	00 00 00
Geographic-style maps with a local novelty	00
distance help navigate in the material space	00
Daniel Widdowson ¹ and Vitaliy Kurlin ^{1^*}	01
¹ Materials Innovation Factory, University of Liverpool, Oxford Street, Liverpool, L7 3NY, United Kingdom.	01 01 01 01
*Corresponding author(s). E-mail(s): vitaliy.kurlin@liverpool.ac.uk;	01 01 01
Abstract	01
With the advent of self-driving labs promising to synthesize large numbers of new materials, new automated tools are required for checking potential duplicates in existing structural databases before a material can be claimed as novel. To avoid duplication, we rigorously define the novelty metric of any periodic material as the smallest distance to its nearest neighbor among already known materials.	02 02 02 02 02
Using ultra-fast structural invariants, all such nearest neighbors can be found within seconds on a typical computer even if a given crystal is disguised by changing a unit cell, perturbing atoms, or replacing chemical elements. This real-time novelty check is demonstrated by finding near-duplicates of the 43 mate- rials produced by Berkeley's A-lab in the world's largest collections of inorganic structures, the Inorganic Crystal Structure Database and the Materials Project.	02 02 02 02 02 03 03
To help future self-driving labs successfully identify novel materials, we propose navigation maps of the materials space where any new structure can be quickly located by its invariant descriptors similar to a geographic location on Earth.	03 03 03
Keywords: materials space, crystal structure, isometry invariant, continuous metric	03
1 Introduction: how is the materials space defined?	03 03
The chemical space of all possible molecules is often estimated at the scale of 10^{60} [1]. Similar numbers are quoted for potential materials, though many polymorphs such as diamond and graphite have the same chemical composition and hence can only be dis- tinguished by their geometry. When materials are claimed to be novel amongst already known ones, we need to rigorously define what constitutes two materials being the "same or different" [2]. The definition of a <i>crystal structure</i> was finalized in the peri- odic case in [3], so we focus on ideal periodic crystals (briefly, <i>crystals</i>) as formalized below. When a material is disordered, we consider its closest periodic analogue. A <i>crystal</i> is usually given by a basis of vectors v_1, v_2, v_3 in Euclidean space \mathbb{R}^3 and a <i>motif</i> of atoms with chemical elements and fractional coordinates in this basis.	03 04 04 04 04 04 04 04 04

and 050If we forget about chemical elements, the atomic centers p_1, \ldots, p_m can be considered 051zero-sized points in the primitive unit cell $U = \{t_1 v_1 + t_2 v_2 + t_3 v_3 \mid t_1, t_2, t_3 \in [0, 1)\}$ defined by the basis v_1, v_2, v_3 . In dimension 2, the second picture of Fig. 1 highlights 052the square cell U with the orthonormal basis v_1, v_2 . Then the underlying *periodic* 053054point set of any crystal consists of infinitely many points $p_i + c_1 v_1 + c_2 v_2 + c_3 v_3$ for 055 $i = 1, \ldots, m$ and integer coefficients $c_1, c_2, c_3 \in \mathbb{Z}$. Infinitely many different pairs of a 056basis (or a primitive cell) and a motif M generate pointwise identical crystals, see a 057 detailed discussion of this ambiguity of the traditional definition in [3, section 2].

059 Fig. 1 Almost any tiny perturbation discontinuously scales up a primitive cell and makes unreliable any comparison based on cells or motifs. This discontinuity was resolved without relying on cells [4].

$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	000									
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	061	of 🔹	of 🔹	a verv small	•	٠	•	•	another tiny • • • • •	
$\begin{array}{c c} 063 \\ 064 \\ 065 \end{array} & \begin{array}{c} \bullet & \bullet \\ \bullet & \bullet \\ \bullet & \bullet \\ 065 \end{array} & \begin{array}{c} \bullet & \bullet \\ \bullet & \bullet \\ \bullet & \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \\$	062	of 🔹	o n 🔸	perturbation	•	٠	•	•	perturbation / \bullet / \bullet duplicates $\bullet \circ$ \bullet $\bullet \circ$	\langle
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	063	<u>.</u>	e •	$\sim V_2$	•	٠	•	٠	develops the of of of of an have	
000	$\begin{array}{c} 064 \\ 065 \end{array}$	₫ •	 •	symmetry	v ₁	•	•	•	cell volume / / / / / / primitive cells	

066

067

068 Since atoms always vibrate [5, chapter 1], their fractional coordinates are always 069 uncertain and will slightly deviate under repeated measurements even on the same 070 instrument. Almost any displacement of one atom breaks the symmetry and can arbi-071trarily scale up a primitive unit cell as in Fig. 1. This discontinuity of a reduced cell 072[6] was experimentally reported in 1965 [7, p. 80] and remained unresolved until 2022 073 [4] when all periodic crystals in the Cambridge Structural Database (CSD) [8] were 074distinguished within two days (now within an hour) on a modest desktop computer. 075 Several unexpected duplicates with identical geometries (almost to the last decimal 076 place in all cell parameters and atomic coordinates) but with different chemistry are 077 under investigation by five journals for data integrity [9, section 6].

078 Since crystal structures are determined in a rigid form, there is no sense in 079 distinguishing crystal representations related by a *rigid motion* (a composition of 080 translations and rotations in \mathbb{R}^3), which change a basis and atomic coordinates. On 081 the other hand, there is no sense to fix any threshold $\varepsilon > 0$ that would allow us to call 082 crystals the "same" if all their atomic centers (without chemical attributes) can be 083 matched up to ε -perturbations. Indeed, any periodic point sets can be connected by 084 sufficiently many ε -perturbations [9, Proposition 2.10], which makes the classification 085 based on any threshold $\varepsilon > 0$ trivial due to the transitivity axiom saying that if S is 086 equivalent to Q, and Q is equivalent to T, then S is equivalent to T [3, section 1]. 087

Hence a rigorous way to classify crystals under rigid motion, is to define the *crystal* 088 structure as a rigid class of periodic point sets, see [3, Definition 6]. Then any devi-089 ations of atomic positions are not ignored but continuously quantified by a distance 090 metric between different rigid classes. This definition would remain impractical unless 091 we can efficiently separate rigid classes by quickly computable *invariants* that are 092numerical properties preserved under rigid motion. The chemical composition written 093 as percentages of chemical elements is such an invariant but is *incomplete* because 094many polymorphs have the same composition but can not be matched by rigid motion. 095

1096 In the sequel, we will consider the sightly weaker equivalence of *isometry* (any 1097 distance-preserving transformation in \mathbb{R}^3), which is a composition of rigid motion 1098 and mirror reflections. Since mirror images can be distinguished by a suitable sign of 1099 orientation, the main difficulty is to classify periodic point sets under isometry.

When comparing crystals as periodic sets of atomic centers without chemical attributes, it might seem that all chemistry is lost. However, the fact that all (more than 850 thousand) periodic crystals in the CSD (apart from the investigated duplicates) can be distinguished by isometry invariants in section 2 implies that no information is lost so that all chemistry under standard conditions such as temperature and pressure is in principle reconstructable from sufficiently precise atomic geometry.

This Crystal Isometry Principle (CRISP) first appeared in 2022 [9, section 7] and was inspired by Richard Feynman's hint in Fig.1-7 [5, chapter 1], which distinguished 7 cubic crystals by their cube size in the first lecture "Atoms in motion", see Fig. 2 (left).

More importantly, when we consider atoms only as zero-sized points, we can study all periodic structures in a common space similar to the periodic table of all elements.

In the geographic analogy, the chemical composition can be compared to the altitude (the height above the sea level) of any location on Earth. If our geographic map is precise enough, we can determine the average temperature or any other property at every location. If we know the altitude (chemical composition) in addition to geographic coordinates (structural invariants), the property prediction will be easier. **Definition 1** (space of periodic materials). The Crystal Isometry Space $CRIS(\mathbb{R}^3)$ is 117the space of isometry classes of all periodic sets of points without atomic attributes. 118

Fig. 2 Left: the Crystal Isometry Principle says that all chemistry of any real periodic crystal under standard ambient conditions can be reconstructed from (the isometry class of) the periodic set of atomic centers given with precisely enough coordinates [4]. Right: most optimization methods output local optima without exploring the space around. De-fogging this Crystal Isometry Space $CRIS(\mathbb{R}^3)$ beyond known or predicted materials will enable a proper navigation across the crystal universe.



Since (the isometry class of) any periodic point set has a unique location in $CRIS(\mathbb{R}^3)$, all known materials can be considered 'visible stars' in this continuous universe. Any periodic crystal discovered in the future will appear at a new unique location like a 'new star', while all past crystals remain at the same locations. Until recently, optimizing complicated energy functions blindly climbed as a mountaineer to a high peak, illustrated in Fig. 2 (right), and stopping at many (approximations to) isolated peaks without even any reliable method to continuously measure distances between these peaks, while the remaining landscape was covered by clouds. Our vision is to map the full space $CRIS(\mathbb{R}^3)$ to enable a non-blind discovery of materials [10].

