, answering whether a predicted structure is genuinely novel is more subtle than it first appears. Novelty is not an intrinsic property of a compound, but one defined relative to a reference: a structure is novel only with respect to the space against which it is compared. Although ordered structure databases are a natural reference given the data on which many discovery systems are trained, relying on them alone introduces a systematic blind spot. Experimentally synthesized materials frequently exhibit occupational disorder, in which the same crystallographic lattice sites are occupied by multiple species in the statistical average structure (Rühl, 2019). In this setting, a predicted ordered structure may not represent a genuinely new compound, but rather a specific ordering of a known disordered phase, which we term *ordered children of disordered parents* (Cheetham & Seshadri, 2024). This discrepancy is epitomized by the A-Lab study (Szymanski et al., 2023), where subsequent analysis suggested that over 65% of the claimed discoveries were more plausibly interpreted as existing disordered parents rather than the predicted ordered children (Leeman et al., 2024). When disordered phases are excluded from the reference space, these ordered children are routinely labeled as novel, leading to a systematic overestimation of discovery rates relative to experimental reality.

These cases point to a structural gap in how computational predictions are interpreted. Ordered and disordered phases are often not unrelated alternatives, but members of the same crystallographic lineage (Bärnighausen, 1980). This connection arises because atomic ordering directly governs crystal symmetry: the onset of ordering breaks the full symmetry of the underlying disordered lattice, producing a lower-symmetry structure. These symmetry reductions are con-

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

strained by **group-subgroup relationships** (Bärnighausen, 1980; Ivantchev et al., 2000; Han et al., 2025), governing the ordered structures that can descend from a given disordered parent. Multiple ordered structures can descend from the same parent through distinct symmetry-breaking pathways, forming a family of ordered siblings. Such relations are rarely available in a systematic, high-throughput form, and most reported examples have been uncovered only through detailed manual analysis of a small number of predicted structures subjected to experimental validation (Cheetham & Seshadri, 2024; Leeman et al., 2024; Juelscholt, 2026). In principle, one could attempt exhaustive mapping by enumerating all symmetry-inequivalent ordered configurations of known disordered phases, for example using methods such as POCC (Yang et al., 2016), and then matching a query ordered structure against the resulting lookup space. In practice, however, extending such enumeration across the full set of known disordered compounds is computationally prohibitive and unsuitable for large-scale analysis.

To address this gap, we develop a symmetry-based framework that organizes ordered and disordered phases through group-subgroup transitions in space-group hierarchy: **order-(dis)order family trees**. Rather than treating each predicted structure as an isolated point, we view it as an element embedded in a manifold of related phases. This perspective allows us to distinguish between genuinely new families and symmetry-lowered rediscoveries of known phases.

Our contributions are as follows:

- **Order-(dis)order Family Trees:** We formalize a symmetry-based framework that traces crystallographic lineage through group-subgroup transitions, providing a natural language for linking ordered children to disordered parents.
- **High-throughput Family Matching Procedure:** We develop a fast and scalable procedure to identify candidate disordered parents and symmetry-related structures for any ordered query structure.
- **Disorder-aware Novelty Metrics:** We define  $FT_{\text{disorder}}$  and  $FT_{\text{order}}$ , two new Family Tree-based metrics that evaluate the fraction of compounds connected through family trees to experimentally known phases, providing a view of novelty that better reflects experimental reality.
- **Large-scale Benchmark:** We apply these metrics across experimental and simulated structure databases (ICSD, MP-20, Alex-MP-20, GNoME) and 10+ state-of-the-art crystal generative models. Our analysis reveals that a substantial fraction of apparently novel structures belong to existing order-(dis)order family trees, particularly in symmetry-agnostic all-atom models, which exhibit this pattern 2–4× more often than symmetry-constrained models.

## 2. Methods

Systematic navigation of order–disorder family trees requires a crystal representation that is both physically meaningful and computationally tractable. Our framework builds on two components: (i) the symmetry hierarchy of crystals expressed through space groups, Wyckoff positions, and group-subgroup transitions (Bärnighausen, 1980; Wyckoff, 1922); and (ii) a unified crystal representation describing both ordered and disordered structures within the same formal language. Together, these enable family-level reasoning where ordered children, disordered parents, and symmetry-related siblings are analyzed within a common framework.

### 2.1. Preliminaries

**Space groups and Wyckoff positions** The space group  $G$  constitutes an important aspect of a crystal structure, representing the set of symmetry operations that leave the lattice invariant. Each operation  $g \in G$  acts on a point  $\mathbf{r} \in \mathbb{R}^3$  via an affine transformation  $\mathbf{r} \mapsto R\mathbf{r} + \mathbf{t}$ , where  $R$  denotes a rotation or rotoinversion and  $\mathbf{t}$  represents a translation. This group structure dictates the symmetry-equivalence of points within the crystal. For any specific point  $\mathbf{r}$ , the site-symmetry group is defined as the subgroup  $G_{\mathbf{r}} = \{g \in G \mid g\mathbf{r} = \mathbf{r}\}$ , containing all operations that leave  $\mathbf{r}$  fixed. Points are formally categorized into Wyckoff positions, which are sets of symmetry-equivalent sites whose site-symmetry groups are conjugate within  $G$  (Wyckoff, 1922). Each Wyckoff position is uniquely characterized by its multiplicity, site symmetry, and fractional coordinates. Ultimately, Wyckoff positions provide a mathematically rigorous and compact description of crystal structures by grouping atomic sites according to their transformation properties under  $G$ .

**Group-subgroup transitions** Crystallographic relationships between structures are established using group–subgroup relations (Bärnighausen, 1980). For two space groups  $G$  and  $H$ , the relation  $H \leq G$  denotes that  $H$  is a subgroup of  $G$ , such that all symmetry operations in  $H$  are contained within  $G$ . A transition from  $G$  to  $H$  therefore corresponds to a reduction in symmetry. The degree of symmetry reduction is quantified by the subgroup index  $[G : H] = |G|/|H|$ , which counts how many cosets of  $H$  partition  $G$ . A particularly important relation to consider is the *translationengleiche* ( $t$ -type) subgroup relations, which preserve the translation group while lowering the point symmetry. These transitions are particularly relevant for tracing ordering processes in which previously equivalent sites become symmetry-distinct. Because  $t$ -type relations preserve the translational periodicity, this site-splitting accommodates ordering within the bounds of the original primitive cell, without introducing a supercell. Although *klassengleiche* ( $k$ -type) subgroup relations, which involve a loss of

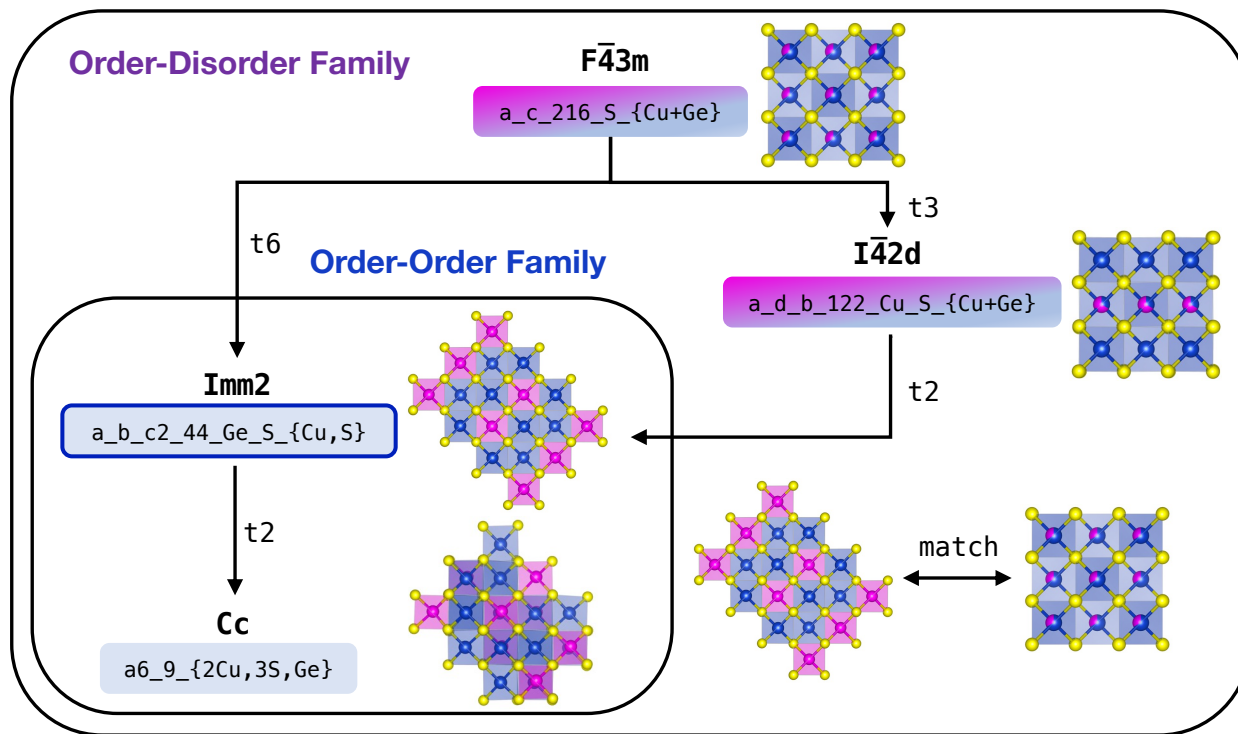


Figure 1. Order-(dis)order family tree recovered from the query structure  $\text{Cu}_2\text{GeS}_3$  ( $\text{Imm}2$ , highlighted in bright blue). Each phase is labeled by its space group and SWORD label, with parent-child links annotated by the subgroup index. The vertical arrangement follows point-group order, with higher-symmetry phases placed near the top. In this example, one of the disordered phases  $\text{a}_c\text{-}216\text{-S}_{\{\text{Cu}+\text{Ge}\}}$  occupies the highest-symmetry position and acts as the parent from which both ordered and disordered descendants are derived through group-subgroup transitions. Ordered children that share the same parent define an order-order subfamily within the broader order-(dis)order family tree. Two structures `match` if they share the same parent. Implementation details are provided in Appendix A.

translational symmetry, can also be important, they are not the main focus here. Under a group-subgroup transition, Wyckoff positions may split, causing previously equivalent sites to become distinguishable. This splitting provides the crystallographic basis for linking parent, child, and sibling structures within the same family tree.

