# Average Minimum Distances continuously quantify similarities between geometries of any crystal structures

Marco Mosca[a] and Vitaliy Kurlin [a]*

[a]*Materials Innovation Factory, University of Liverpool, UK.*

*E-mail: vitaliy.kurlin@gmail.com*

**crystal structure prediction, isometry classification, distance-based invariants, homometric crystals**

## Abstract

Modern tools of Crystal Structure Prediction output thousands of simulated structures, though only few of them can be really synthesized. This embarrassment of over-prediction can be resolved only if crystals are compared for similarity by invariants that are independent of crystal representations and preserved by rigid motions. To continuously quantify a similarity between crystals with the same chemical composition, these invariants should be also stable under atomic vibrations, while discrete invariants such as symmetry groups discontinuously change under perturbations. We define the infinite sequence of Average Minimum Distances (AMDs) that satisfy the above conditions and can distinguish homometric crystals, hence are stronger than powder diffraction patterns. The AMDs can visualize geometric similarities between any periodic patterns or crystals. The classification power of AMDs is validated on the dataset of almost 6K simulated T2 crystals reported in Nature 543 (2017), 657-664.

## 1. Importance of isometry invariants for Crystal Structure Prediction

A periodic crystal consists of periodically repeated unit cells (possibly non-rectangular parallelepipeds) containing a finite motif of atoms or molecules, see Fig. 1 and mathematical details in Definition 1. The decomposition of a crystal into a sum of a unit cell and a motif is highly ambiguous, because one choose infinitely many bases or cells that define with suitable motives an equivalent crystal. Since most crystals are rigid, the most natural equivalence of crystals is a rigid motion, a composition of translations and rotations in $\mathbb{R}^3$. Hence a crystal is an equivalence class of infinitely many decompositions cell+motif modulo all changes of a basis and rigid motions:

$$\text{a periodic crystal} = \left\{ \frac{(\text{a unit cell}) + (\text{a motif of atoms or ions})}{(\text{a change of a basis}) \times (\text{a rigid motion})} \right\}.$$

The Crystal Structure Prediction (CSP) aims to predict a thermodynamically stable arrangement of given atoms or molecules. The CSP was pioneered in 1960s when ball models of atoms were physically shaken in a box until they settle in a stable configuration (Kitaigorodsky, 2012). Nowadays this physical shaking is simulated by supercomputers, which start from millions of almost random arrangements and minimize a complicated energy function. This energy has no simple expression and depends on (theoretically infinite) interactions between atoms within a periodic crystal.

A typical CSP software outputs thousands of approximate local minima of this energy, i.e. simulated crystals whose local perturbations are unlikely to produce more stable arrangements. Simulated crystals are visualized by an *energy landscape* representing each crystal as a dot with two coordinates (density,energy), see Fig. 1. The key CSP challenge is the *embarrassment of over-prediction* when the state-of-the-art optimization outputs too many approximate local minima (Price, 2018). Materials scientists expect only few stable crystals (*metastable polymorphs*) at deep local minima of the energy function on a continuous space of all potential crystals, see Fig. 2.

Any energy landscape should be post-processed to (1) remove numerous near duplicate crystals not to waste more time on predicting other properties by further simulations; and (2) identify really different crystals separated by high energy barriers. If all local minima are shallow, there is no chance to synthesize a stable crystal and one should look for a better chemical composition (Pulido *et al.*, 2017).

Since both problems above remain unresolved, even more supercomputer's time (12 weeks (Pulido *et al.*, 2017) in the case of Fig. 1) is spent on predicting target properties of crystals for applications. Since all simulated crystals have the same chemical compositions, they can be distinguished only by their geometry. Geometrically different crystals often have different properties such as solubility, which is vitally important in the pharmaceutical industry. In 1998 manufacturing the HIV drug ritonavir (branded as Kaletra) accidentally produced a more stable but much less soluble polymorph, which has made the drug useless and put many lives at risk (Morissette *et al.*, 2003).

The energy landscape in the last picture of Fig. 1 was lucky due to downward spikes that hinted at deep local minima. This landscape is only a discrete sample from a continuous space of all crystals with a fixed composition. To parameterize such a space as a geographic map of a new planet, we need to uniquely name each crystal so that different crystals have different names, and similar crystals have similar names.

The above mapping problem is important and hard not only when all crystals in question have the same chemical composition, but also to study pure geometric patterns with zero-sized points representing all atoms. Such geometric patterns or templates can be used for put different atoms at known positions, hence generating new crystals without starting from scratch. Though the main results in the paper are stated in mathematical terms, section 7 has an experimental validation on the T2 database of 5799 crystals that were impossible to reliably compare by past tools.

## 2. The isometry classification problem for periodic crystals

This section first introduces periodic sets that model all crystals and then states the algorithmic problems for their stable-under-noise classification modulo isometries.

In the Euclidean space $\mathbb{R}^n$, any point $p \in \mathbb{R}^n$ can be represented by the vector $\vec{p}$ from the origin of $\mathbb{R}^n$ to the point $p$. The symbol $\vec{p}$ will also denote the class of all equal vectors that have equal coordinates. The *Euclidean* distance between points $p, q \in \mathbb{R}^n$ is denoted by $|pq| = |\vec{p} - \vec{q}|$. For a standard orthonormal basis $\vec{e}_1, \ldots, \vec{e}_n$, the lattice $\mathbb{Z}^n \subset \mathbb{R}^n$ consists of all points with integer coordinates.

Definition 1 below models all atoms in crystals as zero-sized points, which is enough for their isometry classification. To model real atoms, one can labels for elements such as C for carbon, O for oxygen etc. Geometrically, atoms can be modeled as weighted points, i.e. balls of different (usually van der Waals) radii. This paper uses the term *periodic set* in any dimension, while *periodic crystals* refer only to dimension 3.

**Definition 1** (a lattice, a periodic set). A *lattice* $\Lambda$ in $\mathbb{R}^n$ consists of all linear combinations $\sum\limits_{i=1}^{n} \lambda_i \vec{v}_i$ with integer coefficients $\lambda_i \in \mathbb{Z}$. Here the vectors $\vec{v}_1, \ldots, \vec{v}_n$ should form a *basis* so that if $\sum\limits_{i=1}^{n} \lambda_i \vec{v}_i = \vec{0}$ for some real $\lambda_i$, then all $\lambda_i = 0$. A *periodic set* (or a *crystal*) consists of a basis $\vec{v}_1, \ldots, \vec{v}_n$ and a *motif $M$* of finitely many points $p_1, \ldots, p_m$ (representing molecules, atoms or ions) in the *unit cell* $U(\vec{v}_1, \ldots, \vec{v}_n) = \left\{ \sum\limits_{i=1}^{n} \lambda_i \vec{v}_i : \lambda_i \in [0, 1] \right\}$, which is the parallepiped spanned by $\vec{v}_1, \ldots, \vec{v}_n$. ■

The two pictures in the top left of Fig. 3 show two lattices with a square unit cell and a single black point in a motif. Though the lattices look different, they are related by a rotation through $\frac{\pi}{4}$, hence are isometric, see Definition 17. Any periodic set can be considered as the *Minkowski* sum of a lattice and a motif, i.e. $S = \Lambda + M = \{\vec{u} + \vec{v} : u \in \Lambda, v \in M\}$. Any periodic set is a finite union of translates of $\Lambda$.