If we do not restrict the motif size, the space CRIS is infinitely dimensional. How-142ever, if we consider all periodic sets with m points in a motif, the resulting subspace 143 $CRIS(\mathbb{R}^3; m)$ has dimension 3m + 3 due to m triples x, y, z of atomic coordinates and 1446 parameters of a unit cell, of which 3 are neutralized by translations along basis vec-145tors. Alternatively, we can define a unit cell by 3 basis vectors with 3 coordinates, of 146which 6 are neutralized by 3+3 parameters of translations and rotations in \mathbb{R}^3 . 147

In the partial case m = 1, $CRIS(\mathbb{R}^3; 1)$ is a continuous 6-dimensional space of 3D lattices, which was previously cut in 14 disjoint subspaces of Bravais classes [11] but is now parametrized by complete invariants [12, 13]. Continuous maps of the simpler 3-dimensional space $CRIS(\mathbb{R}^2; 1)$ of 2D lattices recently appeared in [14], [15], [16].

152The full space $\operatorname{CRIS}(\mathbb{R}^3) = \bigcup_{m=1}^{+\infty} \operatorname{CRIS}(\mathbb{R}^3; m)$ is a union of infinitely many sub-153spaces for $m = 1, 2, 3, \ldots$ such that any periodic set with m points in a cell is infinitesimally close to infinitely many subspaces of sets with $2m, 3m, \ldots$ points in a 155primitive cell. Indeed, perturbations in Fig. 1 arbitrarily extend any given cell and make the extended cell primitive by a tiny displacement of any atom and all its translational copies. Crystals should be continuously compared only across multiple subspaces, not within one subspace $CRIS(\mathbb{R}^3; m)$ for a fixed number m of atoms. Any database of periodic crystals is a finite sample from the continuous space $CRIS(\mathbb{R}^3)$. 160

The first contribution of this work is the local novelty distance based on generically complete invariants, which identify closest neighbors of the 43 A-lab crystals in the Inorganic Crystal Structure Database (ICSD) [17] and Materials Project (MP) [18] within seconds on a desktop computer. The second contribution is the geographic-style maps showing how the ICSD and MP populate $CRIS(\mathbb{R}^3)$ in invariant coordinates.

Methods: invariant-based novelty distance metric $\mathbf{2}$

This section introduces a new metric LND (Local Novelty Distance) that satisfies all metric axioms and continuously quantifies in real time a deviation of any newly synthesized crystal from its nearest neighbor in an existing structural database.

170171

119120

121

122

123

124

125126

127

128

129130

131132

133

134

135

136

137

138

139

140

141

148

149

150

151

154

156

157

158

159

161

162

163

164

165166

167168

169

172173

175 2.1 Generically complete and continuous structural invariants

176 177 Definition 2 reminds us of the Pointwise Distance Distribution (PDD), which suffices 178 together with a lattice to reconstruct any generic periodic point set $S \subset \mathbb{R}^3$ up to 179 isometry by [4, Theorem 4.4] and [19, Theorem 5.8]. Generic means any set apart from 180 a singular subspace of measure 0, e.g. almost any noise makes every crystal generic.

The PDD is a matrix of inter-point distances and is stronger than the Pair Distribution Function (PDF) [20] in the sense that PDD can be simplified to PDF but distinguishes homometric structures [21] that have the same PDF [4, section 3].

Definition 2 (isometry invariant PDD(S;k)). Let $S \subset \mathbb{R}^n$ be a periodic point set 184with a motif $M = \{p_1, \ldots, p_m\}$. Fix an integer $k \ge 1$. For every point $p_i \in M$, let 185 $d_1(p) \leq \cdots \leq d_k(p)$ be the distances from p to its k nearest neighbors within the full 186infinite set S not restricted to any cell. The matrix D(S;k) has m rows consisting 187of the distances $d_1(p_i), \ldots, d_k(p_i)$ for $i = 1, \ldots, m$. If any $l \ge 1$ rows are identical to 188 each other, we collapse them into a single row and assign the weight l/m to this row. 189The resulting matrix of maximum m rows and k+1 columns including the extra (say, 1900-th) column of weights is called the *Pointwise Distance Distribution* PDD(S; k). 191

192 In Definition 2, any point $p_i \in M$ can have several different neighbors at the same 193 distance but the k smallest distances (without any indices or types of neighbors) are 194 always well-defined. The matrix PDD(S;k) has ordered columns (according to the 195 index of neighbors) but unordered rows because points of a motif of S are unordered. 196 The appendix computes a weighted PDD with atomic masses as extra weights of 197 rows in PDD(S;k). Using only on atomic centers detects duplicates where chemical 198 elements were artificially replaced without changing geometry, see [3, Table 1].

199 If the number k of neighbors increases to infinity, the asymptotic behavior of 200 distances to neighbors is described in terms of the Point Packing Coefficient below.

201 **Definition 3** (Point Packing Coefficient PPC). Let $S \subset \mathbb{R}^3$ be a periodic point set 203 with *m* atoms in a unit cell *U*. The *Point Packing Coefficient* is $PPC(S) = \sqrt[3]{\frac{\text{vol}(U)}{mV_3}}$,

where $\operatorname{vol}(U)$ is the volume of U, $V_3 = \frac{4}{3}\pi$ is the volume of the unit ball in \mathbb{R}^3 .

The distances in each row of PDD(S; k) asymptotically increase as $PPC(S)\sqrt[3]{k}$ by [9, Theorem 13]. This asymptotic behavior motivates the simplified invariants below. **Definition 4** (invariants AMD, ADA, PDA). The Average Minimum Distance AMD_k(S) is the weighted average of the k-th column of PDD(S; k). The Average Deviation from Asymptotic is $ADA_k(S) = AMD_k(S) - PPC(S)\sqrt[3]{k}$ for $k \ge 1$. The Pointwise Deviation from Asymptotic is the matrix PDA(S; k) obtained from PDD(S; k) by subtracting $PPC(S)\sqrt[3]{k}$ from any distance in row i and column k for $i, k \ge 1$.

214 215 Fig. 3 The average invariants AMD_k and ADA_k from Definition 4 for k = 1, ..., 25 and five simple crystals from the Materials Project, see more details and perovskite examples in the appendix.



The invariants AMD_k and ADA_k form vectors of length k, e.g. set $AMD(S;k) = (AMD_1(S), \ldots, AMD_k(S))$ and $ADA(S;k) = (ADA_1(S), \ldots, ADA_k(S))$. These vectors can be compared by many metrics. The metric $L_{\infty}(u, v) = \max_{i=1,\ldots,k} |u_i - v_i|$ for any vectors $u, v \in \mathbb{R}^k$ preserves the intuition of atomic displacements in the following sense. If S is obtained from Q by perturbing every point up to a small ε , then $L_{\infty}(AMD(S;k), AMD(Q;k)) \leq 2\varepsilon$ by [9, Theorem 9]. Other distances such as Euclidean can be considered but will accumulate a larger deviation depending on k.

All invariants above and metrics on them are measured in the same units as original 233coordinates, i.e. in Angstroms for crystals given by Crystallographic Information Files 234(CIFs). The Point Packing Coefficient PPC(S) was defined as the cube root of the 235cell volume per atom (of the same radius 1Å) and can be interpreted as an average 236radius of balls 'packed' in a unit cell. So PPC(S) is roughly inverse proportional to 237the physical density but they are exactly related only when materials have the same 238 average atomic mass (total mass of atoms in a unit cell divided by the cell volume). 239

240While $AMD_k(S)$ monotonically increases in k, the invariants $ADA_k(S)$ can be 241positive or negative as deviations around the asymptotic $PPC(S)\sqrt[3]{k}$. Fig. 4 reveals 242geometric differences between the mainly organic databases CSD and Crystallography 243Open Database (COD) [22] versus the more inorganic collections ICSD and MP. 244



Fig. 4 The averages of ADA_k and standard deviations (1 sigma shaded) vs $\sqrt[3]{k}$ for four databases. CSD

245

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281282283

284

285

286

287

The first average of ADA₁ $\in [-0.25, -0.17]$ in the top images of Fig. 4 can be explained by the presence of many hydrogen atoms, which have distances smaller than PPC(S) to their first neighbor in most organic materials. Indeed, hydrogens are usually bonded at distances less than 1.2\AA , while PPC(S) is often larger than 1.2\AA because most chemical elements have van der Waals radii above 1.2Å [23].

For inorganic materials, metal atoms or ions have relatively large distances to their first neighbors, so the average ADA₁ is in [0.58, 0.62] in the bottom images of Fig. 4.

If we increase k, the matrix PDD(S;k) and hence the vector ADA(S;k) become longer by including distance data to further neighbors but all initial values remain the same. Hence we consider k not as a parameter that changes the output but as a degree of approximation similarly to the number of decimal places on a calculator.

The experimental convergence $ADA_k \rightarrow 0$ as $k \rightarrow +\infty$ in Fig. 4justifies computing the distance L_{∞} between ADA vectors up to a reasonable k. We use k = 100 because all ADA_k for k > 100 are close to 0 (the range of 1 sigma between ± 0.2 Å) in Fig. 4.