**Disorder** Experimental structures often deviate from an ideal fully ordered configuration and may exhibit disorder in either site occupancy or atomic position (Antypov et al., 2025). A broad and practically important class is *occupational disorder*, in which the occupation of a crystallographic site is statistical rather than uniquely assigned. This includes cases where multiple chemical species share the same site with fractional occupancies, as well as vacancy-containing cases where a site is only partially occupied. Another major form is *positional disorder*, in which atoms are distributed over multiple nearby positions instead of occupying a single well-defined site in the average structure. In this work, we focus primarily on occupational disorder and do not consider positional disorder, which constitutes less than 10% of ICSD entries (Huang et al., 2026).

## 2.2. SWORD Representation

Here, we use Symmetry and Wyckoff-sequence of Ordered and Disordered crystals (SWORD) (Huang et al., 2026), a unified and scalable representation where both ordered and disordered crystal structures are encoded. This enables consistent evaluation of family relations through space group symmetry-based analysis. SWORD expresses a crystal structure through symmetry-aware site assignments and occupancy information, making it particularly well suited for tracking how ordering, partial ordering, and disorder are distributed across Wyckoff sites. A SWORD label is written in the order of Wyckoff positions, space group number, and the atomic species occupying those positions. Using the disordered parent at the root of Figure 1 as an example,  $\text{a}_c\text{-}216\text{-S}_{\{\text{Cu}+\text{Ge}\}}$  denotes a structure in space group 216 in which Wyckoff position WP  $a$  is fully occupied by Element  $S$ , while Wyckoff position WP  $c$  is compositionally disordered between Element  $\text{Cu}$  and Element  $\text{Ge}$ . Likewise, one of its ordered children,  $\text{a}_b\text{-}c2\text{-}44\text{-Ge}_a\text{-S}_b\text{-}\{\text{Cu}, S\}$ , denotes a structure in space group 44 in which WP  $a$  is occupied by Element  $\text{Ge}$ , WP  $b$  by Element  $S$ , and the two distinct crystallo-

graphic orbits of  $WP_c$  are occupied by  $Element_{Cu}$  and  $Element_S$ , respectively. This form of description is essential in our setting because the object of interest is not an isolated ordered phase, but the structural relation between an ordered derivative and its disordered or ordered relatives. By adopting a unified representation for both ordered and disordered phases, they can be compared at the level where order-disorder relationships are actually manifested.

### 2.3. Order-(Dis)Order Family Trees

Family trees are constructed to recover, for a given structure (query), the set of crystallographically related phases connected by group-subgroup transitions and occupancy splitting, along which stoichiometry may vary. Starting from the query structure, candidate parent descriptions are systematically generated, and experimentally reported parent phases are subsequently identified using our `SWORDFamilyMatcher` module.

Figure 1 illustrates the concept of order-(dis)order family mapping using the family tree recovered from an ordered  $Cu_2GeS_3$  query structure ( $I_{mm2}$ , highlighted in bright blue) as an example. Each phase in this particular family tree corresponds to an experimentally reported ICSD entry and is labeled by its space group symbol and SWORD representation, while parent-child relations are defined by group-subgroup transitions and annotated with the corresponding subgroup index. The vertical arrangement follows point-group order, so that higher-symmetry phases appear higher in the tree. In this example, the disordered phase  $a_c216_S\{Cu+Ge\}$  occupies the highest-symmetry position and acts as the parent from which both ordered and disordered descendants are reached through group-subgroup transitions. From Figure 1, it is clear that apparently distinct structures cannot be treated as isolated phases: an ordered phase may be placed within a broader family containing higher-symmetry disordered parents and symmetry-related ordered siblings. Our framework turns this intuition into a symmetry-guided mapping problem, in which a query structure is assigned to a crystallographic lineage rather than evaluated only as an isolated structure. The full connected set defines an order-(dis)order family tree, while ordered descendants sharing the same parent form order-order subfamilies. The same family tree would be recovered regardless of which phase in the tree is used as the starting query. The query therefore serves only as an entry point to the family. Additional implementation details are provided in Appendix A.

This matching framework captures several cases within a single rule: a predicted ordered structure may match an experimentally synthesized disordered parent, and two or more ordered structures may be linked through a shared parent, such that structures that are not identical at the child level

are still recognized as members of the same broader family. In this way, the mapping lifts the structure matching from the level of isolated crystals to the level of symmetry-related families. Rather than asking only whether a predicted structure has been seen before in exactly the same ordered or disordered form, we also ask whether it already resides within an existing order-(dis)order lineage.

## 3. Results

### 3.1. Validation on A-Lab

We first validate that our framework can recover disordered parents for experimentally synthesized structures. Applying it to the 35 GNoME phases reported as "successfully" synthesized by A-Lab (Szymanski et al., 2023), we find 22 have experimentally reported disordered parent phases, either recovered by our framework or manually identified by Leeman et al. (2024). Table S1 compares the ICSD Collection Codes of these previously reported disordered parents against those assigned by our framework for all 22 cases.

For 16 of the 22 (white-shaded  $\checkmark$  rows), our framework systematically recovers the same disordered parents as Leeman et al. (2024), directly validating the approach. Grey-shaded  $\square$  rows involve positional disorder, which lies outside the current scope since our method targets occupational disorder; nevertheless, even in several of these cases our framework recovers alternative substitutionally disordered parents not previously reported.

Most notably, for the three green-shaded \* phases ( $K_2TiCr(PO_4)_3$ ,  $KPr_9(Si_3O_{13})_2$ ,  $Mn_2VPO_7$ ), prior analyses by Leeman et al. (2024) suggested the possible existence of disordered parents, but did not identify any explicitly. In contrast, our framework recovers corresponding disordered parents for all three cases. Figure 2 compares the target GNoME A-Lab ordered structures with their corresponding experimentally reported disordered parent phases in ICSD for these cases. Together, these results demonstrate the framework's capacity to systematically navigate order-disorder family trees, including cases that are difficult to identify through manual inspection alone.

### 3.2. Order-(Dis)Order Family Tree Benchmark

Having validated the framework, we next use it as a benchmarking tool for family-level novelty assessment of ordered structures relative to experimental order-disorder family trees in ICSD. Rather than evaluating a structure set only through isolated child-level matches, we assess it through the family relations recovered by the proposed matching framework. To this end, we introduce two complementary Family Tree based metrics,  $FT_{disorder}$  and  $FT_{order}$ , that quantify how a given set of ordered structures is situated with respect to existing family trees.  $FT_{disorder}$  measures the