A lattice $\Lambda$ of a periodic set $S = M + \Lambda \subset \mathbb{R}^n$ is not unique in the sense that

$S$ can be generated by a sublattice of $\Lambda$ and a motif larger than $M$. For example, if $U$ is any unit cell of $\Lambda$, the sublattice $2\Lambda$ has the $2^n$ times larger unit cell $2^nU$ (twice larger along each of $n$ basis vectors of $U$), hence contains $2^n$ times more points than $M$. Such an extended unit cell $2^nU$ is superfluous, because $S$ remains invariant under translations along not only integer linear combinations $\sum\limits_{i=1}^{n} \lambda_i \vec{v}_i$ with $\lambda_i \in \mathbb{Z}$, but also for half-integer coefficients $\lambda_i \in \frac{1}{2}\mathbb{Z}$. The two periodic sets in the bottom left of Fig. 3 look even more different than square lattices above. However, they are also isometric and actually represent the same hexagonal lattice, because every black point has exactly 6 nearest neighbors that form a regular hexagon.

The key obstacle to compare crystals modulo isometries is the enormous ambiguity or non-uniqueness of a crystal representation illustrated in Fig. 3. A standard Crystallographic Information File (CIF) contains parameters of a unit cell spanned by a linear basis in $\mathbb{R}^3$ and fractional coordinates of atoms from a motif in this basis. If we change a basis as in the bottom left of Fig. 3, the same hexagonal lattice will have a new CIF with a different unit cell possibly containing a different number of points with new fractional coordinates. Hence cell-dependent descriptors of a crystal can not be justified for comparing crystals modulo isometries. For example, humans should be not be compared or identified by the average color of their clothes, though such colors are easily accessible in photos. Justified comparisons should use only *invariant* features, e.g. biometric data of a human. Any machine learning algorithm can be confused until this representation problem is properly resolved.

The data representation challenge is stated below as the problem to classify periodic sets modulo isometries (or rigid motions) in $\mathbb{R}^n$, see details in Definition 17.

**Problem 2** (algorithmic classification of periodic sets modulo isometries)**.** Find a function $I$ on periodic sets in $\mathbb{R}^n$ satisfying the following conditions:

(2a) *invariance* : $I$ is preserved by isometries: $I(S) = I(Q)$ for any isometric $S, Q$;

(2b) *completeness* : if the invariants coincide $I(S) = I(Q)$, then $S, Q$ are isometric;

(2c) *continuity* : $I(S)$ continuously changes under perturbations of points in $S$;

(2d) *computability* : $I(S)$ is computable in a polynomial time in the size of a motif. ∎

Though a minimal (by volume) unit cell isn't invariant under a change of basis, the volume is invariant. The second set in Fig. 4 is a slight perturbation of the square lattice, but has a rectangular minimal unit cell, not a square. Hence the volume of a minimal cell is unstable under atomic vibrations. Condition (2c) is needed to continuously quantifying a similarity between crystals. Algorithmic condition (2d) is added to guarantee fast time processing for large crystal datasets.

Condition (2b) means that a complete invariant is sufficient to unambiguously identify a periodic crystal in the same way as a DNA code identifies a human. Many claimed 'fingerprints of materials' distinguish usually about 90% of crystals in certain datasets. The *density* of a crystal defined as the molecular weight (or simply the number of points) within a unit cell divided by the cell volume satisfies conditions (2a,c,d), but there was no complete invariant that provably satisfies completeness (2b).

### 3. Closely related past work on comparisons of point sets and crystals

This section discusses the closest work for finite and periodic sets. The excellent book (Liberti & Lavor, 2017) reviews the wider area of distance geometry. The full distribution of all pairwise Euclidean distances $|ab|$ between points $a, b$ in a finite set $S \subset \mathbb{R}^m$ is a well-known isometry invariant. This invariant is almost complete (Boutin & Kemper, 2004). The last picture of Fig. 4 shows non-isometric sets that are not distinguishable by all pairwise distances, i.e. a 4-point set can not be uniquely reconstructed modulo an isometry of $\mathbb{R}^2$ from the distances $\{\sqrt{2}, \sqrt{2}, 2, \sqrt{10}, \sqrt{10}, 4\}$. This example can be extended to any number of points, see Fig. 6. Our methods are similar

to the work (Lai & Zhao, 2014) for finite point clouds.

For periodic sets such as crystals usually given as a CIF file with a unit cell and a motif, it is inevitable to start from a unit cell. However, if output descriptors still depend on a unit cell (Himanen *et al.*, 2020), they are non-invariants modulo isometries. Though the average color can sometimes distinguish all people in a meeting, non-invariants aren't reliable for identifying humans.

The past approach was to try to find a unique unit cell of a crystal. The best example is Niggli's reduced cell in Hahn *et al.* (1983, section 9.3), so Niggli's reduction should be the first step. In 1980 Niggli's cell was shown to be unstable in the sense that a reduced cell of a perturbed lattice can have a basis that substantially differs from that of a non-perturbed lattice, see (Andrews *et al.*, 1980), (Andrews *et al.*, 2019).

More than 40 years since (Andrews *et al.*, 1980), the difficulty of comparing periodic sets is highlighted at http://roninstitute.org/research-scholars/larry-andrews: "...find a measure of the difference between pairs of lattices. Surprisingly, this is not a mathematical problem with a well-defined solution". Our recent work (Mosca & Kurlin, 2020) has resolved this problem for lattices by introducing two distances that satisfy the metric axioms so that the distance between any isometric lattices is 0.

Though there is still no justified distance that satisfies metric axioms for any periodic crystals, the COMPACK algorithm (Chisholm & Motherwell, 2005) in the Mercury software is widely used for a pairwise comparison of crystals as follows. Within given tolerances ($20°$ for angles and $20\%$ for distances), up to a given number (15 by default) of molecules from two crystals are matched by a rigid motion that minimizes the Root Mean Square deviation of $n$ matched atoms RMS $= \sqrt{\frac{1}{n} \sum_{i=1}^{n} |p_i - q_i|^2}$. Table 1 shows how this RMS depends on the maximum number of attempted molecules to match by a rigid motion. A final number of matched molecules seems rather unpredictable.

The newer COMPSTRU algorithm (Flor *et al.*, 2016) like COMPACK predicts a similarity between a reference crystal $S$ and other available crystals whose unit cell parameters are close to those of $S$. The default thresholds are $5°$ for angles and $0.5\text{Å}$ for distances ($1\text{Å} = 10^{10}\text{m}$). The COMPSTRU comparison is restricted to crystals that have the same space-group type. Crystals are compared by powder diffraction patterns up to a cut off radius (Oliynyk *et al.*, 2016), which introduces an extra parameter without resolving the underlying instability under perturbations.