2.2 Novelty distance based on practically complete invariants

This subsection introduces the Local Novelty Distance LND(S; D) of a periodic crystal S as a distance to the closest neighbor Q of S in a given dataset D. The LND will be measured as a distance between PDA(S; k) and PDA(Q; k) for $Q \in D$.

We can compare PDD matrices that have the same number of columns and possibly 288 different numbers of rows by interpreting PDD(S;k) as a distribution of unordered 289rows (or points in \mathbb{R}^k) with weights or probabilities. Many similarities between discrete 290 distributions such as the Cramer - von Mises distance and the Kullback - Leibler (KL) divergence fail the axioms of a metric, which are prerequisites for convergence guarantees. If the triangle axiom fails with any positive error, outputs of widely used clustering algorithms such as k-means and DBSCAN may not be trustworthy [24].

Though the KL divergence can be symmetrized to the Jensen-Shannon divergence whose square root becomes a metric [25], these divergences work best for distributions on the same finitely many discrete values, while PDD matrices of different crystals are unlikely to share any common rows (points in \mathbb{R}^k) consisting of k continuous distances.

We use the Earth Mover's Distance (EMD) on PDDs, see Definition 7 in the appendix because EMD is invariant under permutations of rows in PDDs and EMD satisfies all metric axioms [26, Appendix]. The EMD can be extended to more complicated Wasserstein metrics [27] but the simplest EMD behaves most nicely under bounded noise, motivated by atomic vibrations. If any point is perturbed up to ε , any inter-point distance (a value in PDD) can become smaller or larger but only up to 2ε , which will allow us to prove the same upper bound 2ε for EMD in Theorem 6.

307 **Definition 5** (Local Novelty Distance LND(S; D)). Let D be a finite dataset of 308 periodic point sets. Fix an integer $k \ge 1$. For any periodic point set S, the *Local Novelty* 309 *Distance* $\text{LND}(S; D) = \min_{\substack{Q \in D}} \text{EMD}(\text{PDA}(S; k), \text{PDA}(Q; k))$ is the shortest L_{∞} -based

310 EMD distance from S to its nearest neighbor Q in the given crystal dataset D.

³¹¹ If S is already contained in the dataset D, then LND(S; D) = 0, so S cannot be considered novel. Conversely, if LND(S; D) = 0 then S highly likely belongs to S, because PDD(S; 100) distinguished all non-duplicate periodic crystals in the CSD.

Also, for a generic periodic set S (away from a measure 0 subspace), PDD(S;k)315 with a big enough k and a lattice of S suffices to reconstruct S uniquely under isometry 316 in \mathbb{R}^n by [4, Theorem 4.4]. LND(S; D) is based on PDAs instead of PDDs because 317 distances to k-th neighbors in PDD(S; k) asymptotically increase as $PPC(S)\sqrt[3]{k}$ by 318 [9, Theorem 4.4]. If crystals S, Q have $PPC(S) \neq PPC(Q)$, the distance L_{∞} between 319 rows of PDDs equals the largest absolute difference of *i*-th distances, which likely 320 happens for i = k. So subtracting $PPC(S)\sqrt[3]{k}$ in Definition 4 makes any metric on 321 PDAs more informative than on PDDs. If a newly synthesized periodic crystal S is 322 a near-duplicate of some known $Q \in D$, then LND(S; D) is small as justified below. 323 The packing radius r(Q) is the minimum half-distance between any points of Q. 324

325 **Theorem 6.** If S is obtained from a crystal Q in a dataset D by perturbing every 326 point of Q up to $\varepsilon < r(Q)$, then $\text{LND}(S; D) \le 2\varepsilon$. To get S from a crystal $Q \in D$ with 327 LND(S; D) < 2r(Q), some atom of Q should be perturbed by at least 0.5LND(S; D).

Theorem 6 is proved in Appendix A. The distance LND(S; D) is called *local* because Definition 5 uses the first nearest neighbor of S in D. Another novelty of S can be characterized with respect to a global distribution of all crystals in D, which we will explore in a forthcoming work. The local novelty is more urgently needed to tackle the growing crisis of duplication in experimental and simulated databases, some of which were publicly rebutted in [28], [29], and [30], [3, Tables 1-2 in section 6], respectively.

³³⁵ 2.3 Insufficiency of past invariants and similarities of crystals ³³⁶

This subsection briefly reviews the past approaches to classify crystals. Some widely used similarities such as the Root Mean Square Deviation (RMSD) [31] deserve their own detailed discussions in another forthcoming work. Conventional settings were thoroughly developed to uniquely represent any periodic crystal in a reduced cell [32] and can be theoretically considered complete under rigid motion but discontinuously change under almost any perturbation of atoms in practice as shown in Fig. 1.

Indeed, perturbations in Fig. 1 apply to any crystal and can arbitrarily extend a reduced cell to a larger cell whose size cannot be reduced. Searching for a small perturbation (pseudo-symmetry) to make a cell smaller [33] inevitably uses thresholds and leads to a trivial classification due to the transitivity axiom, see [3, section 1].

The COMPACK algorithm [31] outputs an RMSD quantity by comparing finite 349portions of only molecular crystals. Its implementation in Mercury also uses thresh-350 olds for acceptable deviations of atoms and angles. Even if these thresholds are ignored 351(made large), the algorithm chooses one molecule in a unit cell and 14 (by default) clos-352est molecules around it. The resulting molecular group depends on a central molecule. 353 Even for simple crystals based on a single molecule as is often the case in Crystal 354Structure Prediction [34], the choice of 14 (or any other number of) neighbors can 355 be discontinuous when a central molecule has 14th and 15th neighbors at the same 356 distance. Selected clusters of molecules in two crystals require an optimal alignment, 357which is a hard problem because atomic sets can contain numerous indistinguishable 358 atoms, so the optimization must consider many potential permutations. This prob-359lem of exponentially many permutations was resolved in [35] but a choice of a single 360 atomic environment in a periodic crystal remains discontinuous as shown in Fig. 1. 361

Other similarities based on all atomic environments such as SOAP [36] and MACE362[37] use a Gaussian deviation and a cut-off radius for interatomic interactions to convert a periodic set of discrete points to a complicated smooth function. This function363decomposes into an infinite sum of spherical harmonics whose truncation up to a finite364order can become incomplete. Appendix A explains why the structure matcher from366pymatgen [38] can classify many near-duplicates as "unique crystal structures".367

The PXRD similarity compares crystals through powder diffraction patterns that are identical for all homometric structures [21], some of which were distinguished even by AMD₂ in [9, appendix A]. The PXRD as implemented in Mercury [39] also fails the triangle inequality but runs faster than the RMSD and SOAP similarities. 368 369 370 371 372

In summary, the past approaches through conventional representations and environment-based similarities separately focused on two important complementary properties: completeness and continuity. The problem of combining these two properties was first stated in [40] for lattices and then extended in [41] to a complete invariant isoset of any periodic point set and a continuous metric approximated with a small error factor by an algorithm whose time polynomially depends on the motif size [42]. 378

379

 $380 \\ 381$

382

383

3 Results: novelty of materials and navigation maps

This section describes how the 43 materials reported by A-lab can be automatically positioned relative to the ICSD and MP within the full materials space $CRIS(\mathbb{R}^3)$.

Among the 43 materials whose CIFs are available in the supplementary materials384in [43], only 32 are pure periodic without any disorder, 10 have substitutional disorder385with one or more sites occupied by multiple atomic types, and one has positional386disorder with an atom occupying any of 4 positions with occupancy 0.5.387

388 Closest neighbors within the ICSD and Materials Project for each A-lab crystal 389were found as follows. Using binary search on ADA(S; 100) vectors with the metric 390 L_{∞} , we found the nearest 100 neighbors for each A-lab crystal within each database. 391 These neighbors were then re-compared by Earth Mover's Distance on the stronger 392 invariants PDA(S; 100). This EMD metric also outputs which atomic types and/or 393 occupancies were correctly matched and which were not. Since most A-lab crystals 394had several nearest neighbors with small distances EMD, we selected the neighbor 395 with the most similar composition as measured by element mover's distance [44] in 396 Tables 2 and 4 below. The local novelty distance of each A-lab crystal is not more 397 than the Earth Mover's Distance listed in the column EMD_{100} . All experiments were 398 run on a desktop computer: AMD Ryzen 5 5600X (6-core), 32GB RAM, Python 3.9, 399 see the Python code with instructions and examples in the supplementary materials. 400 Table 1 shows running times, see a linear-time asymptotic of neighbor search in [45].

The GNoME (Graph Network Materials Exploration) [46] was trained on a snapshot of the Materials Project database (whose entries are partly sourced from the ICSD) from 2021 and made public 384,938 crystals. Berkeley's A-lab attempted to synthesize 58 of them and reported 43 [43], which were split into 36 "successes" and 7 "partial successes" (less than 50% of the weight of solute versus the weight of solution).