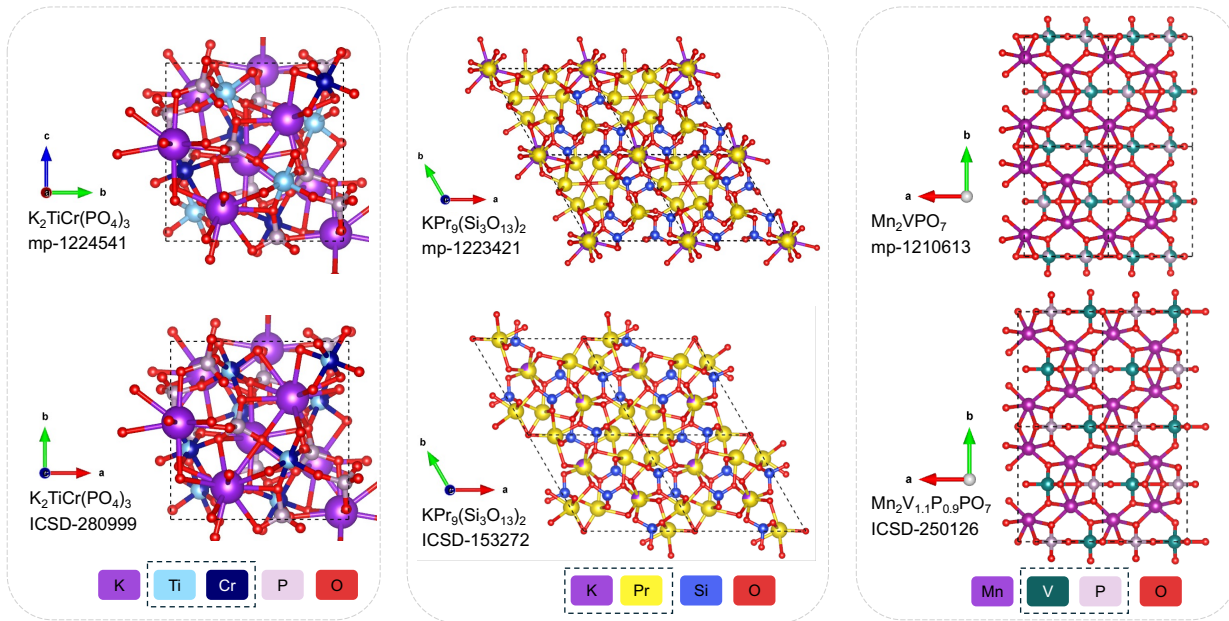


Figure 2. Recovery of experimentally reported disordered parents for three GNoME A-Lab target structures. Ordered GNoME structures are compared with the corresponding disordered parent phases reported in ICSD and recovered by our framework for  $\text{K}_2\text{TiCr}(\text{PO}_4)_3$ ,  $\text{KPr}_9(\text{Si}_3\text{O}_{13})_2$ , and  $\text{Mn}_2\text{VPO}_7$ . The colored labels at the bottom denote the elemental species, and elements enclosed by dashed outlines indicate disordered species.

fraction of ordered structures that fall within family trees rooted in existing disordered ICSD parents. A high value of  $FT_{\text{disorder}}$  indicates that many ordered structures can be traced to known disordered phases, suggesting that their apparent ordering could be artificial and that experimentally realized compounds instead favor occupational disorder. As an example, in the set of 35 GNoME A-Lab structures, 21 are found by our framework to have existing disordered parents, corresponding to  $FT_{\text{disorder}} = 60\%$ . Meanwhile,  $FT_{\text{order}}$  measures the fraction of ordered query structures that fall within family trees already spanned by existing ordered ICSD structures, i.e., those sharing the same disordered parent as their root, whether or not such a parent is experimentally synthesized. Unlike  $FT_{\text{disorder}}$ , a high  $FT_{\text{order}}$  does not imply experimental suspicion; rather, it quantifies how much of the set of ordered structures lies within the known ordered family manifold formalized by group-subgroup relations, that is, structures that can in principle be derived from already synthesized ordered phases by rearranging atoms over specific crystallographic sites within a common host structural network, with accompanying symmetry lowering or raising. Conversely, a low  $FT_{\text{order}}$  suggests genuinely new family trees absent from experimentally known crystallographic records.

**Benchmark on Databases.** We evaluate  $FT_{\text{disorder}}$  and  $FT_{\text{order}}$  across both experimental and computational structure databases, including ICSD (Zagorac et al., 2019), MP-

Table 1. Order-(dis)order family tree benchmark on known databases. For each dataset,  $FT_{\text{disorder}}$  represents the percentage of ordered structures that fall within order-(dis)order families containing an existing disordered parent, and  $FT_{\text{order}}$  represents the percentage that fall within families already spanned by existing ordered structures through group-subgroup relations.  $n_{\text{disorder}}^{FT}$  and  $n_{\text{order}}^{FT}$  are the corresponding numbers of structures assigned to these two categories, and  $n_{\text{total}}$  is the total number of ordered structures evaluated. Lower values indicate weaker overlap with experimentally known family manifolds and therefore greater family-level novelty.

Dataset	$FT_{\text{disorder}}$ (%)	$FT_{\text{order}}$ (%)	$n_{\text{disorder}}^{FT}$	$n_{\text{order}}^{FT}$	$n_{\text{total}}$
ICSD <sub>order</sub>	6.13	10.37	3564	6026	58116
MP-20 <sub>train+val</sub>	23.27	11.16	8421	4039	36183
Alex-MP-20 <sub>train+val</sub>	3.75	2.39	1876	1196	50000
GNoME <sub>stable</sub>	0.73	0.40	367	199	50000
GNoME <sub>A-Lab</sub>	60.00	0.00	21	0	35

20 (Xie et al., 2021), Alex-MP-20 (Zeni et al., 2025), and GNoME (Merchant et al., 2023). Table 1 summarizes these statistics alongside the corresponding numbers of structures assigned to each category.

Ordered structures in ICSD, the largest experimental crystal structure database to date, exhibit  $FT_{\text{disorder}} = 6.13\%$ . This reveals a nontrivial regime in which a single underlying structural network has been realized in both ordered and disordered forms, connected through group-subgroup relations. One possible explanation is that a subset of these ordered structures reflects incomplete structural refinement,

where compositional disorder was not fully resolved or reported (Spek, 2018). More plausibly, however, these cases represent genuine order-disorder transitions within a group-subgroup crystallographic lineage (Gratias & Quiquandon, 2020). In either case, the result highlights that even within the experimental record, a fraction of ordered compounds are not isolated from disorder-rooted families. A more pronounced effect is observed in MP-20 (Xie et al., 2021), a legacy Materials Project subset derived from ICSD structures with at most 20 atoms per unit cell. Here,  $FT_{\text{disorder}}$  increases sharply to 23.27%, indicating that a substantial fraction of structures coincide with experimentally realized disordered parent phases. We defer a more detailed analysis of this behavior to the Appendix C.

By contrast, computational structure databases exhibit substantially lower  $FT_{\text{disorder}}$  values. For 50k subsamples of Alex-MP-20 (Zeni et al., 2025) and GNoME stable structures on the MP hull (Merchant et al., 2023), we obtain  $FT_{\text{disorder}} = 3.75\%$  and  $0.73\%$ , respectively. This reduction should not be interpreted as evidence that computational structures are intrinsically less prone to disorder. Rather,  $FT_{\text{disorder}}$  is best understood as a lower bound on disorder propensity: it only counts cases for which a potential disordered parent phase can be found in existing experimental references. Experimental databases are concentrated in well-explored chemical systems, where both ordered and disordered phases have often been identified and recorded. Computational datasets, in contrast, extend into sparsely explored regions of composition space. In these regions, an ordered structure could still belong to a broader disorder-rooted family, but the corresponding disordered parent might simply be absent from current databases. As a result, the lower  $FT_{\text{disorder}}$  observed for computational datasets more likely reflects incomplete coverage of disorder families in experimental references than a true absence of disorder. Consistent with this view, manual inspection of GNoME structures reveals that many candidates can be associated with plausible disordered parents in ICSD (Cheetham & Seshadri, 2024); moreover, recent predictive studies also suggest that a majority of GNoME structures are in fact prone to occupational disorder (Jakob et al., 2025).

A complementary perspective emerges from  $FT_{\text{order}}$ . In ICSD, more than 10% of ordered structures belong to order-order families, indicating that a noticeable fraction of experimentally reported phases within the same composition space are related through group-subgroup transitions on a common structural network. MP-20 shows a similar level of overlap, whereas Alex-MP-20 and GNoME exhibit lower fractions, potentially because their broader compositional sampling reduces the likelihood that multiple ordered members of the same family are simultaneously represented. At the same time, the lower values of  $FT_{\text{order}}$ , which explicitly account for chemical identity, may indicate that although a

Table 2. Order-(dis)order family tree benchmark on 10,000 samples from crystal generative models trained on MP-20 and Alex-MP-20. For each model,  $FT_{\text{disorder}}$  represents the percentage of ordered structures that fall within order-(dis)order families containing an existing disordered parent, and  $FT_{\text{order}}$  represents the percentage that fall within families already spanned by existing ordered structures through group-subgroup relations. Lower values indicate weaker overlap with experimentally known family manifolds and therefore greater family-tree-based novelty.

Dataset	Model	$FT_{\text{disorder}} \downarrow (\%)$	$FT_{\text{order}} \downarrow (\%)$
MP-20	DiffCSP	4.88	10.13
	DiffCSP++	4.48	4.44
	FlowMM	4.23	7.79
	MiAD	14.24	25.43
	SymmCD	3.83	<b>3.26</b>
	CrystalDiT	5.05	13.69
	ADiT	12.61	28.12
	MatterGen	6.88	10.11
	Chemeleon2-RL	6.65	6.04
	WyFormer	<b>3.40</b>	4.00
Alex-MP-20	MatterGen	3.67	3.01
	WyFormer	<b>0.96</b>	<b>1.17</b>

structure is new for a particular combination of elements, its underlying structure type has already been explored elsewhere (Eckert et al., 2024).