### 4. A fast algorithm to detect sets with identical diffraction patterns

This section discusses homometric crystals that were hard to distinguish, because they have identical diffraction patterns depending only on the *difference* set below.

**Definition 3** (difference multi-set $\text{Dif}(S)$, distance multi-set $\text{Dist}(S)$)**.** Let $S \subset \mathbb{R}^n$ be a finite or a periodic set. The *difference* multi-set is $\text{Dif}(S) = \{\vec{a} - \vec{b}$ for all points $a, b \in S\}$. The *distance* multi-set is $\text{Dist}(S) = \{|\vec{a} - \vec{b}|$ for all points $a, b \in S\}$. ∎

If a set of points $S \subset \mathbb{R}^n$ is finite, then so is the difference set $\text{Dif}(S)$, hence the vector differences $\vec{a} - \vec{b}$ can be counted with multiplicities. For any periodic set $S$, any vector difference or a distance will be repeated infinitely many times due to periodicity, hence all values in $\text{Dif}(S)$ and $\text{Dist}(S)$ have the same infinite (countable) multiplicity.

**Example 4** (Patterson's homometric 1D periodic sets)**.** Patterson in Patterson (1944, p. 197, Fig. 2) has suggested the 1D periodic sets $S = \{0, 1, 3, 4\} + 8\mathbb{Z}$ and $Q = \{0, 3, 4, 5\} + 8\mathbb{Z}$, see Fig. 5 and 6. Theorem 10 will justify why $S, Q$ are non-isometric.

The vector differences of the 4-point motives of the periodic sets $S, Q$ in Fig. 5 differ:

| $S$ | 0 | 1 | 3 | 4 |
|---|---|---|---|---|
| 0 | 0 | −1 | −3 | −4 |
| 1 | 1 | 0 | −2 | −3 |
| 3 | 3 | 2 | 0 | −1 |
| 4 | 4 | 1 | 3 | 0 |

and

| $Q$ | 0 | 3 | 4 | 5 |
|---|---|---|---|---|
| 0 | 0 | −3 | −4 | −5 |
| 3 | 3 | 0 | −1 | −2 |
| 4 | 4 | 1 | 0 | −1 |
| 5 | 5 | 2 | 1 | 0 |

, but they coincide modulo 8 (with infinite multiplicities): $\text{Dif}(S) \equiv \{0, 1, 2, 3, 4, 5, 6, 7\} \equiv \text{Dif}(Q) \bmod 8$.

The equivalence modulo 8 gives rise to a bijection between all 16 elements of the distance matrices above, hence to a bijection between the differences multi-sets $D(S) \to D(Q)$, e.g. the difference $(8i + 1) - (8j + 4) = 8(i - j) - 3 \equiv 5 \pmod 8$ in $S$ can be bijectively mapped to $(8i + 5) - 8j = 8(i - j) + 5$ in $Q$. Fig. 6 shows a generic pair from the family of homometric sets $S(r), Q(r)$, where $r = 1$ is for the sets $S, Q$ on the left. The mirror image of $S(r) = \{0, r, r + 2, 4\} + 8\mathbb{Z}$ under $t \mapsto 4 - t$ coincides with $S(2 - r) = \{0, 2 - r, 4 - r, 4\} + 8\mathbb{Z}$, so they are equivalent modulo all isometries including reflections. Similarly, $Q(r)$ and $Q(2 - r)$ are isometric by the reflection $t \mapsto -t$. To distinguish all these sets modulo an isometry in section 5, we can assume that $0 < r \leq 1$.

**Definition 5** (homometric sets)**.** Finite or periodic sets $S, Q \subset \mathbb{R}^n$ are called *homometric* if there is a bijection between their multi-sets $\mathrm{Dif}(S) \to \mathrm{Dif}(Q)$ from Definition 3. So if $S, Q \subset \mathbb{R}^3$ are crystals, they have identical diffraction patterns. ∎

The following result makes the experimental concept of a homometric crystal verifiable in an algorithmic way. Theorem 6 and all others are proved in Appendix B.

**Theorem 6** (a fast criterion of homometric sets)**.** Any periodic sets $S, Q \subset \mathbb{R}^n$ are homometric in the sense of Definition 5 if and only if $S, Q$ have a common lattice $\Lambda$ such that their sets of vector differences are equal modulo this lattice: $\mathrm{Dif}(S) \equiv \mathrm{Dif}(Q)$ $(\mathrm{mod}\ \Lambda)$. Given a common unit cell containing $m$ points of sets $S, Q \subset \mathbb{R}^n$, there is an algorithm of complexity $O(m^2)$ to determine whether $S, Q$ are homometric. ∎

## 5. Point-wise distributions of distances in finite and periodic sets

This section introduces new invariants of periodic crystals: point-wise distributions of distances in Definition 7 and their simplified averages (AMDs) in Definition 9.

**Definition 7** (point-wise distribution of distances PDD)**.** Let a periodic set $C = M + \Lambda$ have a motif $M$ of $m$ points $p_1, \ldots, p_m$. For a fixed integer $k \geq 1$, the *point-wise distribution of distances* is the $m \times k$ matrix $\mathrm{PDD}(C; k)$, whose $i$-th row corresponds to the point $p_i$, $i = 1, \ldots, m$. The $i$-th row consists of the ordered distances $d_{i1} \leq \cdots \leq d_{ik}$ measured from $p_i$ to its first $k$ nearest neighbors within $C$. ∎

The sets $S, Q$ in Fig. 5 have these point-wise distribution of distances for $k = 3$.

| $S$ | 1st distance | 2nd distance | 3rd | $Q$ | 1st distance | 2nd distance | 3rd |
|---|---|---|---|---|---|---|---|
| $p_1 = 0$ | $\|0-1\| = 1$ | $\|0-3\| = 3$ | 4 | $p_1 = 0$ | $\|0-3\| = 3$ | $\|0-(-3)\| = 3$ | 4 |
| $p_2 = 1$ | $\|1-0\| = 1$ | $\|1-3\| = 2$ | 3 | $p_2 = 3$ | $\|3-4\| = 1$ | $\|3-5\| = 2$ | 3 |
| $p_3 = 3$ | $\|3-4\| = 1$ | $\|3-1\| = 2$ | 3 | $p_3 = 4$ | $\|4-3\| = 1$ | $\|4-5\| = 1$ | 4 |
| $p_4 = 4$ | $\|4-3\| = 1$ | $\|4-1\| = 3$ | 4 | $p_4 = 5$ | $\|5-4\| = 1$ | $\|5-3\| = 2$ | 3 |