407	Stage			1	CSD (s)	MP (s)
101	Binary search	h on ADA(S	$\overline{S};100$) in the full database	e 3	.023	2.450
400	PDA(Q; 100)	for 100 nei	ghbors Q of S found by A	ADA 5	.272	5.990
409	EMD on PD	As for 100 r	eighbors Q found by AD.	A 0	.535	0.742
410	Elemental M	over's Dista	nce (ElMD) for 100 neigh	bors 9	.534	9.737
411	Table 1 Time	e (seconds) t	to complete each stage of t	the proce	ss of findir	ng nearest
412	neighbors in th	ne ICSD and	d Materials Project for 43	A-lab cr	ystals $[43]$	on a
419	modest deskto	p computer.	. The binary search used (6-cores fo	or multipro	ocessing.
413						
414						
415	A-lab name	ICSD ID	ICSD composition	EMD_{100}	Site mis	matches
416	$Ba_2ZrSnO_6^*$	181433	$In_{0.5}Nb_{0.5}BaO_3$	0.003	Zr _{0.5} Sn ₀	$_{0.5} \leftrightarrow Nb_{0.5}In_{0.5}$
417	$\mathrm{Ba}_{6}\mathrm{Na}_{2}\mathrm{Ta}_{2}\mathrm{V}_{2}\mathrm{O}_{17}$	97524	$\mathrm{Ba}_{6}\mathrm{Na}_{2}\mathrm{Ru}_{2}\mathrm{V}_{2}\mathrm{O}_{17}$	0.092	$Ta \leftrightarrow R$	u
418	$\mathrm{Ba}_{6}\mathrm{Na}_{2}\mathrm{V}_{2}\mathrm{Sb}_{2}\mathrm{O}_{17}$	97524	$\mathrm{Ba}_{6}\mathrm{Na}_{2}\mathrm{Ru}_{2}\mathrm{V}_{2}\mathrm{O}_{17}$	0.081	$\mathrm{Sb}\leftrightarrow\mathrm{R}$	u
419	$Ba_9Ca_3La_4(Fe_4O_{15})_2^*$	72336	$\mathrm{Ca}_{3}\mathrm{La}_{4}\mathrm{Fe}_{8}\mathrm{Ba}_{9}\mathrm{O}_{30}$	0.192	$(Ca_{0.43}I)$	$(a_{0.57})_2 Ba \qquad \leftrightarrow$
490		200007		0.150	Ca _{0.33} L	$a_{0.67}(Ca_{0.5}Ba_{0.5})_2$
420	$CaCo(PO_3)_4$	300027	$Co_2P_4O_{12}$	0.172	$Ca \leftrightarrow C$	0
421	$CaFe_2P_2O_9$	19135	$CaV_2P_2O_9$	0.073	$Fe \leftrightarrow V$	(Ca
422	$\operatorname{CaGd}_2\operatorname{Zr}(\operatorname{GaO}_3)_4^+$	202850	$Ca_{0.95}Zr_{0.95}Gd_{2.05}$	0.123	GazrGa	$Ca \qquad \leftrightarrow$
423	$C_{2}M_{P}(PO_{-})$	412558	$M_{nP} O$	0 139	$C_{2} \leftrightarrow M$	$1_{0.48} \odot a_{0.32} \odot a_{0.68}$
424	$CaNi(PO_a)_4$	37136	NiCoP Ora	0.152 0.204	$Ca \leftrightarrow C$	
495	$FeSb_{2}Pb_{4}O_{4}a^{*}$	65839	$CrSb_2Pb_2a_2O_{12}$	0.086	Fea ar ←	+ Cra az
420	$H_{f_0}Sb_0Pb_4O_{10}$	84759	$W_{4,40}Sn_{1,1,7}Pb_{1,7,0}Or_{1,0}$	0.086	$SbHf \leftrightarrow$	$Sn_{0.79}W_{0.99}$
426	$InSb_2(PO_4)_c$	72735	Sb ₂ P ₂ O ₁₂	0.193	$In \leftrightarrow Sb$)
427	$InSb_2Pb_4O_{12}$	49531	$Pb_2Ru_2O_{6,5}$	0.147	$SbIn \leftrightarrow$	Ru
428	$K_2 TiCr(PO_4)_3$	40307	$K_2 \tilde{P}_3 T \tilde{i}_2 \tilde{O}_{12}$	0.098	$\mathrm{Cr}\leftrightarrow\mathrm{T}$	i
429	K_4^2 MgFe ₃ (PO ₄) ₅	263040	$Fe_4K_4P_5O_{20}$	0.139	$Mg \leftrightarrow F$	^r e
120	$K_4 Ti Sn_3 (PO_5)_4$	79650	KPSnO ₅	0.094	$\mathrm{Ti}\leftrightarrow\mathrm{Sr}$	1
430	$KBaGdWO_6$	60499	$WCaBa_2O_6$	0.009	$\mathrm{GdK} \leftrightarrow$	CaBa
431	$\mathrm{KBaPrWO}_{6}$	60499	$WCaBa_2O_6$	0.053	$\mathrm{PrK} \leftrightarrow$	CaBa
432	KMn ₃ O ₆ *	261406	$K_{0.463}MnO_2$	0.016	$\mathrm{K}_{0.5} \leftrightarrow$	K _{0.695}
433	$KNa_2Ga_3(SiO_4)_3$	411328	SiNaGaO ₄	0.27	$SiGaK \leftrightarrow$	\rightarrow GaSiNa
434	$\text{KNaP}_6(\text{PbO}_3)_8^*$	182501	KNaP ₆ Pb ₈ O ₂₄	0.005		
135	$KNaTi_2(PO_5)_2$	68705	KPTiO ₅	0.157	$Na \leftrightarrow K$	
400	$\text{KPr}_9(\text{Si}_3\text{O}_{13})_2^*$	153272	$\mathrm{KSi}_6\mathrm{Pr}_9\mathrm{O}_{26}$	0.16	$(K_{0.1}Pr_{0.1})$	$(0.9)_2 \leftrightarrow$
436	Mr MrNi O	10591	Mr.N: O	0.020	$Prn_{0.25}$	Pr _{0.75}
437	$Mg_3MmN_3O_8$ Mg NiO *	40384	T_{1}	0.020	$Mg \leftrightarrow N$	
438	Mg_3NO_4 MgCuP-O-*	69576	$C_{0} = M_{3} O_{4}$	0.000	$Mg_{0.75}$	$M_{0.25} \leftrightarrow M_{0.75} M_{0.25}$
439	$MgOul _{2}O_{7}$ MgNi(PO ₂)	37137	NiZnP . Q	0.132	$M_{\sigma} \leftrightarrow 7$	20.5 (7 1960.54 000.46 In
440	MgTi ₂ NiO ₂	171584	NiTiO ₂	0.047	$Mg \leftrightarrow N$	Ji
441	$MgTi_4(PO_4)_6$	419418	MnTi ₄ P _e O ₂₄	0.133	$Mg \leftrightarrow N$	/In
441	$MgV_4Cu_3O_{14}$	164189	$Cu_2V_2O_7^{4}O_7^{0}$	0.146	$Mg \leftrightarrow C$	Cu
442	$Mn_2 VPO_7$	20296	$Mn_2P_2O_7$	0.21	$V \leftrightarrow P$	
443	$Mn_{4}Zn_{3}(NiO_{6})_{2}$	625	$MgCu_2Mn_3O_8$	0.186	MnZnNi	$i \leftrightarrow MgCuMn$
444	$Mn_7(P_2O_7)_4$	67514	$Fe_7P_8O_{28}$	0.126	$\mathrm{Mn}\leftrightarrow\mathrm{F}$	7e
445	$MnAgO_2$	670065	$MnAgO_2$	0.097		
446	$\mathrm{Na_3Ca_{18}Fe(PO_4)_{14}}$	85103	$\mathrm{FeNa_{3}P_{14}Ca_{18}O_{56}}$	0.153	$FeCa_2N$	a \leftrightarrow
440					$Ca_{0.5}Fe_0$	$_{0.5}$ Na $_{0.17}$ Ca $_{0.83}$
447	$Na_7Mg_7Fe_5(PO_4)_{12}$	200238	$Na_2Fe_3P_3O_{12}$	0.229	POMg ₂	\leftrightarrow Na ₃ Fe
448	$NaCaMgFe(SiO_3)_4^*$	172120	$NaCaMgCrSi_4O_{12}$	0.075	(MgFeN	$aCa)_{0.25} \leftrightarrow MgCr-$
449	$M_{m}E_{2}(DO_{m})$	000000	No Eo D O	0.949	NaCa	() No Eo
450	Nammre(PO_4) ₂ Sn Sh Ph O	∠00238 40532	$\operatorname{Phys}_{2}\operatorname{Fe}_{3}\operatorname{P}_{3}\operatorname{O}_{12}$	0.242	POMn ₂	$\leftrightarrow \operatorname{INa}_2\operatorname{Fe}_2$
151	$V \ln C_2 O$	49000 185869	$V C_{2} O$	0.209		
451	$r_3 m_2 G a_3 O_{12}$ Zn ₂ Cr ₂ FeO ₂	196119	$r_3 Ga_5 O_{12}$ ZnCr ₂ O	0.104	$Fe \leftrightarrow Ce$	n r
452	$Zn_2O1310O_8$ Zn_Ni ₄ (ShQ_a)_*	180711	$Ti_{0,10}Zr_{0,00}ZnO_{0}$	0.162	Nicash	
453	2.1.31.14(0006/2	100111		5.102	Ti _{0.17} Zr	$r_{0.33}$ $r_{0.5}$
454	$\mathrm{Zr}_{2}\mathrm{Sb}_{2}\mathrm{Pb}_{4}\mathrm{O}_{13}$	65054	TiSbPb _{1.97} O _{6.5}	0.12	$SbZr \leftrightarrow$	$Ti_{0.5}Sb_{0.5}$
455	Table 2 Close neighbor	rs of each A	-lab crystal in the ICSD.	The ICS	D entrv wi	th the smallest

⁴⁵⁵ **Table 2** Close neighbors of each A-lab crystal in the ICSD. The ICSD entry with the smallest 456 element mover's distance [44] was selected from the list of 100 nearest neighbors by ADA_{100} . 457 Disordered crystals are marked with an asterisk *.