**Benchmark on Generative Models.** We next benchmark two major classes of state-of-the-art crystal generative models trained on MP-20: (i) all-atom diffusion or flow-matching models and (ii) space group symmetry-based models, spanning 10 models in total. Category (i) includes DiffCSP (Jiao et al., 2023), FlowMM (Miller et al., 2024), MiAD (Okhotin et al., 2025), CrystalDiT (Yi et al., 2025), ADiT (Joshi et al., 2025), Chemeleon2-RL (Park & Walsh, 2025), and MatterGen (Zeni et al., 2025), whereas category (ii) includes DiffCSP++ (Jiao et al., 2024), SymmCD (Levy et al., 2025), and WyFormer (Kazeev et al., 2025). Their central difference lies in the space over which generation is parameterized. All-atom models operate directly on the full set of atoms in the unit cell, without explicitly restricting the generative trajectory by crystallographic symmetry. Symmetry-based models instead operate in the asymmetric-unit, or Wyckoff subspace, representing crystals through discrete symmetry elements such as space group and site symmetry within a constrained parameterization. This both regularizes the search space and reduces the effective number of degrees of freedom. The two classes therefore offer a natural comparison between symmetry-agnostic generation and generation that is symmetry-aware by design.

Table 2 reports the order-(dis)order family tree benchmark for the 10 generative models trained on MP-20, together with MatterGen and WyFormer additionally trained on Alex-MP-20. The most immediate observation is the elevated

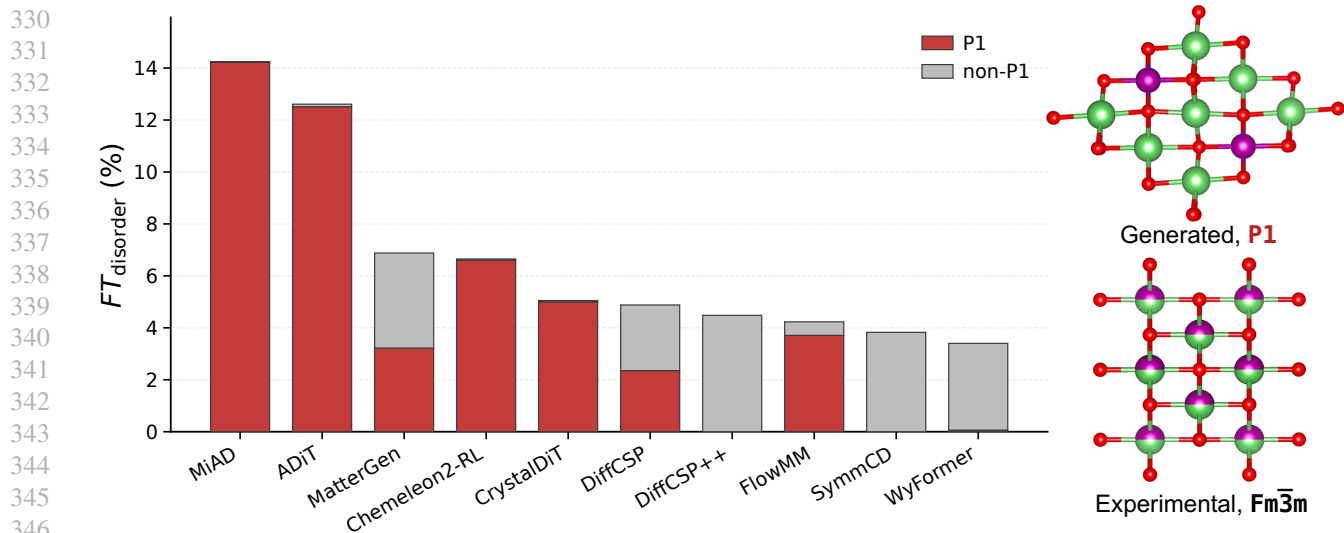


Figure 3. Space group decomposition of  $FT_{\text{disorder}}$  across crystal generative models. For each model,  $FT_{\text{disorder}}$  is decomposed by the space group of the ordered generated structure, with P1 highlighted in red. The dominant contribution of P1 in all-atom models points to the possibility that some of the overproduced P1 structures originate as ordered children of known disordered parent phases. Right, an example of generated P1 child and its matched ICSD disordered parent in  $Fm\bar{3}m$ .

$FT_{\text{disorder}}$  of MiAD and ADiT, at 14.24% and 12.61%, respectively, followed by MatterGen. Notably, all three belong to the symmetry-agnostic all-atom models. These high values indicate that a substantial fraction of generated structures that appear novel at the ordered-child level may in fact fall within family trees whose disordered parent phases are already known experimentally. In other words, apparent novelty at the ordered level can mask a closer relationship to previously realized disorder-rooted families. This trend is particularly notable in light of the reported synthesis case of  $TaCr_2O_6$  generated by MatterGen. Although initially proposed as a new ordered structure, the synthesized product was later found to correspond instead to a Cr-rich variant of an existing disordered parent phase (Juelsholt, 2026). In this context, the high  $FT_{\text{disorder}}$  values of this model class suggest that the reported example may be symptomatic of a more general trend.

By contrast, WyFormer achieves the lowest  $FT_{\text{disorder}}$ , at 3.40%, followed by SymmCD, and both belong to the space group symmetry-based models. A similar pattern is observed when MatterGen and WyFormer are trained on AlexMP-20: WyFormer again exhibits an  $FT_{\text{disorder}}$  nearly four times lower than that of MatterGen. Together, these results suggest that symmetry-regularized models are substantially less likely to generate ordered structures that fall within family trees of known disordered parents, and are therefore less prone to producing ordered candidates shadowed by experimentally reported disordered phases. A similar trend is observed for  $FT_{\text{order}}$ . All-atom models again show consistently higher values, indicating that their generations more frequently occupy regions of already known order-order

family space. By contrast, the lower  $FT_{\text{order}}$  of symmetry-based models suggests that they are more likely to generate ordered structures that are novel not only as individual entries, but also at the family level.

This contrast is consistent with the design of the two model classes. WyFormer, one of the leading symmetry-constrained models, achieves the highest uniqueness and novelty in the Wyckoff subspace relative to symmetry-agnostic models such as ADiT and MatterGen (Betala et al., 2025). Our results suggest that this advantage at the child level may translate into the family level: models designed to operate explicitly in a symmetry-constrained design space appear better positioned to discover genuinely new family-level configurations, rather than merely new ordered structures within already known families. This should be a consideration for future generative models in order to achieve the goal of synthesizing truly novel materials in experiments.

**Stable P1.** Space group P1 is the simplest crystallographic symmetry, containing no symmetry operations beyond identity and lattice translations. In nature, such trivial symmetry is rare: its occurrence in ICSD amounts to less than 1%, and even this fraction may be inflated by misreported cases (Urusov & Nadezhina, 2009; Marsh, 1999). In contrast, all-atom generative models are known to overproduce P1 structures, which in some cases constitute nearly half of generated samples, far exceeding the experimental distribution. Whether these *novel* and *stable* P1 structures represent a byproduct of imperfect computational evaluation or a genuinely unexplored class of experimentally accessible

structure types remains an open question in the community.

To examine this, we analyzed the space group distribution of ordered structures identified as children of existing disordered parents for each model. Figure 3 highlights the  $FT_{\text{disorder}}$  of each model, with  $P1$  contributions marked in red. For all-atom models, including MiAD, ADiT, MatterGen, Chemeleon2-RL, CrystalDiT, DiffCSP, and FlowMM,  $P1$  dominates  $FT_{\text{disorder}}$ , accounting for anywhere from nearly half to essentially all child structures, making it apparent that a part of  $P1$  structures systematically appear as ordered configurations derived from existing disordered parent phases. An example is shown on the right side of Figure 3, where a  $P1$  structure generated by DiffCSP is traced back to an ICSD disordered parent crystallizing in space group  $Fm\bar{3}m$ , using our framework.

Further DFT validation of  $P1$  structures that have ICSD disordered parents, with MiAD and MatterGen taken as representative all-atom models, reveals that roughly 80–90% are metastable after relaxation, with energies above the hull below 100 meV/atom relative to the Materials Project database (Horton et al., 2025). Details of the DFT calculations are provided in Appendix D. These findings, therefore, suggest that the overproduction of apparently stable  $P1$  structures can be explained to some extent by their emergence as ordered children of high-symmetry disordered parent phases: configurations that, while compositionally disordered in experiment, are unlikely to be experimentally synthesized in their ordered  $P1$  form.

**Group-Subgroup Polymorphs.** Order-order family trees, as formalized in Section 2.3, are governed by group-subgroup relations. In general, ordered structures belonging to the same  $FT_{\text{order}}$  family do not necessarily share the same stoichiometry, because symmetry lowering can split Wyckoff orbits in ways that alter the allowed occupancy pattern. A particularly important subset arises when symmetry descent preserves the exact same stoichiometry. We refer to this subset as *group-subgroup polymorphs*.

These polymorphs represent a symmetry-grounded form of polymorphism: one ordered phase can be derived from another through a group-subgroup transition that lowers point-group symmetry while preserving composition. This perspective goes beyond treating polymorphs simply as distinct structures with the same stoichiometry, and instead isolates those connected by an explicit symmetry-descent pathway. A representative example is provided by the two ordered  $\text{Cu}_2\text{GeS}_3$  phases in Figure 1, experimentally reported in space groups  $\text{Imm}2$  and  $\text{Cc}$ , where one can be obtained from the other through symmetry breaking along a group-subgroup transition.