The rows of $\mathrm{PDD}(C; k)$ correspond to an arbitrary order of given points $p_1, \ldots, p_m \in M$. There is a suitable convention to order rows by using columns. The columns of $\mathrm{PDD}(C)$ are naturally ordered by increasing distances to neighbors. Then the rows (hence, the points $p_1, \ldots, p_m$) can *lexicographically* ordered as follows. A row $(d_{i1}, \ldots, d_{ik})$ is smaller than $(d_{j1}, \ldots, d_{jk})$ if the first (possibly none) distances coincide: $d_{i1} = d_{j1}, \ldots, d_{il} = d_{jl}$ for some $l \in \{1, \ldots, k-1\}$ and the next distances satisfy $d_{i,l+1} < d_{j,l+1}$. In this lexicographic order, the periodic sets $S, Q$ from Fig. 5 have

$$\mathrm{PDD}(S; 3) = \begin{pmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \\ 1 & 3 & 4 \\ 1 & 3 & 4 \end{pmatrix} \text{ and } \mathrm{PDD}(Q; 3) = \begin{pmatrix} 1 & 1 & 4 \\ 1 & 2 & 3 \\ 1 & 2 & 3 \\ 3 & 3 & 4 \end{pmatrix}.$$

Notice that $\mathrm{PDD}(S; 3)$ contains two pairs of identical rows, because $S = \{0, 1, 3, 4\} + 8\mathbb{Z}$ is symmetric with respect to the reflection $t \mapsto 4 - t \pmod 8$. The similar reflection $t \mapsto -t \pmod 8$ explains two identical rows in $\mathrm{PDD}(Q; 3)$ of $Q = \{0, 3, 4, 5\} + 8\mathbb{Z}$.

The point-wise distribution of distances (PDD) in Definition 7 depends on a number $m$ of points in a given cell $U$. If we make one of its edges twice longer, the resulting non-primitive unit cell contains $2m$ points and PDD will be twice larger. However, a translated copy of any point $p_i \in U$ will have exactly the same ordered distances to

its neighbors as $p_i$ due to periodicity. After doubling $U$ as above, every row appears twice in PDD. These repetitions can be reduced by using the weights of rows below.

**Definition 8** (weighted point-wise distribution WPD). For a periodic set $S$ with $m$ points in a motif, the *weight* of a row in $\mathrm{PDD}(S; k)$ is the number of times the row appears in $\mathrm{PDD}(S; k)$ divided by $m$ so that all weights sum up to 1. The *weighted point-wise distribution* $\mathrm{WPD}(S; k)$ is obtained from $\mathrm{PDD}(S; k)$ by keeping only one of identical rows and putting the weight of this row into the extra $(k + 1)$-st column. ∎

The rows of $\mathrm{WPD}(C; k)$ are lexicographically ordered as in $\mathrm{PDD}(C; k)$. Then $S, Q$ in Fig. 5 have $\mathrm{WPD}(S; 3) = \left( \begin{array}{ccc|c} 1 & 2 & 3 & 1/2 \\ 1 & 3 & 4 & 1/2 \end{array} \right)$ and $\mathrm{WPD}(Q; 3) = \left( \begin{array}{ccc|c} 1 & 1 & 4 & 1/4 \\ 1 & 2 & 3 & 1/2 \\ 3 & 3 & 4 & 1/4 \end{array} \right)$.

Any isometric crystals have the same lattice, the same Niggli's reduced cell, the same number $m$ of points in a motif, hence the same number of rows in PDD and WPD. However, this isometry invariant (the number of points within a primitive cell) is unstable under perturbations by the following reasons. If we perturb one point within a motif, hence all its periodic copies of this point in a crystal, the perturbed crystal can have a different number of rows in WPD. The weights of rows will help continuously quantify perturbations by using a distance between distributions.

**Definition 9** (average minimum distance $\mathrm{AMD}_k$). For $k \geq 1$, the *average minimum distance* of a periodic set $S$ with $m$ points $p_1, \ldots, p_m$ in a motif is $\mathrm{AMD}_k(S) = \frac{1}{m} \sum_{i=1}^{m} \mathrm{PDD}_{ik}(S; k)$, the average of the last $k$-th column in $\mathrm{PDD}(S; k)$. Alternatively, if $\mathrm{WPD}(S; k)$ has $l$ rows with weights $w_1, \ldots, w_l$ such that $\sum_{i=1}^{l} w_i = 1$, then $\mathrm{AMD}_k(S) = \sum_{i=1}^{l} w_i \mathrm{WPD}_{ik}$ is the weighted average of the $k$-th column in $\mathrm{WPD}(S; k)$. ∎

Any lattice $L \subset \mathbb{R}^n$ has a unit cell with only one point in a motif. After a suitable translation, this point can be assumed to be at the origin $0 \in \mathbb{R}^n$. Then $\mathrm{AMD}_k(L)$ is the $k$-th minimum distance from 0 to another point of $L$, see Fig. 11 in appendix A.

## 6. Invariance of distance distributions and their stability under noise

This section proves that the WPDs and AMDs from section 5 are stable isometry invariants (Theorems 10, 12, 14), complete for generic sets (Theorem 15).

Since the periodic sets $S(r) = \{0, r, r + 2, 4\} + 8\mathbb{Z}$ and $S(2 - r)$ are isometric by the reflection $t \mapsto 4 - t$, their WPDs are identical. The similar conclusion holds for $Q(r) = \{0, r + 2, 4, r + 4\} + 8\mathbb{Z}$ and $Q(2 - r)$ isometric by the reflection $t \mapsto -t$. Theorem 10 will justify that all $S(r), Q(r)$ are not isometric to each other for $0 < r \leq 1$.

**Theorem 10** (isometry invariance of WPDs and AMDs). For any finite or periodic set $S \subset \mathbb{R}^n$, the weighted point-wise distribution $\mathrm{WPD}(S; k)$ and average minimum distance $\mathrm{AMD}_k(S)$ are isometry invariants for any number $k \geq 1$ of neighbors. ∎

The power of WPDs is illustrated by classifying the homometric sets that are impossible to distinguish by diffraction patterns. Table 2 in appendix A shows the detailed computations of $\mathrm{PDD}(S(r); 5)$ and $\mathrm{PDD}(Q(r); 5)$ for the homometric sets $S(r), Q(r)$ in Fig. 6, where the rows are ordered by the given points. The lexicographic re-ordering of the rows of PDDs gives the weighted point-wise distributions:

$$\mathrm{WPD}(S(r); 5) = \left( \begin{array}{ccccc|c} r & 2 & 4 - r & 4 + r & 6 & 1/4 \\ r & 2 + r & 4 & 4 & 6 - r & 1/4 \\ 2 - r & 2 & 2 + r & 6 - r & 6 & 1/4 \\ 2 - r & 4 - r & 4 & 4 & 4 + r & 1/4 \end{array} \right),$$

$$\mathrm{WPD}(Q(r); 5) = \left( \begin{array}{ccccc|c} r & 2 - r & 4 & 4 & 6 - r & 1/4 \\ r & 2 & 4 - r & 4 + r & 6 & 1/4 \\ 2 - r & 2 & 2 + r & 6 - r & 6 & 1/4 \\ 2 + r & 4 - r & 4 & 4 & 4 + r & 1/4 \end{array} \right).$$

The first columns of WPDs for $k = 1$ distinguish all $S(r), Q(r)$ for any $0 < r \leq 1$. $\mathrm{AMD}_k(S(r))$ are independent of $r$, hence don't distinguish $S(r)$ modulo isometries. However, for $0 < r \leq 1$, the first minimum distance between any points in $S(r)$ equals $r$ and implies that $S(r)$ are not isometric to each other for different $r$.