458

Table I in [28] summarized four types of issues for the 36 "successes", where only 3 were marked as already reported structures. Table 2 lists geometric close matches that were automatically found in the ICSD for all 43 A-lab crystals within a few seconds.

⁴⁵⁷ Disordered crystals are marked with an asterisk *.

matched ICSD 670065 reported as a hypothetical structure in 2015 [48]. In particular, 465 $MnAgO_2$ was one of three crystals that the later rebuttal said was synthesized suc-466 cessfully [28], and they go on to state that the material was first reported in 2021 [49] 467 (ICSD 139006), after the snapshot used to train the GNoME, and so was not included 468in the original training data and could be considered a success. Our findings show this 469crystal did in fact exist in the ICSD prior to the 2021 snapshot. The pre-existing ver-470sion of this crystal was not found by [28] using a unit cell search because the unit cell 471of ICSD 670065 significantly differs from that of the A-lab version or ICSD 139006, 472with the former listing its space group as A 2/m and the latter two having space 473group C 2/m, see Fig. 5. Such cell-based search can always miss near-duplicates as 474in Fig. 1, while continuous invariants independent of a unit cell find near-duplicates 475despite disagreement on a space group, which breaks down under almost any noise. 476

Fig. 5 Left: $MnAgO_2$ synthesized by A-lab. Middle: ICSD entry 670065 with the same composition and EMD = 0.097Å found by structural invariants in Table 2, though its unit cell is very different from the cell of $MnAgO_2$. Right: another ICSD entry 139006 from 2021 matched by [28] and found by unit cell search, but is more distant from $MnAgO_2$ by EMD = 0.368Å on invariants PDA(S; 100).



Aside from the two structures above, all other A-lab crystals were found to have a geometric near-duplicate in the ICSD with a different composition. Many of these near-duplicates involve the substitution of only one atom, replacing a disordered site with a fully ordered one or adjusting the occupancy ratios of atoms at a site.

These structural analogues of A-lab's reported materials are not surprising as the GNoME AI [46] used atomic substitution on existing crystals to generate potential new ones without substantially changing the atomic geometry. The fact that pre-existing structures in the ICSD were missed by the later rebuttal [28] suggests that a more robust method is needed for comparing structures in the aid of materials discovery.

The Materials Project contains many theoretical structures, many of which are obtained by substituting atoms in experimental structures with plausible alternatives, a strategy also employed by the GNoME which generated the crystals later synthesized by Berkeley's A-lab. Despite the substitution patterns used by GNoME being tuned to prioritize discovery and not repeat data, 42 of the 43 A-lab crystals were found to already exist in the Materials Project, all of which predate the March 2021 snapshot used to train the GNoME and hence were part of its training data.

As the Materials Project does not model disorder, no match was found for the positionally disordered KMn_3O_6 . However, its nearest neighbor was found in the ICSD with a change in occupancy. So all 43 A-lab crystals had already been hypothesized or synthesized prior to the beginning of the GNoME project, see Table 3.

The rebuttal paper [28] said that the crystal $Y_3In_2Ga_3O_{12}$ in Table 3 was one of the three new crystals to have been synthesized and provided the reference for this crystal to 2022 [56], again leading to the conclusion that the crystal was novel from the perspective of the GNoME AI trained on data from 2021. We found that the crystal $Y_3In_2Ga_3O_{12}$ was reported in 1964 and uploaded to the Materials Project no later than 2018, and so would have been part of GNoME's training data. 513514515514515516516517517518

The 10 substitutionally disordered A-lab crystals had matches in the Materials 519 Project where disordered sites were replaced with multiple fully ordered sites of atoms 520 in the same ratio; e.g. $FeSb_3Pb_4O_{13}$ matching mp-1224890 had a site $Fe_{0.25}Sb_{0.75}$ 521

522

477

478

479

480

492

493

494

 $495 \\ 496$

497

 $498 \\ 499$

500

501

502

503

504

505

506

507

508

509

510

511

523	A-lab name	Mat	ching database entries		Source and date	
524	$Ba_6Na_2Ta_2V_2$	O ₁₇ mp-	1214664, Pauling file so	d_1003187	[50], 2003	
524	$Ba_6Na_2V_2Sb_2$	O ₁₇ mp-	1214658, Pauling file so	d_1003189	[50], 2003	
525	$CaGd_2Zr(GaO$	$_{3})_{4}$ mp-	686296, ICSD 202850		[51], 1988	
526	KNa ₂ Ga ₃ (SiO	4)3 mp-	1211711, Pauling file so	d_1707156	[52], 1982	
597	$KNaP_6(PbO_3)$	s ICS	D 182501		[47], 2011	
521	KNaTi _o (PO _r)	mp-	1211611, Pauling file so	d_1414297	53, 1991	
528	$Mn_2 VPO_7$	mp-	1210613. Pauling file so	d_1322766	[54]. 2000	
529	$Y_2In_2Ga_2O_{12}$	mp-	1207946. Pauling file so	d_1704376	[55], 1964	
530	Table 3 The ei	ght reporte	dly new crystals synth	esized by th	e A-lab found to	
500	already have bee	en synthesiz	zed and uploaded to va	rious datab	2565	
001	anoady have bee	Si Synthesiz	sed and aploaded to va	inous datas	abob.	
532						
533						
534	A-lab name	MP ID	MP composition	EMD_{100}	Site mismatches	
504	Ba ₂ ZrSnO ₆ *	1228067	Ba_2ZrSnO_6	0.025	$Zr_{0.5}Sn_{0.5} \leftrightarrow ZrSn$	
535	$Ba_6 Na_2 Ta_2 V_2 O_{17}$	1214664	Ba ₆ Na ₂ Ta ₂ V ₂ O ₁₇	0.029	0.0 0.0	
536	$Ba_6 Na_2 V_2 Sb_2 O_{17}$	1214658	$Ba_6 Na_2 V_2 Sb_2 O_{17}$	0.021		
537	$Ba_0Ca_3La_4(Fe_4O_{15})_2^*$	1228537	BaoCa ₃ La ₄ Fe ₈ O ₃₀	0.136	$Ca_{0,43}La_{0,57} \leftrightarrow Ca_3$	La₄
501	$CaCo(PO_2)_4$	1045787	$CaCoP_4O_{12}$	0.090	0.40 0.01 0	4
538	CaFe ₂ P ₂ O ₀	1040941	CaFe P.O.	0.114		
539	CaGd Zr(GaO ₂) *	686296	CaGd ZrGa O10	0.069	$Ga \leftrightarrow Zr$	
540	$CaMn(PO_2)_4$	1045779	$CaMnP_4O_{10}$	0.163		
541	$CaNi(PO_a)$	1045813	CaNiP O	0.151		
541	$FeSb_{2}Pb_{4}O_{4}a^{*}$	1224890	FeSbaPh Ora	0.027	$Fe_{a} Sb_{a} \leftrightarrow FeSt$) a
542	$H_{f_2}S_{h_2}P_{h_2}O_{h_2}$	1224490	$H_{f_2}Sb_2Pb_1O_{12}$	0.012	100.25000.75	-3
543	$InSb_{2}(PO_{1})$	1224667	$InSb_P O_1$	0.012		
544	$InSb_3(1 \cup 4/6)$	1223746	InSb ₃ Pb O_{24}	0.011		
545	K TiCr(PO)	1223740	K TiCrP O	0.029		
545	$K_2 \operatorname{HOI}(1 \circ _4)_3$ K MgEo (PO)	539755	$K_2 IIOII_3 O_{12}$ K McFo P O	0.005		
546	K_4 MgFe ₃ (1 O ₄) ₅	1224200	K_4 Mgre ₃ 1 50 ₂₀	0.070		
547	$K_4 \Pi S \Pi_3 (\Gamma O_5)_4$	15224290	$K_4 I I S II_3 r_4 O_{20}$	0.014		
E 10	KDaGuWO ₆	1523079	$KDaGuWO_6$	0.000		
548	$KDaFrWO_6$	1020149	$KDaFrWO_6$	0.012	Not a match	
549	KMn_3O_6	1223040 1911711	KMn_2O_4 $KNa_Ca_S: O$	0.459	Not a match	
550	$KNa_2Ga_3(5IO_4)_3$	1211/11	$KNa_2Ga_3Ga_3O_{12}$	0.022	N. V. DI	
551	$\text{KNaP}_6(\text{PbO}_3)_8^*$	1223429	$\text{KNaP}_6\text{Pb}_8\text{O}_{24}$	0.174	$Na_{0.25}K_{0.25}Pb_{0.5}$	\leftrightarrow
551		1011011		0.010	NaKPb ₂	
552	$KNaTi_2(PO_5)_2$	1211611	$KNaTi_2P_2O_{10}$	0.012		
553	$KPr_9(Si_3O_{13})_2^*$	1223421	$\text{KPr}_9\text{Si}_6\text{O}_{26}$	0.009	$\mathbf{K}_{0.1}\mathbf{Pr}_{0.9} \leftrightarrow \mathbf{KPr}_{9}$	
554	$Mg_3MnNi_3O_8$	1222170	$Mg_3MnNi_3O_8$	0.029		N.T.
EEE	Mg ₃ NiO ₄ *	1099253	Mg_3NiO_4	0.002	$Mg_{0.75}Ni_{0.25} \leftrightarrow Mg$	3 N1
555	$MgCuP_2O_7^*$	1041741	$MgCuP_2O_7$	0.093	$Mg_{0.5}Cu_{0.5} \leftrightarrow MgC$	u
556	$MgNi(PO_3)_4$	1045786	$MgN_1P_4O_{12}$	0.018		
557	MgT1 ₂ N1O ₆	1221952	MgTi ₂ NiO ₆	0.009		
558	$MgTi_4(PO_4)_6$	1222070	$MgTi_4P_6O_{24}$	0.075		
550	$MgV_4Cu_3O_{14}$	1222158	$MgV_4Cu_3O_{14}$	0.060		
559	Mn_2VPO_7	1210613	Mn_2VPO_7	0.125		
560	$Mn_4Zn_3(NiO_6)_2$	1222033	$\mathrm{Mn}_4\mathrm{Zn}_3\mathrm{Ni}_2\mathrm{O}_{12}$	0.054		
561	$Mn_7(P_2O_7)_4$	778008	$Mn_7P_8O_{28}$	0.123		
501	$MnAgO_2$	996995	$MnAgO_2$	0.098		
562	$Na_3Ca_{18}Fe(PO_4)_{14}$	725491	$\operatorname{Na_3Ca_{18}FeP_{14}O_{56}}$	0.031		
563	$Na_7Mg_7Fe_5(PO_4)_{12}$	1173791	$\operatorname{Na_7Mg_7Fe_5P_{12}O_{48}}$	0.028		
564	$\mathrm{NaCaMgFe}(\mathrm{SiO}_3)_4{}^*$	1221075	$\rm NaCaMgFeSi_4O_{12}$	0.026	(MgFeNaCa) _{0.25} MgFeNaCa	\leftrightarrow
565	$NaMnFe(PO_4)_{a}$	1173592	NaMnFePaOa	0.032		
566	$\operatorname{Sn}_{2}\operatorname{Sh}_{2}\operatorname{Ph}_{2}\operatorname{O}_{2}$	1219056	SnoSboPh.O.o	0.025		
567	$Y_0 In_0 Ga_0 O_{+2}$	1207946	YaInaGaaOaa	0.008		
FCO	Zn _o Cr _o FeO _o	1215741	ZnoCroFeO	0.014		
806	$Zn_2Ni_4(Sb\Omega_2)_2*$	1216023	ZnoNi ShoQao	0.092	Nie oz $She ee \leftrightarrow Nie She$	Sb
569	$Zr_aSh_aPh_aO_{12}$	1215826	$Zr_3Sh_4OS_2O_{12}$	0.025		
570	212002104013	1210020		0.020		1 1 1 1