This behavior is not rare in experimental data. In ICSD, among roughly  $5 \times 10^4$  unique ordered compositions, about

10% exhibit polymorphism, and around 40% of those polymorphic systems contain at least one pair related by a group-subgroup transition. This indicates that a meaningful fraction of experimentally reported polymorphs are not merely alternative structural realizations of the same composition, but members of an explicit symmetry-connected lineage, namely an order-order family. Order-order family mapping therefore offers more than a descriptive framework for known polymorphs, but it could also provide a systematic basis for discovering new ones.

## 4. Discussion and Conclusion

The order-(dis)order family tree benchmark introduced here redefines novelty against experimentally grounded crystallographic lineage rather than isolated ordered endpoints. This family-level view provides a more meaningful notion of novelty in systems where multiple ordered phases descend from a common disordered parent. The benchmark is most informative when the relevant disordered parent is already present in the experimental record. For underexplored systems, the corresponding parent may simply remain unseen, and the absence of a family assignment marks regions of family space not yet experimentally charted rather than a genuine absence of disorder tendency.

More fundamentally, these results indicate that disorder is best treated not as a label to predict, but as a manifold to navigate. Extending access to unseen disordered parents is therefore a natural next step: it would improve how novelty is assessed, clarify what synthesis is actually likely to realize, and provide a broader structural basis for evaluating stability across experimentally relevant energetic landscapes. In this sense, the same family-manifold perspective should ultimately inform not only crystallographic novelty, but also how thermodynamic preference is estimated within an order-disorder family (Divilov et al., 2024). The order-(dis)order family tree framework provides a concrete structural basis for that shift.

In this work, we introduced order-(dis)order family trees as a symmetry-grounded framework for organizing ordered and disordered crystal structures through group-subgroup relations. By shifting the unit of analysis from isolated structures to crystallographic families, the framework makes it possible to assess novelty at the level most relevant to experimental realization. We therefore argue that order-(dis)order family matching should be treated as a key requirement in evaluating novelty in materials discovery systems. Overall, these results position disorder-aware family matching as a practical foundation for future discovery pipelines and a necessary step toward identifying genuinely new, experimentally realizable materials.

## Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

## References

Antypov, D., Collins, C. M., Dyer, M. S., Claridge, J. B., and Rosseinsky, M. J. Classification and statistical analysis of structural disorder in crystalline materials. *Journal of Applied Crystallography*, 58(3): 659–677, May 2025. ISSN 1600-5767. doi: 10.1107/s1600576725003000. URL <http://dx.doi.org/10.1107/S1600576725003000>.

Bärnighausen, H. Group-subgroup relations between space groups: A useful tool in crystal chemistry. *MATCH Communications in Mathematical and in Computer Chemistry*, 9:139–175, 1980.

Betala, S., Gleason, S. P., Ramlaoui, A., Xu, A., Channing, G., Levy, D., Fourier, C., Kazeev, N., Joshi, C. K., Kaba, S.-O., Therrien, F., Hernandez-Garcia, A., Mercado, R., Krishnan, N. M. A., and Duval, A. Lemat-genbench: A unified evaluation framework for crystal generative models, 2025. URL <https://arxiv.org/abs/2512.04562>.

Cheetham, A. K. and Seshadri, R. Artificial intelligence driving materials discovery? perspective on the article: Scaling deep learning for materials discovery. *Chemistry of Materials*, 36(8):3490–3495, April 2024. ISSN 1520-5002. doi: 10.1021/acs.chemmater.4c00643. URL <http://dx.doi.org/10.1021/acs.chemmater.4c00643>.

Divilov, S., Eckert, H., Hicks, D., Oses, C., Toher, C., Friedrich, R., Esters, M., Mehl, M. J., Zettel, A. C., Lederer, Y., Zurek, E., Maria, J.-P., Brenner, D. W., Campilongo, X., Filipović, S., Fahrenholtz, W. G., Ryan, C. J., DeSalle, C. M., Creales, R. J., Wolfe, D. E., Calzolari, A., and Curtarolo, S. Disordered enthalpy–entropy descriptor for high-entropy ceramics discovery. *Nature*, 625(7993):66–73, January 2024. ISSN 1476-4687. doi: 10.1038/s41586-023-06786-y. URL <http://dx.doi.org/10.1038/s41586-023-06786-y>.

Eckert, H., Divilov, S., Mehl, M. J., Hicks, D., Zettel, A. C., Esters, M., Campilongo, X., and Curtarolo, S. The aflow library of crystallographic prototypes: Part 4. *Computational Materials Science*, 240:112988, May 2024. ISSN 0927-0256. doi: 10.1016/j.commatsci.2024.112988. URL <http://dx.doi.org/10.1016/j.commatsci.2024.112988>.

Ganose, A., Sahasrabudde, H., Asta, M., Beck, K., Biswas, T., Bonkowski, A., Bustamante, J., Chen, X., Chiang, Y., Chrzan, D., Clary, J., Cohen, O., Ertural, C., Gallant, M., George, J., Gerits, S., Goodall, R., Guha, R., Hautier, G., Horton, M., Kaplan, A., Kingsbury, R., Kuner, M., Li, B., Linn, X., McDermott, M., Mohanakrishnan, R. S., Naik, A., Neaton, J., Persson, K., Petretto, G., Purcell, T., Ricci, F., Rich, B., Riebesell, J., Rignanese, G.-M., Rosen, A., Scheffler, M., Schmidt, J., Shen, J.-X., Sobolev, A., Sundararaman, R., Tezak, C., Trinquet, V., Varley, J., Vigil-Fowler, D., Wang, D., Waroquiers, D., Wen, M., Yang, H., Zheng, H., Zheng, J., Zhu, Z., and Jain, A. Atomate2: Modular workflows for materials science. January 2025. doi: 10.26434/chemrxiv-2025-tcr5h. URL <http://dx.doi.org/10.26434/chemrxiv-2025-tcr5h>.

Gratias, D. and Quiquandon, M. Bicrystallography and beyond: Example of group–subgroup phase transformations. *Crystals*, 10(7):560, July 2020. ISSN 2073-4352. doi: 10.3390/cryst10070560. URL <http://dx.doi.org/10.3390/cryst10070560>.

Han, Y., Ding, C., Wang, J., Gao, H., Shi, J., Yu, S., Jia, Q., Pan, S., and Sun, J. Efficient crystal structure prediction based on the symmetry principle. *Nature Computational Science*, 5(3):255–267, February 2025. ISSN 2662-8457. doi: 10.1038/s43588-025-00775-z. URL <http://dx.doi.org/10.1038/s43588-025-00775-z>.

Horton, M. K., Huck, P., Yang, R. X., Munro, J. M., Dwaraknath, S., Ganose, A. M., Kingsbury, R. S., Wen, M., Shen, J. X., Mathis, T. S., Kaplan, A. D., Berket, K., Riebesell, J., George, J., Rosen, A. S., Spotte-Smith, E. W. C., McDermott, M. J., Cohen, O. A., Dunn, A., Kuner, M. C., Rignanese, G.-M., Petretto, G., Waroquiers, D., Griffin, S. M., Neaton, J. B., Chrzan, D. C., Asta, M., Hautier, G., Cholia, S., Ceder, G., Ong, S. P., Jain, A., and Persson, K. A. Accelerated data-driven materials science with the materials project. *Nature Materials*, 24(10):1522–1532, 2025. ISSN 1476-4660. doi: 10.1038/s41563-025-02272-0. URL <http://dx.doi.org/10.1038/s41563-025-02272-0>.

Huang, Y., Nong, W., Yamazaki, S., Petersen, M. H., Wang, J., Zhu, R., and Hippalgaonkar, K. Sword: Symmetry and wyckoff-sequence of ordered and disordered crystals, 2026. URL <https://arxiv.org/abs/2604.17994>.

Ivantchev, S., Kroumova, E., Madariaga, G., Pérez-Mato, J. M., and Aroyo, M. I. Subgroupgraph: a computer program for analysis of group–subgroup relations between space groups. *Journal of Applied Crystallography*, 33(4):1190–1191, August 2000. ISSN 0021-8898. doi: 10.1107/s0021889800007135. URL <http://dx.doi.org/10.1107/S0021889800007135>.