Since all atoms vibrate above the absolute zero temperature, the bottleneck distance between crystals in Definition 11 naturally quantifies crystal similarities.

**Definition 11** (bottleneck distance BND). For a fixed bijection $g : S \to Q$ between finite or periodic sets $S, Q \subset \mathbb{R}^n$, the *maximum deviation* is the supremum $\sup\limits_{a \in S} |a - g(a)|$ of Euclidean distances. The *bottleneck distance* is $\text{BND}(S, Q) = \inf\limits_{g:S \to Q} \sup\limits_{a \in S} |a - g(a)|$ is the infimum of maximum deviations over all bijections $g : S \to Q$. ∎

The key obstacle in applying the bottleneck distance to real crystals is the minimization over bijections between infinite sets. Any reduction to a finite set is hard to justify because of the instability of a unit cell under perturbations (Andrews *et al.*, 1980).

**Theorem 12** (stability of AMD under noise). For any number $k \geq 1$ of neighbors, all finite or periodic sets $S, Q$ satisfy $|\text{AMD}_k(S) - \text{AMD}_k(Q)| \leq 2\text{BND}(S, Q)$. ∎

Weighted point-wide distributions are matrices that can have different sizes, hence are harder to compare than AMD vectors of a fixed length $k$. Definition 13 (Rubner *et al.*, 2000) introduces a suitable distance between distributions of different sizes.

**Definition 13** (earth mover's distance EMD). Fix two finite or periodic sets $S, Q \subset \mathbb{R}^n$ and a number $k$ of nearest neighbors for each point. Let $\text{WPD}(S; k)$ consist of $m_S$ rows $R_i(S) \in \mathbb{R}^k$. Each row $R_i(S)$ has a weight $w_i(S)$, $i = 1, \ldots, m(S)$ so that $\sum\limits_{i=1}^{m(S)} w_i(S) = 1$. Using the similar notations for the set $Q$, we quantify by a parameter $0 \leq f_{ij} \leq 1$ a move from each row $R_i(S) \in \mathbb{R}^k$ to another row $R_j(Q) \in \mathbb{R}^k$ of a weight $w_j(Q)$, where $j = 1, \ldots, m(Q)$. The distance between such rows (vectors in $\mathbb{R}^k$) is Euclidean. The *earth mover's distance* is defined as the minimum value of the cost flow $\text{EMD}(S, Q) = \sum\limits_{i=1}^{m(S)} \sum\limits_{j=1}^{m(Q)} f_{ij} |R_i(S) - R_j(Q)|$ over all $0 \leq f_{ij} \leq 1$ subject to $\sum\limits_{j=1}^{m(Q)} f_{ij} = w_i(S)$ for $i = 1, \ldots, m(S)$, and $\sum\limits_{i=1}^{m(S)} f_{ij} = w_j(Q)$ for $j = 1, \ldots, m(Q)$. ∎

The first condition $\sum\limits_{j=1}^{m(Q)} f_{ij} = w_i(S)$ means that the full weight $w_i(S)$ of the row $R_i(S)$ 'flows' into the rows $R_j(Q)$, each via a 'flow' $f_{ij}$, $j = 1, \ldots, m(Q)$. Similarly,

the second condition $\sum_{i=1}^{m(S)} f_{ij} = w_j(Q)$ means that all 'flows' $f_{ij}$ from rows $R_i(S)$ for $i = 1, \ldots, m(S)$ 'flow' into the row $R_j(Q)$ and sum up to the full weight $w_j(Q)$.

The earth mover's distance (EMD) has more than one advantage over the bottleneck distance (BND) for periodic sets. First, the EMD uses the isometry invariant WPD, whose stability in the EMD is proved in Theorem 14. The BND between infinite sets can be computed only on finite subsets, e.g. on points in an extended cell, which is unstable (Andrews *et al.*, 1980). Second, even for finite subsets, the fastest algorithm in Kerber *et al.* (2017, Theorem 3.1) computes the BND (for 2D set of $m$ points) in time $O(m^{1.5} \log m)$. The EMD can be approximated (Pele & Werman, 2008) in a time linear in the size of any k-dimensional distributions.

**Theorem 14** (stability of weighted point-wise distribution WPD)**.** For any number $k \geq 1$ of neighbors, any finite or periodic sets $S, Q$ satisfy $\mathrm{EMD}(S, Q) \leq 2\sqrt{k}\mathrm{BND}(S, Q)$. So any small perturbation of positions in the bottleneck distance (BND) yields a small change of the weighted point-wise distribution in the earth mover's distance. ∎

After satisfying the invariance and stability conditions in Problem 2, Theorem 15 proves a generic completeness of the weighted point-wise distributions (WPDs).

**Theorem 15** (unique reconstruction of a finite set from WPD)**.** Let a finite set $S \subset \mathbb{R}^n$ have $m$ points such that all pairwise distances between points of $S$ are distinct. Then $S$ can be uniquely reconstructed modulo an isometry of $\mathbb{R}^n$ from $\mathrm{WPD}(S; m-1)$. ∎

We conjecture that $\mathrm{WPD}(S; k)$ are complete isometry invariants for sufficiently large $k$ depending on a complexity of $S \subset \mathbb{R}^n$. If $S$ is a finite set of $m$ points, then $k = m-1$ should be enough. The sets $A, B \subset \mathbb{R}^2$ in the last picture of Fig. 4 have $\mathrm{WPD}(A; 3) = \left( \begin{array}{ccc|c} \sqrt{2} & 2 & \sqrt{10} & 1/2 \\ \sqrt{2} & \sqrt{10} & 4 & 1/2 \end{array} \right)$ and $\mathrm{WPD}(B; 3) = \left( \begin{array}{ccc|c} \sqrt{2} & \sqrt{2} & 4 & 1/4 \\ \sqrt{2} & 2 & \sqrt{10} & 1/2 \\ \sqrt{10} & \sqrt{10} & 4 & 1/4 \end{array} \right)$, which distinguish $A, B \subset \mathbb{R}^2$ modulo isometries. Actually, $k = 1$ is enough in this case.

If $m = 1$, any set $S \subset \mathbb{R}^n$ is a lattice and $k$ needs to be at least $n(n+1)$. For example, any lattice in $\mathbb{R}^2$ can be reconstructed from the distribution of 6 minimum distances from the origin 0 to 3 pairs of 6 neighbors symmetric with respect to 0.

## 7. Computations, applications to crystal comparisons and a discussion

Theorem 16 covers final computability condition in Isometry Classification Problem 2.