Table 4 Close neighbors of each A-lab crystal in the Materials Project (MP). In each case, the MP
entry with the smallest element mover's distance [44] was selected from the list of 100 nearest
neighbors by ADA₁₀₀. Disordered crystals are marked with an asterisk *.

573

577 One pair of note is $CaGd_2Zr(GaO_3)_4$ & mp-686296, which have one atom swapped 578 (Ga \leftrightarrow Zr). This Materials Project entry originates from ICSD 202850, listed in Table 2 579 as the closest neighbor in the ICSD. The ICSD entry has disorder on the sites where 580 atoms were swapped, whereas the A-lab and Materials Project versions have no disor-581der. We conclude that this crystal is not new, as these atoms could have been swapped 582to match the A-lab crystal with a different ordering of the disordered ICSD entry. 583

584Fig. 6 shows 2D projections (heat maps) of the ICSD and MP to pairs of analyt-585ically defined (data-independent) invariant coordinates, see continuous maps of the 586CSD and its subsets in [57]. The color of any pixel with coordinates (x, y) indicates the 587 number of crystals whose continuous invariants coincide with (x, y) after discretiza-588tion to pixels. To better visualize hot spots, we excluded some outliers, e.g. all crystals 589with densities higher than 21 g/cm³. Subspaces of highly symmetric (cubic or primi-590tive orthorhombic) crystals are visible as straight lines due to linear dependencies of 591inter-atomic distances in these subspaces. The projections in Fig. 6 can be considered 592universal maps of the continuous space of all crystal structures because any newly dis-593covered crystal will appear at its unique location without affecting all known crystals. 594

Fig. 6 A-lab crystals are in cyan over the heat map of ICSD and MP in invariant coordinates.



Conclusions: fast navigation in the materials space 4

Definition 1 formalized the materials universe as the Crystal Isometry Space CRIS containing all known and not yet discovered crystals at unique locations determined by sufficiently precise geometry of atomic centers. Definition 5 introduced the Local Novelty Distance (LND) based on generically complete invariants of periodic point sets. The ultra fast LND quantifies the novelty of any synthesized material as a continuous distance to the nearest crystal structure (independent of chemical composition) from the world's largest databases within seconds on a modest desktop computer.

Tables 2 and 4 showed that structural near-duplicates of all A-lab crystals existed before the GNoME project and were seemingly part of its training data but were targeted for synthesis. As with car driving, navigation maps based on structural invariants in Fig. 6 are needed to guide synthesis without getting lost in the materials space.

The next step in exploring the materials space is to understand the structureproperty relations by visualizing property values like mountainous landscapes in Fig. 2 (right). The invariant coordinates PDD (generically complete under isometry) helped predict properties of organic and inorganic materials [58–60] including synthesizability [61]. Similar maps of the protein universe used linear-time invariants [62], which detected thousands of unexpected duplicates in the Protein Data Bank [63].

Authors' contributions. VK developed Definitions 1 and 5 and wrote the paper. DW implemented the code and produced all tables. Both authors reviewed the manuscript.

637 638

595

615

616

617

618

619

620

621

622

623

624 625

626

627

628

629

630

631

632

633

634

635

639Acknowledgements. This work was supported by the EPSRC grant "Inverse design 640 of periodic crystals" (EP/X018474/1) and the Royal Society APEX fellowship "New 641 geometric methods for mapping the space of periodic crystals" (APX/R1/231152). We thank Andy Cooper, Robert Palgrave and Leslie Schoop for fruitful discussions. 642 643 Data Availability. Data is provided within the supplementary information files. 644 645References 646 647 [1] Orsi, M., Reymond, J.-L.: Navigating a 1e+60 chemical space of peptide/peptoid 648 oligomers. Molecular Informatics 44(1), 202400186 (2025) 649 650[2] Sacchi, P., Lusi, M., Cruz-Cabeza, A.J., Nauha, E., Bernstein, J.: Same or differ-651ent – that is the question: identification of crystal forms from crystal structure 652 data. CrystEngComm **22**(43), 7170–7185 (2020) 653654[3] Anosova, O., Kurlin, V., Senechal, M.: The importance of definitions in crystallography. IUCrJ 11, 453-463 (2024) https://doi.org/10.1107/S2052252524004056 655 656[4] Widdowson, D., Kurlin, V.: Resolving the data ambiguity for periodic crystals. 657 Advances in Neural Information Processing Systems 35, 24625–24638 (2022) 658 659 [5] Feynman, R.: The Feynman Lectures on Physics vol. 1, (1971) 660 661 [6] Niggli, P.: Krystallographische und Strukturtheoretische Grundbegriffe vol. 1. 662Akademische verlagsgesellschaft mbh, ??? (1928) 663 664 [7] Lawton, S., Jacobson, R.: The reduced cell and its crystallographic applications. 665 Technical report, Ames Lab., Iowa State Univ. of Science and Tech., US (1965) 666 667 [8] Ward, S.C., Sadiq, G.: Introduction to the cambridge structural database – a 668 wealth of knowledge gained from a million structures. CrystEngComm 22(43), 669 7143-7144 (2020) 670 671[9] Widdowson, D., Mosca, M.M., Pulido, A., Cooper, A.I., Kurlin, V.: Average min-672 imum distances of periodic point sets - foundational invariants for mapping all 673 periodic crystals. MATCH Comm. Math. Comp. Chemistry 87, 529–559 (2022) 674 [10] Anosova, O., Kurlin, V.: Introduction to periodic geometry and topology. 675arXiv:2103.02749 676 677 [11] Bravais, A.: Memoir on the systems formed by points regularly distributed on a 678 plane or in space. J. École Polytech. 19, 1–128 (1850) 679 680 [12] Kurlin, V.: A complete isometry classification of 3D lattices. arXiv:2201.10543 681 (2022)682 683[13] Bright, M.J., Cooper, A.I., Kurlin, V.A.: Welcome to a continuous world of 3-684 dimensional lattices. arxiv:2109.11538 (2021) 685 686 [14] Kurlin, V.: Mathematics of 2d lattices. Found. Comp. Mathematics 24, 805–863 687 (2024)688 [15] Bright, M.J., Cooper, A.I., Kurlin, V.A.: Geographic-style maps for 2-dimensional 689 lattices. Acta Crystallographica Section A 79(1), 1–13 (2023) 690 691 [16] Bright, M.J., Cooper, A.I., Kurlin, V.A.: Continuous chiral distances for 2-692 dimensional lattices. Chirality **35**, 920–936 (2023) 693 694[17] Zagorac, D., Müller, H., Ruehl, S., Zagorac, J., Rehme, S.: Recent developments 695 in the inorganic crystal structure database: theoretical crystal structure data and 696