- 495 Jakob, K. S., Walsh, A., Reuter, K., and Margraf, J. T.  
 496 Learning crystallographic disorder: Bridging prediction  
 497 and experiment in materials discovery. *Advanced Ma-*  
 498 *terials*, 38(5), October 2025. ISSN 1521-4095. doi:  
 499 10.1002/adma.202514226. URL [http://dx.doi.](http://dx.doi.org/10.1002/adma.202514226)  
 500 [org/10.1002/adma.202514226](http://dx.doi.org/10.1002/adma.202514226).
- 501 Jiao, R., Huang, W., Lin, P., Han, J., Chen, P., Lu, Y., and  
 502 Liu, Y. Crystal structure prediction by joint equivariant  
 503 diffusion, 2023. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2309.04475)  
 504 [2309.04475](https://arxiv.org/abs/2309.04475).
- 505 Jiao, R., Huang, W., Liu, Y., Zhao, D., and Liu, Y. Space  
 506 group constrained crystal generation, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2402.03992)  
 507 [2402.03992](https://arxiv.org/abs/2402.03992).
- 508 Joshi, C. K., Fu, X., Liao, Y.-L., Gharakhanyan, V., Miller,  
 509 B. K., Sriram, A., and Ulissi, Z. W. All-atom diffusion  
 510 transformers: Unified generative modelling of molecules  
 511 and materials, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2503.03965)  
 512 [2503.03965](https://arxiv.org/abs/2503.03965).
- 513 Juelsholt, M. Continued challenges in high-throughput  
 514 materials predictions: Mattergen predicts compounds  
 515 from the training dataset. January 2026. doi: 10.26434/  
 516 chemrxiv-2025-mkls8/v3. URL [http://dx.doi.](http://dx.doi.org/10.26434/chemrxiv-2025-mkls8/v3)  
 517 [org/10.26434/chemrxiv-2025-mkls8/v3](http://dx.doi.org/10.26434/chemrxiv-2025-mkls8/v3).
- 518 Kazeev, N., Nong, W., Romanov, I., Zhu, R., Ustyuzhanin,  
 519 A., Yamazaki, S., and Hippalgaonkar, K. Wyckoff trans-  
 520 former: Generation of symmetric crystals. *Proceedings*  
 521 *of the 42nd International Conference on Machine Learn-*  
 522 *ing*, 2025. doi: 10.48550/ARXIV.2503.02407. URL  
 523 <https://arxiv.org/abs/2503.02407>.
- 524 Kresse, G. and Furthmüller, J. Efficient iterative schemes  
 525 forab initio total-energy calculations using a plane-wave  
 526 basis set. *Physical Review B*, 54(16):11169–11186, Oc-  
 527 tober 1996. ISSN 1095-3795. doi: 10.1103/physrevb.  
 528 54.11169. URL [http://dx.doi.org/10.1103/](http://dx.doi.org/10.1103/PhysRevB.54.11169)  
 529 [PhysRevB.54.11169](http://dx.doi.org/10.1103/PhysRevB.54.11169).
- 530 Kresse, G. and Joubert, D. From ultrasoft pseudopotentials  
 531 to the projector augmented-wave method. *Physi-*  
 532 *cal Review B*, 59(3):1758–1775, January 1999. ISSN  
 533 1095-3795. doi: 10.1103/physrevb.59.1758. URL <http://dx.doi.org/10.1103/PhysRevB.59.1758>.
- 534 Leeman, J., Liu, Y., Stiles, J., Lee, S. B., Bhatt, P., Schoop,  
 535 L. M., and Palgrave, R. G. Challenges in high-throughput  
 536 inorganic materials prediction and autonomous synthesis.  
 537 *PRX Energy*, 3(1), March 2024. ISSN 2768-5608. doi:  
 538 10.1103/prxenergy.3.011002. URL [http://dx.doi.](http://dx.doi.org/10.1103/PRXEnergy.3.011002)  
 539 [org/10.1103/PRXEnergy.3.011002](http://dx.doi.org/10.1103/PRXEnergy.3.011002).
- 540 Levy, D., Panigrahi, S. S., Kaba, S.-O., Zhu, Q., Lee,  
 541 K. L. K., Galkin, M., Miret, S., and Ravanbakhsh, S.  
 542 Symmcd: Symmetry-preserving crystal generation with  
 543 diffusion models, 2025. URL [https://arxiv.org/](https://arxiv.org/abs/2502.03638)  
 544 [abs/2502.03638](https://arxiv.org/abs/2502.03638).
- 545 Marsh, R. E. P1 or p1<sup>-</sup>? or something else? *Acta*  
 546 *Crystallographica Section B Structural Science*, 55(6):  
 547 931–936, December 1999. ISSN 0108-7681. doi:  
 548 10.1107/s0108768199009441. URL [http://dx.doi.](http://dx.doi.org/10.1107/s0108768199009441)  
 549 [org/10.1107/s0108768199009441](http://dx.doi.org/10.1107/s0108768199009441).
- 550 Merchant, A., Batzner, S., Schoenholz, S. S., Aykol, M.,  
 551 Cheon, G., and Cubuk, E. D. Scaling deep learn-  
 552 ing for materials discovery. *Nature*, 624(7990):80–85,  
 553 November 2023. ISSN 1476-4687. doi: 10.1038/  
 554 s41586-023-06735-9. URL [http://dx.doi.org/](http://dx.doi.org/10.1038/s41586-023-06735-9)  
 555 [10.1038/s41586-023-06735-9](http://dx.doi.org/10.1038/s41586-023-06735-9).
- 556 Miller, B. K., Chen, R. T. Q., Sriram, A., and Wood, B. M.  
 557 Flowmm: Generating materials with riemannian flow  
 558 matching. 2024. doi: 10.48550/ARXIV.2406.04713.  
 559 URL <https://arxiv.org/abs/2406.04713>.
- 560 Okhotin, A., Nakhodnov, M., Kazeev, N., Ustyuzhanin,  
 561 A. E., and Vetrov, D. Miad: Mirage atom diffusion for de  
 562 novo crystal generation, 2025. URL [https://arxiv.](https://arxiv.org/abs/2511.14426)  
 563 [org/abs/2511.14426](https://arxiv.org/abs/2511.14426).
- 564 Ong, S. P., Richards, W. D., Jain, A., Hautier, G., Kocher,  
 565 M., Cholia, S., Gunter, D., Chevrier, V. L., Persson, K. A.,  
 566 and Ceder, G. Python materials genomics (pymatgen): A  
 567 robust, open-source python library for materials analysis.  
 568 *Computational Materials Science*, 68:314–319, February  
 569 2013. ISSN 0927-0256. doi: 10.1016/j.commatsci.2012.  
 570 10.028. URL [http://dx.doi.org/10.1016/j.](http://dx.doi.org/10.1016/j.commatsci.2012.10.028)  
 571 [commatsci.2012.10.028](http://dx.doi.org/10.1016/j.commatsci.2012.10.028).
- 572 Park, H. and Walsh, A. Guiding generative models to un-  
 573 cover diverse and novel crystals via reinforcement learn-  
 574 ing, 2025. URL [https://arxiv.org/abs/2511.](https://arxiv.org/abs/2511.07158)  
 575 [07158](https://arxiv.org/abs/2511.07158).
- 576 Perdew, J. P., Burke, K., and Ernzerhof, M. Generalized gra-  
 577 dient approximation made simple. *Physical Review Let-*  
 578 *ters*, 77(18):3865–3868, October 1996. ISSN 1079-7114.  
 579 doi: 10.1103/physrevlett.77.3865. URL [http://dx.](http://dx.doi.org/10.1103/PhysRevLett.77.3865)  
 580 [doi.org/10.1103/PhysRevLett.77.3865](http://dx.doi.org/10.1103/PhysRevLett.77.3865).
- 581 Riebesell, J., Goodall, R. E. A., Benner, P., Chiang, Y.,  
 582 Deng, B., Ceder, G., Asta, M., Lee, A. A., Jain, A., and  
 583 Persson, K. A. A framework to evaluate machine learning  
 584 crystal stability predictions. *Nature Machine Intelligence*,  
 585 7(6):836–847, 2025. ISSN 2522-5839. doi: 10.1038/  
 586 s42256-025-01055-1. URL [http://dx.doi.org/](http://dx.doi.org/10.1038/s42256-025-01055-1)  
 587 [10.1038/s42256-025-01055-1](http://dx.doi.org/10.1038/s42256-025-01055-1).