**Theorem 16** (algorithm for computing new invariants WPDs and AMDs). Let a periodic crystal $S \subset \mathbb{R}^n$ have $m$ points in a unit cell whose extension by a factor $\mu$ covers all $k$ of neighbors of the given points. Then the matrix $\mathrm{WPD}(S; k)$ and all $\mathrm{AMD}_i(S)$ for $i = 1, \ldots, k$ can be computed in time $O(m(n\mu^n + k)\log(\mu^n m))$. ∎

Though we have no exact value of the factor $\mu$, our experiments show that $\mu = O(n)$. So in the practical case of $m = 3$ the time is near linear in the number $m$ of points.

The only input for computing the new invariants is a crystal itself (a unit cell with a finite set of points) without parameters. The number of neighbors $k$ is independent of a crystal and reflects our desire to extra more distance information. For example, vectors of 1000 AMDs will better differentiate crystals than vectors of 100 AMDs.

The Nature paper (Pulido *et al.*, 2017) has reported 4 experimental crystals T2-$\alpha$, T2-$\beta$, T2-$\gamma$, T2-$\delta$ (one more T2-$\varepsilon$ was synthesized after the publication), see Fig. 10 in Appendix A. The synthesis in a lab started only after an energy landscape in Fig. 1 of 5679 simulated crystals produced by 12-week simulations on a supercomputer hinted at potential stable crystals in downward spikes (imaginable deep minima). To validate this approach, the synthesized crystals should be matched with closest crystals from the simulated dataset of 5679. If there was no close match, the expensive simulations missed a real crystal, which is always possible, because the continuous space of all potential crystals in Fig. 2 is randomly and discretely sampled.

Until now the density was practically used as a stable isometry invariant of crystals. The density in the horizontal axis in Fig. 1 can separate nano-porous organic crystals, while inorganic crystals are much denser and can not be well-separated by densities.

Using the density $\Delta$ of an experimental crystal, chemists look for a corresponding simulated crystal in a vertical strip of the energy landscape in Fig. 1 over a small interval around $\Delta$ to allow for errors. From this strip one takes the crystal with the lowest energy as the best guess, which depends on a strip. A final match is confirmed by the RMS deviation between finite portions, which is also uncertain, see Table 1.

For the experimental crystal T2-$\delta$, the past method above found crystal 14 in the simulated dataset of 5679. However, another crystal 15 has a much closer AMD curve in Fig. 7. Though both crystals have almost identical energy and density in Fig. 1 in Table 3. They are separated by their AMD curves (red: dotted vs dashed) in Fig. 7.

Fig. 8 shows that that most stable 100 simulated crystals split into two clusters, which merge at about 17Å. The 1st large cluster on the left contains two simulated matches of the experimental crystals T2-$\varepsilon$ and T2-$\delta$, though the past match 14 is in another subcluster than the new match 15. The 2nd small cluster has all matches of T2-$\alpha$, T2-$\beta$, T2-$\gamma$. The thresholds between the largest subclusters are about 3Å.

Fig. 7 and 8 show that the new invariants from Definitions 8 and 9 continuously quantify similarities between periodic crystals, which was impossible by past non-invariant descriptors or unstable discrete invariants such as symmetry groups.

We have resolved the following challenges in the geometry for periodic sets.

• The criterion in Theorem 6 detects crystals with identical diffraction patterns.

• Theorems 10, 12, 14 have proved that the weighted point-wise distributions $\mathrm{WPD}(S; k)$ and average minimum distances $\mathrm{AMD}_k(S)$ are stable isometry invariants of a finite or a periodic set $S \subset \mathbb{R}^n$ and are computable fast enough by algorithmic Theorem 16.

• Completeness Theorem 15 proves that any set $S \subset \mathbb{R}^n$ of $m$ points with distinct distances can be uniquely reconstructed from its $\text{WPD}(S; m-1)$ modulo isometries.

# Appendix A
## Background on isometries and isometry invariants

We first remind key facts about isometries in $\mathbb{R}^n$ and then give the proofs of all results from sections 4 and 6. The strongest possible equivalence on rigid materials is defined by isometries (or rigid motions) that preserve interpoint distances.

**Definition 17** (isometries). An *isometry* of $\mathbb{R}^n$ is any map $f : \mathbb{R}^n \to \mathbb{R}^n$ that preserves the Euclidean distance, i.e. $|pq| = |f(p)f(q)|$ for any points $p, q \in \mathbb{R}^n$. If $f$ also preserves the *orientation*, i.e. the matrix whose columns are images under $f$ of the standard basis vectors $\vec{e}_1, \ldots, \vec{e}_n$ has a positive determinant, then $f$ can be called a *rigid motion*, because $f$ is included into a continuous family of isometries $f_\lambda : \mathbb{R}^n \to \mathbb{R}^n$, $\lambda \in [0, 1]$, where $f_1 = f$ and $f_0$ is the identity map $f_0(p) = p$ for any $p \in \mathbb{R}^n$. ∎

Any isometry of $\mathbb{R}^n$ can be decomposed into at most $n + 1$ reflections over hyperspaces, hence is bijective and can be inverted. A composition of isometries is also an isometry, which defines the operation in the group $\text{Iso}(\mathbb{R}^n)$ of all isometries in $\mathbb{R}^n$. Rigid motions are orientation-preserving isometries and form the smaller subgroup $\text{Iso}^+(\mathbb{R}^n) \subset \text{Iso}(\mathbb{R}^n)$. Materials are compared modulo rigid motions or the isometries from $\text{Iso}(\mathbb{R}^n)$, because mirror images of materials can have different properties. Examples in $\mathbb{R}^3$ are translations by vectors and rotations around straight lines.

For any $n \times n$ matrix $A$, recall that $A^T$ denotes the *transpose* matrix with elements $A_{ij}^T = A_{ji}$, $i, j = 1, \ldots, n$. A matrix $A$ is *orthogonal* if the inverse matrix $A^{-1}$ equals the

transpose $A^T$. Orthogonality of a matrix $A$ means that $\vec{v} \mapsto A\vec{v}$ maps any orthonormal basis to another orthonormal basis. All orthogonal matrices $A$ have the determinant $\det A = \pm 1$. If $\det A = 1$, then the map $\vec{v} \mapsto A\vec{v}$ preserves an orientation of $\mathbb{R}^n$.

All orthogonal matrices $A$ with $\det A = 1$ form the special orthogonal group $\mathrm{SO}(\mathbb{R}^n)$, where the operation is the matrix multiplication. The group $\mathrm{SO}(\mathbb{R}^2)$ consists of rotations about the origin in the plane. The group $\mathrm{SO}(\mathbb{R}^3)$ consists of rotations about axes passing through the origin in $\mathbb{R}^3$. In general, $\mathrm{SO}(\mathbb{R}^n)$ consists of all isometries from $\mathrm{Iso}^+(\mathbb{R}^n)$ that preserve the origin. Any objects should be classified by invariants that are independent of a given representation of an object. Many machine learning algorithms struggle when features or descriptors include non-invariants of crystals, e.g. parameters of an ambiguous unit cell or atomic coordinates in an arbitrary basis.