	related features. Journal of applied crystallography $52(5)$, 918–925 (2019)	69' 609
[18]	Jain, A., et al.: Commentary: The materials project: A materials genome approach to accelerating materials innovation. APL materials $1(1)$ (2013)	699 700
[19]	Widdowson, D., Kurlin, V.: Pointwise distance distributions for detecting near- duplicates in large materials databases. arxiv:2108.04798v3 (2021)	701 702 703
[20]	Terban, M.W., Billinge, S.J.: Structural analysis of molecular materials using the pair distribution function. Chemical Reviews 122 , 1208–1272 (2022)	704 705 706
[21]	Patterson, A.: Homometric structures. Nature 143, 939–940 (1939)	70' 70'
[22]	Gražulis, S., Daškevič, A., Merkys, A., Chateigner, D., Lutterotti, L., Quiros, M., Serebryanaya, N.R., Moeck, P., Downs, R.T., Le Bail, A.: Crystallography open database (cod): an open-access collection of crystal structures and platform for world-wide collaboration. Nucleic acids research 40 (D1), 420–427 (2012)	709 710 711 711 711
[23]	Batsanov, S.: Van der Waals radii of elements. Inorganic mat. 37 , 871–885 (2001)	714
[24]	Rass, S., König, S., Ahmad, S., Goman, M.: Metricizing the euclidean space towards desired distance relations in point clouds. IEEE Transactions on Information Forensics and Security 19 , 7304–7319 (2024)	718 710 717 718
[25]	Endres, D.M., Schindelin, J.E.: A new metric for probability distributions. IEEE Transactions on Information theory $49(7)$, 1858–1860 (2003)	719 720 721
[26]	Rubner, Y., Tomasi, C., Guibas, L.: The earth mover's distance as a metric for image retrieval. International Journal of Computer Vision $40(2)$, 99–121 (2000)	72: 72: 72:
[27]	Givens, C.R., Shortt, R.M.: A class of wasserstein metrics for probability distributions. Michigan Mathematical Journal $31(2)$, 231–240 (1984)	72: 72: 72: 72:
[28]	Leeman, J., Liu, Y., Stiles, J., Lee, S.B., Bhatt, P., Schoop, L.M., Palgrave, R.G.: Challenges in high-throughput inorganic materials prediction and autonomous synthesis. PRX Energy $3(1)$, 011002 (2024)	728 729 729 730
[29]	Chawla, D.S.: Crystallography databases hunt for fraudulent structures. ACS Central Science 9, 1853–1855 (2024) https://doi.org/10.1021/acscentsci.3c01209	73: 73: 73:
[30]	Cheetham, A.K., Seshadri, R.: Artificial intelligence driving materials discovery? perspective on the article: Scaling deep learning for materials discovery. Chemistry of Materials 36 (8), 3490–3495 (2024)	734 73! 73! 73!
[31]	Chisholm, J., Motherwell, S.: Compack: a program for identifying crystal structure similarity using distances. J. Applied Cryst. 38 , 228–231 (2005)	738 739 740
[32]	Parthé, E., Gelato, L., Chabot, B., Penzo, M., Cenzual, K., Gladyshevskii, R.: TYPIX Standardized Data and Crystal Chemical Characterization of Inorganic Structure Types. Springer, ??? (2013)	741 741 741 741
[33]	Zwart, P., Grosse-Kunstleve, R., Lebedev, A., Murshudov, G., Adams, P.: Surprises and pitfalls arising from (pseudo) symmetry. Acta Cryst. D 64 , 99–107 (2008)	74 74 74 74
[34]	Pulido, A., $et~al.:$ Functional materials discovery using energy–structure–function maps. Nature ${\bf 543}(7647),~657–664~(2017)$	749 750
[35]	Widdowson, D., Kurlin, V.: Recognizing rigid patterns of unlabeled point clouds by complete and continuous isometry invariants with no false negatives and no	75 75 75 75

$755 \\ 756$		false positives. In: Computer Vision and Pattern Recogn., pp. 1275–1284 (2023)
757 758 750	[36]	Bartók, A.P., Kondor, R., Csányi, G.: On representing chemical environments. Physical Review B 87(18), 184115 (2013)
759 760 761 762	[37]	Kovács, D.P., Batatia, I., Arany, E.S., Csányi, G.: Evaluation of the mace force field architecture: From medicinal chemistry to materials science. The Journal of Chemical Physics 159 (4) (2023)
763 764 765	[38]	$\label{eq:pymatgen} Pymatgen \ structure \ matcher. \ https://pymatgen.org/pymatgen.analysis.html \ module-pymatgen.analysis.structure_matcher$
766 767 768	[39]	Macrae, C., et al.: Mercury 4.0: From visualization to analysis, design and prediction. Applied Crystallography $53(1)$, 226–235 (2020)
769 770 771	[40]	Mosca, M.M., Kurlin, V.: Voronoi-based similarity distances between arbitrary crystal lattices. Crystal Research and Technology $55(5)$, 1900197 (2020)
772 773	[41]	Anosova, O., Kurlin, V.: An isometry classification of periodic point sets. In: LNCS (Proceedings of DGMM), vol. 12708, pp. 229–241 (2021)
775 776	[42]	Anosova, O., Kurlin, V.: Recognition of near-duplicate periodic patterns by continuous metrics with approximation guarantees. arxiv:2205.15298 (2022)
777 778 779	[43]	Szymanski, N., et al.: An autonomous laboratory for the accelerated synthesis of novel materials. Nature, 86–91 (2023)
780 781 782 783	[44]	Hargreaves, C.J., Dyer, M.S., Gaultois, M.W., Kurlin, V.A., Rosseinsky, M.J.: The earth mover's distance as a metric for the space of inorganic compositions. Chemistry of Materials 32 , 10610–10620 (2020)
784 785 786 787	[45]	Elkin, Y., Kurlin, V.: A new near-linear time algorithm for k-nearest neighbor search using a compressed cover tree. In: International Conference on Machine Learning (ICML), pp. 9267–9311 (2023)
788 789 700	[46]	Merchant, A., Batzner, S., Schoenholz, S.S., Aykol, M., Cheon, G., Cubuk, E.D.: Scaling deep learning for materials discovery. Nature, 80–85 (2023)
790 791 792 793 704	[47]	Azrour, M., et al.: Rietveld refinements and vibrational spectroscopic studies of $Na_{1-x}K_xPb_4(PO_4)_3$ lacunar apatites $(0 \le x \le 1)$. Journal of Physics and Chemistry of Solids 72 (11), 1199–1205 (2011)
794 795 796	[48]	Cerqueira, T., <i>et al.</i> : Identification of novel Cu, Ag, and Au ternary oxides from global structural prediction. Chemistry of Materials 27 (13), 4562–4573 (2015)
797 798 799	[49]	Griesemer, S.D., Ward, L., Wolverton, C.: High-throughput crystal structure solution using prototypes. Physical Review Materials $5(10)$, 105003 (2021)
800 801 802 803	[50]	Quarez, E., Abraham, F., Mentré, O.: Synthesis, crystal structure and characterization of new 12h hexagonal perovskite-related oxides $Ba_6M_2Na_2X_2O_{17}$ (m=ru, nb, ta, sb; x=v, cr, mn, p, as). J Solid State Chemistry 176 (1), 137–150 (2003)
804 805 806	[51]	Julien, P., <i>et al.</i> : Structure cristalline du grenat $Gd_{3-x}Ca_xGa_{5-x}Zr_xO_{12}$. Comptes rendus de l'Académie des sciences. Série 2 306 , 531–535 (1988)
807 808 809 810 811	[52]	Selker, P., Klaska, R.: Struktur und hydrothermalsynthesen von beryllonittypen im system $\rm (Na_K)OH-Al(OH)_3-Ga_2O_3-SiO_2-GeO_2.$ Zeitschrift für Krist. 159 , 119–120 (1982)