- 550 Rühl, S. The inorganic crystal structure database  
 551 ( $\text{\textasciitilde{scpd}}/\text{\textasciitilde{scpd}}$ ): A tool for materials sciences, Au-  
 552 gust 2019. URL [http://dx.doi.org/10.1002/](http://dx.doi.org/10.1002/9783527802265.ch2)  
 553 [9783527802265.ch2](http://dx.doi.org/10.1002/9783527802265.ch2).  
 554
- 555 Spek, A. L. What makes a crystal structure report  
 556 valid? *Inorganica Chimica Acta*, 470:232–237, Jan-  
 557 uary 2018. ISSN 0020-1693. doi: 10.1016/j.ica.2017.  
 558 04.036. URL [http://dx.doi.org/10.1016/j.](http://dx.doi.org/10.1016/j.ica.2017.04.036)  
 559 [ica.2017.04.036](http://dx.doi.org/10.1016/j.ica.2017.04.036).  
 560
- 561 Szymanski, N. J., Rendy, B., Fei, Y., Kumar, R. E., He, T.,  
 562 Milsted, D., McDermott, M. J., Gallant, M., Cubuk, E. D.,  
 563 Merchant, A., Kim, H., Jain, A., Bartel, C. J., Persson, K.,  
 564 Zeng, Y., and Ceder, G. An autonomous laboratory for  
 565 the accelerated synthesis of inorganic materials. *Nature*,  
 566 624(7990):86–91, November 2023. ISSN 1476-4687.  
 567 doi: 10.1038/s41586-023-06734-w. URL [http://dx.](http://dx.doi.org/10.1038/s41586-023-06734-w)  
 568 [doi.org/10.1038/s41586-023-06734-w](http://dx.doi.org/10.1038/s41586-023-06734-w).  
 569
- 570 Togo, A., Shinohara, K., and Tanaka, I. Spglib: a software  
 571 library for crystal symmetry search. *Science and Tech-*  
 572 *nology of Advanced Materials: Methods*, 4(1), October  
 573 2024. ISSN 2766-0400. doi: 10.1080/27660400.2024.  
 574 2384822. URL [http://dx.doi.org/10.1080/](http://dx.doi.org/10.1080/27660400.2024.2384822)  
 575 [27660400.2024.2384822](http://dx.doi.org/10.1080/27660400.2024.2384822).  
 576
- 577 Urusov, V. S. and Nadezhina, T. N. Frequency distribu-  
 578 tion and selection of space groups in inorganic crys-  
 579 tal chemistry. *Journal of Structural Chemistry*, 50  
 580 (S1):22–37, December 2009. ISSN 1573-8779. doi:  
 581 10.1007/s10947-009-0186-9. URL [http://dx.doi.](http://dx.doi.org/10.1007/s10947-009-0186-9)  
 582 [org/10.1007/s10947-009-0186-9](http://dx.doi.org/10.1007/s10947-009-0186-9).  
 583
- 584 Wyckoff, R. W. G. *The Analytical Expression of the Results*  
 585 *of the Theory of Space-Groups*, volume 318. Carnegie  
 586 Institution of Washington, 1922.  
 587
- 588 Xie, T., Fu, X., Ganea, O.-E., Barzilay, R., and Jaakkola, T.  
 589 Crystal diffusion variational autoencoder for periodic ma-  
 590 terial generation, 2021. URL [https://arxiv.org/](https://arxiv.org/abs/2110.06197)  
 591 [abs/2110.06197](https://arxiv.org/abs/2110.06197).  
 592
- 593 Yang, K., Oses, C., and Curtarolo, S. Modeling off-  
 594 stoichiometry materials with a high-throughput ab-initio  
 595 approach. *Chemistry of Materials*, 28(18):6484–6492,  
 596 September 2016. ISSN 1520-5002. doi: 10.1021/acs.  
 597 chemmater.6b01449. URL [http://dx.doi.org/](http://dx.doi.org/10.1021/acs.chemmater.6b01449)  
 598 [10.1021/acs.chemmater.6b01449](http://dx.doi.org/10.1021/acs.chemmater.6b01449).  
 599
- 600 Yi, X., Xu, G., Xiao, X., Zhang, Z., Liu, L., Bian, Y., and  
 601 Zhao, P. Crystaldit: A diffusion transformer for crystal  
 602 generation, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2508.16614)  
 603 [2508.16614](https://arxiv.org/abs/2508.16614).  
 604
- Zagorac, D., Müller, H., Ruehl, S., Zagorac, J., and Rehme,  
 S. Recent developments in the inorganic crystal struc-  
 ture database: theoretical crystal structure data and re-  
 lated features. *Journal of Applied Crystallography*, 52  
 (5):918–925, September 2019. ISSN 1600-5767. doi:  
 10.1107/s160057671900997x. URL [http://dx.doi.](http://dx.doi.org/10.1107/s160057671900997x)  
[org/10.1107/s160057671900997x](http://dx.doi.org/10.1107/s160057671900997x).
- Zeni, C., Pinsler, R., Zügner, D., Fowler, A., Horton, M.,  
 Fu, X., Wang, Z., Shysheya, A., Crabbé, J., Ueda, S.,  
 Sordillo, R., Sun, L., Smith, J., Nguyen, B., Schulz, H.,  
 Lewis, S., Huang, C.-W., Lu, Z., Zhou, Y., Yang, H., Hao,  
 H., Li, J., Yang, C., Li, W., Tomioka, R., and Xie, T. A  
 generative model for inorganic materials design. *Nature*,  
 639(8055):624–632, January 2025. ISSN 1476-4687. doi:  
 10.1038/s41586-025-08628-5. URL [http://dx.doi.](http://dx.doi.org/10.1038/s41586-025-08628-5)  
[org/10.1038/s41586-025-08628-5](http://dx.doi.org/10.1038/s41586-025-08628-5).

## A. Order-(Dis)Order Family Matching Framework

Our order-(dis)order family matching framework is implemented in the `SWORDFamilyMatcher` module. For a query crystal structure  $S$ , the procedure first computes its SWORD child label,

$$\ell_{\text{child}}(S) = \text{SWORD}(S),$$

using symmetry analysis with `spglib` (Togo et al., 2024) through the SWORD labeling routine. In practice, the structure is first parsed into a `pymatgen Structure` (Ong et al., 2013), and partially occupied sites are renormalized when needed so that occupancies sum to unity before label construction. The overall procedure is summarized in Algorithm 1.

The central idea is to infer possible higher-symmetry parent descriptions by selectively removing chemical distinctions and recomputing the corresponding symmetry-aware label. If an ordered or partially ordered structure becomes compatible with a more symmetric description once certain species distinctions are masked, then that higher-symmetry description is taken as a plausible parent candidate.

Let  $E(S)$  denote the set of distinct element types in  $S$ , with  $n = |E(S)|$ . For a chosen subset  $M \subseteq E(S)$ , we define a masking operation  $\mathcal{M}_M(S)$  that replaces every element in  $M$  by a common dummy species  $X$ . Recomputing the SWORD label after masking gives a candidate parent label,

$$\ell_{\text{parent}}^{(M)}(S) = \text{SWORD}(\mathcal{M}_M(S)).$$

Because the raw label obtained after masking is expressed in terms of the dummy species  $X$ , we then restore the chemically meaningful mixed-site description by replacing each occurrence of  $X$  with the corresponding masked element set. This yields the final parent label used for matching.

For an ordered query, all nontrivial masks with  $2 \leq |M| \leq n$  are enumerated, giving

$$\sum_{k=2}^n \binom{n}{k} = 2^n - n - 1$$

possible masking patterns. For a query that is already disordered, masking is restricted to the mixed-element groups already present in the SWORD child label. This restriction is important because an element may appear simultaneously on ordered and disordered sites, so the relevant family relation is not obtained by arbitrary recombination of all species. Instead, the existing mixed-occupancy pattern encoded in the child label is treated as the chemically meaningful starting point, and only those disorder groups are further abstracted when generating parent candidates.

Accordingly, the parent-label set is written as

$$\mathcal{P}(S) = \left\{ \ell_{\text{parent}}^{(M)}(S) \right\},$$

and the family label set of the query is

$$\mathcal{F}(S) = \{ \ell_{\text{child}}(S) \} \cup \mathcal{P}(S).$$

Given a reference structure  $R$  processed in the same way, we regard  $R$  as family-related to  $S$  whenever their family label sets intersect,

$$\mathcal{F}(S) \cap \mathcal{F}(R) \neq \emptyset.$$

For ordered queries, the implementation includes an additional vacancy-ordering augmentation step before parent enumeration. If the child label contains no mixed occupancy marker, the structure is analyzed for symmetry-compatible latent vacancy ordering using the `find_vacancy_ordered` routine. This routine identifies candidate void sites from a periodic Voronoi construction, evaluates their local environments against species-specific coordination templates, and fills only those sites whose geometry and nearest-neighbor spacing are consistent with the insertion of an existing species. When such a filled structure is obtained, parent labels are generated both from the original query and from the vacancy-filled variant. This step allows the matcher to recover family relations that may otherwise remain hidden when an experimentally relevant parent description is more naturally expressed after restoring an ordered vacancy configuration.

In the pairwise setting, two structures  $S_1$  and  $S_2$  are assigned to the same family if

$$\mathcal{F}(S_1) \cap \mathcal{F}(S_2) \neq \emptyset.$$

In database matching, the same criterion is used against every reference entry. If the matched reference child label contains mixed occupancy, the hit is recorded as a disordered-family match; otherwise, if the matched reference child label is ordered and distinct from the query child label, it is recorded as an order-order family match.

As a concrete example, for the ordered  $\text{Cu}_2\text{GeS}_3$  query structure used in Figure 1, `SWORDFamilyMatcher` returns the following SWORD family dictionary:

```
{'child_label': 'a_b_c2_44_Ge_S_{Cu,S}',
 'parent_labels': ['a_c_216_S_{Cu+Ge}',
 'a_b_c2_44_Ge_{Cu+S}_{2(Cu+S)}',
 'a_b_c2_44_{Ge+S}_{Ge+S}_{Cu,Ge+S}',
 'a_227_{Cu+Ge+S}'],
 'matched_ordered_labels': ['a6_9_{2Cu,3S,Ge}'],
 'matched_ordered_ids': [85138, 146583],
 'matched_disordered_labels': ['a_d_b_122_Cu_S_{Cu+Ge}', 'a_c_216_S_{Cu+Ge}'],
 'matched_disordered_ids': [627781, 43531, 102963, 100955, 627773]}
```

This concrete example shows how the matcher encodes the query structure, its candidate parent descriptions, and the ordered and disordered ICSD relatives that together define the family recovered in Figure 1.

## B. A-Lab

*Table S1.* Comparison between ICSD Collection Codes of the previously reported disordered parents and those assigned by our framework for 22 A-Lab target phases. White rows (✓) denote agreement with previously reported parent assignments. Grey rows (□) indicate positional-disorder cases beyond the scope of the present framework. Green rows (\*) highlight newly identified disordered parents, demonstrating that our framework both recovers known parents and identifies previously unreported ones.