**Definition 18** (isometry invariant). An *isometry class* is a set of all materials that are isometric to each other, i.e. any materials $S, Q$ from the same class are related by an isometry $S \to Q$. An *isometry invariant* is a function $I$ that maps all materials from a certain class, e.g. all periodic crystals, to a simpler set (e.g. numbers, matrices) so that $I(S) = I(Q)$ for any isometric materials $S, Q$. An invariant $I$ is called *complete* if the converse is also true: if $I(S) = I(Q)$, then the materials $S, Q$ are isometric. ■

The original definition (Patterson, 1939) said that homometric crystals should be non-isometric. Definition 5 does not have this restriction and defines an equivalence relation on sets satisfying the three axioms:

*reflexivity*: any set $S$ is equivalent to itself;

*symmetry*: if a set $S$ is equivalent to $Q$, then $Q$ is equivalent to $S$;

*transitivity*: if $S$ is equivalent to $Q$ that is equivalent to $T$, then $S$ is equivalent to $T$.

The three axioms above guarantee that all sets can split (or classified) into disjoint equivalence classes (consisting of all sets equivalent to each other) and a classification

modulo an equivalence relation makes sense. The even better equivalence relation is the isometry combined with the homometry saying that one set $S$ is equivalent to a set $Q$ if $\mathrm{Dif}(f(S)) = \mathrm{Dif}(Q)$ for a suitable isometry $f$.

Fig. 11 shows the AMD graphs for six 2D lattices specified by their basis vectors and also shown in Fig. 11 in the same order. Though the orange and red AMD graphs clearly differ Fig. 11 up to $k = 50$, their asymptotic behaviors are very similar for $k$ close to 1000, also for the black and blue graphs, which is interesting to study further.

<h2 style="text-align:center">Appendix B<br/>Proofs of all theorems from the main paper</h2>

*Proof of Theorem 6.* For the 'only if' part, assume that $\mathrm{Dif}(S) = \mathrm{Dif}(Q)$ as infinite sets. Let $S = M + \Lambda$ be any representation of the crystal $S$ in terms of its arbitrarily unit cell. If the motif $M$ consists of $m$ points $p_1, \ldots, p_m$ within the unit cell with a basis $\vec{v}_1, \ldots, \vec{v}_n$, then

$$\mathrm{Dif}(S) = \{p_i - p_j + \sum_{k=1}^{n} \lambda_k \vec{v}_k \ : \ 1 \le i, j \le m, \ \lambda_1, \ldots, \lambda_n \in \mathbb{Z}\}.$$

Hence $\mathrm{Dif}(S)$ is also a periodic set with the same unit cell. The similar conclusion for $\mathrm{Dif}(Q)$ implies that $\mathrm{Dif}(Q) = \mathrm{Dif}(S)$, hence $S$ and $Q$, have the same lattice $\Lambda$. Then the difference sets should be equal modulo this lattice: $\mathrm{Dif}(S) \equiv \mathrm{Dif}(Q) \pmod{\Lambda}$.

For the 'if' part, we start from a common lattice $\Lambda$ of $S, Q$. Any lattice has a unique Niggli's reduced cell, so we assume that given crystals $S, Q$ has a common unit cell $U$ with a basis $\vec{v}_1, \ldots, \vec{v}_n$. The equality $\mathrm{Dif}(S) \equiv \mathrm{Dif}(Q) \pmod{\Lambda}$ means that for any pair of points $q_i, q_j \in Q$, there is a unique pair $p_i, p_j \in S$ and unique coefficients $\lambda_1, \ldots, \lambda_n \in \mathbb{Z}$ such that $q_i - q_j = p_i - p_j + \sum_{k=1}^{n} \lambda_k \vec{v}_k$ and vice versa. We extend this

1-1 correspondence to the infinite set $\text{Dif}(Q)$. For any $q_i - q_j + \sum_{k=1}^{n} \mu_k \vec{v}_k \in \text{Dif}(Q)$, the corresponding difference in $\text{Dif}(S)$ is $p_i - p_j + \sum_{k=1}^{n} (\lambda_k + \mu_k)\vec{v}_k$, which extends the bijection $\text{Dif}(S) \rightarrow \text{Dif}(Q)$ to full crystals.

To determine if $S, Q$ are homometric, one can start with a common cell $U$ containing $m$ points of $S, Q$, which is needed by the above criterion. First compute all $O(m^2)$ pairwise differences (translated to $U$ if necessary) for both $S, Q$. To check if these vector sets coincide, we could lexicographically order them using coordinates in the basis of the cell $U$. Then a single pass over $O(m^2)$ vector differences is enough to decide if $\text{Dif}(S) \equiv \text{Dif}(Q) \pmod{\Lambda}$. $\square$

*Proof of Theorem 10.* Any isometry $f : S \rightarrow Q$ between sets $S, Q \subset \mathbb{R}^n$ establishes a 1-1 correspondence between points of $S$ and $Q$. If $S, Q$ are periodic, $f$ bijectively maps a unit cell $U$ of $S$ to a unit cell $U(Q)$ of $Q$. Hence $f$ allows us to order points $p_1, \ldots, p_m \in U(S)$ according to the order of their images $f(p_1), \ldots f(p_m) \in U(Q)$.

Since the isometry $f$ preserves distances between points, every $i$-th row of $\text{PDD}(S; k)$, which contains the ordered distances from $p_i$ to its first $k$ nearest neighbors, coincides the $i$-th row of $\text{PDD}(Q; k)$, $i = 1, \ldots, m$. These coincidence of rows gives rise to the equality of the matrices $\text{WPD}(S; k) = \text{WPD}(Q; k)$ of Weighted Point-wise Distributions, which are independent of of point ordering. $\square$

Lemma 19 is needed to prove Stability Theorem 12.

**Lemma 19** (perturbed distances)**.** For some $\varepsilon > 0$, let $g : S \rightarrow Q$ be a bijection between finite or periodic sets such that $|a - g(a)| \leq \varepsilon$ for all $a \in S$. Then, for any $i \geq 1$, let $a_i \in S$ and $b_i \in Q$ be the $i$-nearest neighbors of points $a \in S$ and $b = g(a) \in Q$, respectively. Then the Euclidean distances from the points $a, b$ to their $i$-th neighbors $a_i, b_i$ are $2\varepsilon$-close to each other, i.e. $||a - a_i| - |b - b_i|| \leq 2\varepsilon$.

*Proof.* Shifting the point $g(a)$ back to $a$, assume that $a = g(a)$ is fixed and all other points change their positions by at most $2\varepsilon$. Assume by contradiction that the distance from $a$ to its new $i$-th neighbor $b_i$ is less than $|a - a_i| - 2\varepsilon$. Then all first new $i$ neighbors $b_1, \ldots, b_i$ of $a$ within $Q$ belong to the open ball with the center $a$ and the radius $|a - a_i| - 2\varepsilon$. Since the bijection $g$ shifted every point $b_1, \ldots, b_i$ by at most $2\varepsilon$, their preimages $g^{-1}(b_1), \ldots, g^{-1}(b_i)$ belong to the open ball with the center $a$ and the radius $|a - a_i|$. Then the $i$-th neighbor of $a$ within $S$ should be among these $i$ preimages, i.e. the distance from $a$ to its $i$-th nearest neighbor should be strictly less than the assumed value $|a - a_i|$. A similar contradiction is obtained from the assumption that the distance from $a$ to its new $i$-th neighbor $b_i$ is more than $|a - a_i| + 2\varepsilon$. $\square$

*Proof of Theorem 12.* By Lemma 19 each element of $\mathrm{PDD}(S; k)$ changes by at most $2\varepsilon$. Then the average of the $k$-th column changes by at most $2\varepsilon$ as required. $\square$

Lemma 20 is needed to prove Stability Theorem 14.