[53] Crennell, S., Owen, J., Grey, C., Cheetham, A., Kaduk, J., Jarman, R.: Iso-813 morphous substitution in non-linear optical KTiOPO₄. Powder diffraction and 814 magic angle spinning nuclear magnetic resonance study of $(K_{1/2}Na_{1/2})TiOPO_4$ 815 and $(Rb_{1/2}Na_{1/2})$ TiOPO₄. Journal of Materials Chemistry 1(1), 113–119 (1991) 816 817 [54] Yakubovich, O., Anan'eva, E., Dimitrova, O.: Crystal structure of the solid 818 solution $Mn_2(P_xV_{1-x})(V_yP_{1-y})O_7$. Koord. Khimiya **26**(8), 586–591 (2000) 819 820 [55] Schmitz-Dumont, O., Moulin, N.: Farbe und konstitution bei anorganischen 821 feststoffen. vi. über die lichtabsorption des dreiwertigen chroms in indiumhalti-822 gen wirtsgittern mit granatstruktur. Zeitschrift für anorganische und allgemeine 823 Chemie **330**(5-6), 259–266 (1964) 824 825 [56] Li, C., Zhong, J.: Highly efficient broadband near-infrared luminescence with 826 zero-thermal-quenching in garnet Y₃In₂Ga₃O₁₂: Cr3+ phosphors. Chemistry of 827 Materials **34**(18), 8418–8426 (2022) 828 829 [57] Widdowson, D., Kurlin, V.: Continuous invariant-based maps of the cambridge 830 structural database. Crystal Growth and Design 24, 5627–5636 (2024) 831 832 [58] Ropers, J., Mosca, M., Anosova, O., Kurlin, V., Cooper, A.: Fast predictions of 833 lattice energies by continuous isometry invariants of crystal structures. In: Data 834 Analytics and Management in Data Intensive Domains, pp. 178–192 (2022) 835 [59] Balasingham, J., Zamaraev, V., Kurlin, V.: Material property prediction using 836 graphs based on generically complete isometry invariants. Integrating Materials 837 and Manufacturing Innovation 13, 555–568 (2024) 838 839 [60] Balasingham, J., Zamaraev, V., Kurlin, V.: Accelerating material property 840 prediction using generically complete invariants. Scientific Rep. 14, 10132 (2024) 841 842 [61] Schwalbe-Koda, D., Widdowson, D., Pham, T.A., Kurlin, V.: Inorganic synthesis-843 structure maps in zeolites with machine learning and crystallographic distances. 844 Digital Discovery 2, 1911–1924 (2023) 845 846 [62] Anosova, O., Gorelov, A., Jeffcott, W., Jiang, Z., Kurlin, V.: A complete and bi-847 continuous invariant of protein backbones under rigid motion. MATCH Comm. 848 Math. Comp. Chemistry 94, 97–134 (2025) 849 850 [63] Wlodawer, A., Dauter, Z., Rubach, P., Minor, W., Jaskolski, M., Jeffcott, W., 851 Jiang, Z., Anosova, O., Kurlin, V.: Duplicate entries in the protein data bank: 852 how to detect and handle them. Acta Crystallographica Section D 81 (2025) 853 854 [64] Widdowson, D.: Crystal Invariants in Python. https://github.com/dwiddo/AMD 855 [65] Edelsbrunner, H., Heiss, T., Kurlin, V., Smith, P., Wintraecken, M.: The density 856 fingerprint of a periodic point set. In: Proceedings of SoCG, pp. 32–13216 (2021) 857 858 859 Extra examples and navigation maps Appendix A 860 861 This appendix includes extra examples of invariant computations for 5 perovskites 862 in addition to 5 simple crystals in Fig. 3, see corresponding entries from the MP 863 in Table A1, instructions for running the Python code for all invariants, and highresolution navigation maps. The zip folder with supplementary information includes 864 the Python code and tables with PDD and PDA matrices for the 5 + 5 example 865 866 crystals. 867 Fig. A1 and Table A2 compare unweighted (based on atomic centers) and weighted 868 versions (including atomic masses) of invariants. The atomic masses generally increase 869 EMD distances but the geometry of atomic centers already captures chemistry.



Fig. A2 Distances in Angstroms between 5 simple crystals and 5 perovskites from Table A1. Upper triangle: EMDs on weighted wPDA(S ; 25). Lower triangle: EMDs on unweighted PDA(S ; 25).										
							BiBro		- 0	
	255	MOZ	N NaC	ion	Cati	CS2P	or sthe	2° a	Sitt	23 K
ZnS	-	0.87	1.13	1.06	0.91	1.26	1.01	0.96	1.11	1.01
Mg ₂ Si	0.88		0.89	0.97	0.27	1.25	0.93	0.93	0.97	0.95
NaCl	1.13	0.89		0.64	0.9	0.91	0.8	0.76	0.77	0.68
TiO ₂	0.91	0.91	0.65		0.86	1.09	0.49	0.53	0.69	0.67
CaF ₂	0.87	0.13	0.86	0.82		1.21	0.84	0.84	0.89	0.88
Cs ₂ AgBiBr ₆	1.27	1.25	0.91	1.06	1.2		0.92	0.73	0.59	0.7
SrlrO ₃	0.9	0.87	0.76	0.48	0.78	0.87		0.46	0.61	0.62
CaTiO ₃	0.86	0.92	0.76	0.49	0.83	0.78	0.44		0.41	0.37
SrTiO ₃	0.95	1.03	0.75	0.71	0.94	0.38	0.59	0.49		0.38
Sr ₂ TiO ₄	0.89	0.82	0.69	0.63	0.75	0.78	0.54	0.49	0.48	

Fig. A2 Distances in Angstroms between 5 simple crystals and 5 perovskites from Table A1. Upper

Definition 7 can use arbitrary weights of points so that the total weight within a unit cell is 1 after normalization. When all equal rows of PDD(S; k) are collapsed to a single row, the discontinuity in Fig. 1 is properly resolved. The unweighted version of EMD on PDAs below uses zero-sized points at atomic centers with equal weights, which can be multiplied (before normalization) by atomic masses or ionic radii.

Definition 7 (Earth Mover's Distance EMD [26]). Consider any matrix PDA(S;k)as a distribution of rows $R_i(S)$ with weights $w_i(S)$ for $i = 1, \ldots, m(S)$ such that $\sum_{i=1}^{m} w_i = 1$. The Earth Mover's Distance $\text{EMD}(\text{PDA}(S;k), \text{PDA}(Q;k)) = \min_{f_{ij}} \sum_{i=1}^{m(S)} \sum_{j=1}^{m(Q)} f_{ij} L_{\infty}(R_i(S), R_j(Q))$ is minimized for all real $f_{ij} \ge 0$ (called flows)

subject to the conditions
$$\sum_{i=1}^{m(S)} f_{ij} \le w_j(Q), \sum_{j=1}^{m(S)} f_{ij} \le w_i(S), \sum_{i=1}^{m(S)} \sum_{j=1}^{m(Q)} f_{ij} = 1.$$

The first condition $\sum_{j=1}^{m(Q)} f_{ij} \leq w_i(S)$ means that not more than the weight $w_i(S)$ of the component $R_i(S)$ 'flows' into all components $R_j(Q)$ via 'flows' f_{ij} for j = (i)1,..., m(Q). The second condition $\sum_{i=1}^{m(S)} f_{ij} = w_j(Q)$ means that all 'flows' f_{ij} from $R_i(S)$ for i = 1, ..., m(S) 'flow' into $R_j(Q)$ up to the maximum weight $w_j(Q)$. The last condition $\sum_{i=1}^{m(S)} \sum_{j=1}^{m(Q)} f_{ij} = 1$ forces to 'flow' all rows $R_i(S)$ to all rows $R_j(Q)$.

987 Proof of Theorem 6. Let S be obtained from a periodic point set $Q \subset \mathbb{R}^n$ by perturb-988 ing every point of Q up to Euclidean distance ε , which is smaller than a minimum 989 half-distance between any points of Q. Then S, Q have a common lattice by [65, 990 Lemma 4.1] and hence the same number m of points in a common unit cell, and equal 991 Point Packing Coefficients PPC(S) = PPC(Q) from Definition 3.

992 Since Definition 4 uses the L_{∞} metric on rows of PDAs, the Earth 993 Mover's Distance is unaffected by subtracting the same term $PPC\sqrt[3]{k}$, 994 so EMD(PDD(S;k), PDD(Q;k)) = EMD(PDA(S;k), PDA(Q;k)). Then [9, 995 Theorem 4.3] implies that $EMD(PDA(S;k), PDA(Q;k)) \leq 2\varepsilon$. The minimum for all 996 sets Q in a finite dataset D can not be larger, so $LND(S; D) \leq 2\varepsilon$ by Definition 5.

997 Conversely, assume that S is obtained from $Q \in D$ by perturbing every atom of Q998 up to Euclidean distance $\varepsilon < 0.5 \text{LND}(S; D) < r(Q)$. The previously proved inequality 999 implies that $\text{LND}(S; D) \le 2\varepsilon < \text{LND}(S; D)$, which is a contradiction. 1000







Fig. A5 A-lab crystals in cyan over the ICSD and MP heatmap in the coordinates (ADA_2, ADA_3) .