#	Target Phase	MP ID	ICSD Code <sup>1</sup>	ICSD Code <sup>2</sup>	Remark
1	$\text{Ba}_2\text{ZrSnO}_6$	mp-1228067	[43137]	[43137, 176330, 176331, 176329]	✓
2	$\text{FeSb}_3\text{Pb}_4\text{O}_{13}$	mp-1224890	[60805]	[60805]	✓
3	$\text{Hf}_2\text{Sb}_2\text{Pb}_4\text{O}_{13}$	mp-1224490	[62723]	[62723]	✓
4	$\text{InSb}_3\text{Pb}_4\text{O}_{13}$	mp-1223746	[41119]	[41119]	✓
5	$\text{KMn}_3\text{O}_6$	mp-1016190	[240249]	NaN	□
6	$\text{KNaP}_6(\text{PbO}_3)_8$	mp-1223429	[182501]	[182500, 182501, 182502]	✓
7	$\text{KNaTi}_2(\text{PO}_5)_2$	mp-1211611	[59284]	[67539, 71239]	□
8	$\text{K}_2\text{TiCr}(\text{PO}_4)_3$	mp-1224541	NaN	[280999]	*
9	$\text{KPr}_9(\text{Si}_3\text{O}_{13})_2$	mp-1223421	NaN	[153272]	*
10	$\text{K}_4\text{MgFe}_3(\text{PO}_4)_5$	mp-532755	[161484]	[161484]	✓
11	$\text{K}_4\text{TiSn}_3(\text{PO}_5)_4$	mp-1224290	[250088]	[250088, 72720, 91534, 250087]	✓
12	$\text{NaCaMgFe}(\text{SiO}_3)_4$	mp-1221075	[75294]	[263074, 263075, 263076, 263077, 263078, ...]	□
13	$\text{Mg}_3\text{MnNi}_3\text{O}_8$	mp-1222170	[80306]	[80302, 80303, 80304, 80305, 80306, 80307]	✓
14	$\text{Mg}_3\text{NiO}_4$	mp-1099253	[290603]	[290603, 13774]	✓
15	$\text{MgTi}_2\text{NiO}_6$	mp-1221952	[171583]	[171583]	✓
16	$\text{MgTi}_4(\text{PO}_4)_6$	mp-1222070	[74287]	[74287, 290277]	✓
17	$\text{MgV}_4\text{Cu}_3\text{O}_{14}$	mp-1222158	[69731]	[69731, 69732]	✓
18	$\text{Mn}_2\text{VPO}_7$	mp-1210613	NaN	[250126]	*
19	$\text{Mn}_4\text{Zn}_3(\text{NiO}_6)_2$	mp-1222033	[92223]	[92222, 92223, 92224]	✓
20	$\text{Na}_3\text{Ca}_{18}\text{Fe}(\text{PO}_4)_{14}$	mp-725491	[85103]	[85103]	✓
21	$\text{Sn}_2\text{Sb}_2\text{Pb}_4\text{O}_{13}$	mp-1219056	[62722]	[62722]	✓
22	$\text{Zr}_2\text{Sb}_2\text{Pb}_4\text{O}_{13}$	mp-1215826	[62721]	[62721]	✓

<sup>1</sup>Leeman et al. (Leeman et al., 2024)

<sup>2</sup>This work

## C. MP-20

MP-20 has long been treated as one of the most representative training datasets for crystal generative models. It is defined as an ICSD-derived subset of Materials Project structures restricted to compounds with at most 20 atoms per unit cell, and contains roughly 45 k structures across the train, validation, and test splits. At the same time, when ordered ICSD entries are counted at the level of unique structures with no more than 20 atoms per unit cell, the total is only around 20, k. This discrepancy is difficult to reconcile with a strict one-to-one view of MP-20 as a subset of experimentally reported ordered ICSD phases.

One plausible explanation is that the ICSD provenance inherited in *legacy* Materials Project data is not always cleanly recoverable at the level of unique current ICSD entries, and that MP-20 also includes Materials Project processed or relaxed versions of experimentally derived structures. Under this interpretation, MP-20 should be understood not as a strict snapshot of unique ordered ICSD phases, but as a legacy computational dataset with partial ICSD provenance and additional expansion introduced during downstream curation.

This distinction matters for interpreting its family-tree statistics. In particular, the fact that  $FT_{\text{disorder}}$  reaches 23.27% indicates that more than one in five ordered structures in MP-20 have existing disordered parents in ICSD. This suggests that a non-negligible portion of MP-20 lies in crystallographic families for which disorder is already an experimentally established realization. As a consequence, models trained on MP-20, particularly all-atom generative models, may in part be learning to reproduce ordered children of existing disordered parents rather than structures that are genuinely novel at the family level. At a high level, this observation further motivates disorder-aware dataset construction and generative modeling in crystal structure generation.

## D. DFT details

We use DFT settings from Materials Project <https://docs.materialsproject.org/methodology/materials-methodology/calculation-details/gga+u-calculations/parameters-and-convergence> for structure relaxation and energy computation. In particular, we do GGA and GGA+U calculations with `atomate2.vasp.flows.mp.MPGGADoubleRelaxStaticMaker` (Ganose et al., 2025), which in turn relies on `pymatgen.io.vasp.sets.MPRelaxSet` and `pymatgen.io.vasp.sets.MPStaticSet` (Ong et al., 2013). Computations themselves were done with VASP (Kresse & Furthmüller, 1996) version 5.4.4. with the plane-wave basis set (Kresse & Furthmüller, 1996). The electron-ion interaction is described by the projector augmented wave (PAW) pseudo-potentials (Kresse & Joubert, 1999). The exchange-correlation of valence electrons is treated with the Perdew-Burke-Ernzerhof (PBE) functional within the generalized gradient approximation (GGA) (Perdew et al., 1996). The raw total energies computed by DFT were corrected with `MaterialsProject2020Compatibility` before putting into the `PhaseDiagram` to obtain the DFT  $E_{\text{hull}}$ . We used the MP convex hull `2023-02-07-ppd-mp.pkl.gz` distributed by `matbench-discovery` (Riebesell et al., 2025) as the reference hull.

```

770
771
772
773
774
775
776 Algorithm 1 SWORDFAMILYMATCHER
777
778 1: Input: Query structure  $S$  or two structures  $S_1, S_2$ 
779 2: Output: Family label set for  $S$ , or a Boolean family match for  $S_1, S_2$ 
780 3:
781 4: Function GETFAMILYLABELS( $S$ )
782 5: parse  $S$  into a crystal structure object
783 6: renormalize partial occupancies if needed
784 7:  $\ell_{\text{child}} \leftarrow \text{SWORD}(S)$ 
785 8:  $\mathcal{S}_{\text{proc}} \leftarrow [S]$ 
786 9: if  $\ell_{\text{child}}$  contains no mixed-occupancy marker then
787 10:    $S_{\text{fill}} \leftarrow \text{FINDVACANCYORDERED}(S)$ 
788 11:   if  $S_{\text{fill}}$  is valid then
789 12:     append  $S_{\text{fill}}$  to  $\mathcal{S}_{\text{proc}}$ 
790 13:   end if
791 14: end if
792 15:  $\mathcal{P} \leftarrow \emptyset$ 
793 16: for each  $\tilde{S} \in \mathcal{S}_{\text{proc}}$  do
794 17:    $\tilde{\ell}_{\text{child}} \leftarrow \text{SWORD}(\tilde{S})$ 
795 18:   extract mixed-element groups from  $\tilde{\ell}_{\text{child}}$ 
796 19:   if mixed-element groups are present then
797 20:     define mask sets from those existing mixed groups only
798 21:   else
799 22:     define all element subsets  $M \subseteq E(\tilde{S})$  with  $2 \leq |M| \leq |E(\tilde{S})|$ 
800 23:   end if
801 24:   for each mask set  $M$  do
802 25:     construct  $\mathcal{M}_M(\tilde{S})$  by replacing elements in  $M$  with dummy species  $X$ 
803 26:      $\ell_{\text{raw}} \leftarrow \text{SWORD}(\mathcal{M}_M(\tilde{S}))$ 
804 27:      $\ell_{\text{parent}} \leftarrow$  restore masked element group into  $\ell_{\text{raw}}$ 
805 28:     add  $\ell_{\text{parent}}$  to  $\mathcal{P}$ 
806 29:   end for
807 30: end for
808    $\mathcal{F}(S) = \{\ell_{\text{child}}\} \cup \mathcal{P}$ 
809 31:
810 32: Function FIT( $S_1, S_2$ )
811 33:  $\mathcal{F}_1 \leftarrow \text{GETFAMILYLABELS}(S_1)$ 
812 34:  $\mathcal{F}_2 \leftarrow \text{GETFAMILYLABELS}(S_2)$ 
813 35: if  $\mathcal{F}_1 \cap \mathcal{F}_2 \neq \emptyset$  then
814   True
815 36: else
816   False
817 37: end if
818
819
820
821
822
823
824

```