**Lemma 20** (perturbed rows). For some $\varepsilon > 0$, let $g : S \to Q$ be a bijection between finite or periodic sets such that $|a - g(a)| \leq \varepsilon$ for all $a \in S$. Then, for any $k \geq 1$, the bijection $g$ changes the vector $\vec{R}_a(S) = (|a - a_1|, \ldots, |a - a_k|) \in \mathbb{R}^k$ of the first $k$ minimum distances from any point $a \in S$ to its $k$ nearest neighbors $a_1, \ldots, a_k \in S$ by a Euclidean distance at most $2\varepsilon\sqrt{k}$. So if $b_1, \ldots, b_k \in Q$ are the $k$ nearest neighbors of $b = g(a)$ within $Q$ and $\vec{R}_b(S) = (|b - b_1|, \ldots, |b - b_k|) \in \mathbb{R}^k$ is the vector of the first $k$ minimum distances from the point $b = g(a)$, then $|\vec{R}_a(S) - \vec{R}_b(Q)| \leq 2\varepsilon\sqrt{k}$.

*Proof.* By Lemma 19 every coordinate of $\vec{R}_a(S) \in \mathbb{R}^k$ changes by at most $2\varepsilon$. Hence the Euclidean distance from $\vec{R}_a(S)$ to the perturbed $\vec{R}_b(Q)$ is at most $2\varepsilon\sqrt{k}$. $\square$

*Proof of Theorem 14.* $\mathrm{BND}(S,Q) = \inf\limits_{g:S\to Q} \sup\limits_{a\in S} |a - g(a)|$ by Definition 11 means, for

any $\delta > 0$, there is a bijection $g : S \to Q$ such that $\sup\limits_{a\in S} |a - g(a)| \leq \mathrm{BND}(S,Q) + \delta$. If

the sets $S, Q$ are finite, one can set $\delta = 0$. Indeed, there are finitely many bijections

$S \to Q$, hence the infimum in Definition 11 will be achieved for one of them.

If $S, Q$ are periodic, the chosen bijection $g$ restricts to a bijection between all points

in corresponding unit cells of $S, Q$, so set $m = m(S) = m(Q)$. For any fixed $k \geq 1$,

we will design a flow from the rows of $\mathrm{WPD}(S;k)$ to the rows of $\mathrm{WPD}(Q;k)$ with $f_{ij}$

satisfying Definition 13. We start from a 1-1 flow with $f_{ij} = 0$ for $i \neq j$.

If not all rows $R_i(S)$ in $\mathrm{PDD}(S;k)$ are distinct, we make them symbolically distinct

so that $\mathrm{WPD}(S;k)$ is obtained from $\mathrm{PDD}(S;k)$ by adding the column of equal weights

$\frac{1}{m}$, similarly for $\mathrm{WPD}(Q;k)$. Identifying equal rows later will mean that flows to (or

from) equal (symbolically different) rows are combined into a many-to-one (or one-to-

many, respectively) flow. Since we have the same number $m$ of rows in both matrices

$\mathrm{WPD}(S;k)$ and $\mathrm{WPD}(Q;k)$, we set $f_{ij} = 0$ for $i \neq j$ and $f_{ii} = \frac{1}{m}$, $i = 1, \ldots, m$.

Then $\mathrm{EMD}(S,Q) \leq \frac{1}{m} \sum\limits_{i=1}^{m} |\vec{R}_i(S) - \vec{R}_i(Q)|$, because EMD minimizes the cost over all

flows in Definition 13. Since each $|\vec{R}_i(S) - \vec{R}_i(Q)| \leq 2\sqrt{k}(\mathrm{BND}(S,Q)+\delta)$ by Lemma 20,

we conclude that $\mathrm{EMD}(S,Q) \leq \frac{1}{m} \sum\limits_{i=1}^{m} 2\sqrt{k}(\mathrm{BND}(S,Q) + \delta) = 2\sqrt{k}(\mathrm{BND}(S,Q) + \delta)$.

Since the last inequality holds for any small $\delta > 0$, we get the Lipschitz continuity

$\mathrm{EMD}(S,Q) \leq 2\sqrt{k}\mathrm{BND}(S,Q)$. $\qquad\square$

The following proof will convert a weighted point-wise distribution (in a generic

case when all distances are distinct) into a distance matrix on ordered points.

*Proof of Theorem 15.* Since all pairwise distances between $m$ points of $S$ are distinct,

every distance appears in the matrix $\mathrm{WPD}(S;m-1)$ exactly twice, once as the distance

from a point $p_i$ to its neighbor $p_j$, and once more as the distance from $p_j$ to $p_i$, though

these equal entries are not symmetric. We will convert $\mathrm{WPD}(S;m-1)$ into the distance

matrix $D(S)$ as follows. Let $d_1 < d_2 < \cdots < d_{m-1}$ be all strictly increasing distances from a (say) first point $p_1$ of $S$ to the $m-1$ others.

Each distance $d_i$ from the first row appears exactly once more in another (say, $i'$-th) row of WPD$(S; m-1)$. Then $d_i$ is the distance between the points $p_1$ and $p_{i'}$ numbered as the $i'$-th row. The map of indices $i \mapsto i'$ is a permutation of $\{2, \ldots, m\}$. We set $D_{11} = 0$ and $D_{1,i'} = d_i$ for each $i = 2, \ldots, m$. Then we similarly permute indices in the 2nd row of WPD$(S; m-1)$, starting from the 3rd index due to the symmetry of $D(S)$, and so on. The full distance matrix $D(S)$ uniquely determines a set with ordered points $S \subset \mathbb{R}^n$ modulo isometries by the classical multi-dimensional scaling in Liberti & Lavor (2017, Section 8.5.1). $\qquad\square$

*Proof of Theorem 16.* We build a k-d tree (Brown, 2015) on $\mu^n m$ points in the extended unit cell in time $O(n(\mu^n m) \log(\mu^n m))$. Then all $k$ neighbors of $m$ initial points can be found in time for each $O(km \log(\mu^n m))$. After completing the $m \times k$ matrix PDD, we lexicographically sort its $m$ rows. Each row comparison needs $O(k)$ time. The matrix WPD$(S; k)$ can be obtained in time $O(km \log m)$. All $\text{AMD}_i(S)$ are computed as column averages in time $O(km)$. The total time is $O(m(n\mu^n + k) \log(\mu^n m))$. $\qquad\square$

To guarantee that all $k$ neighbors are correctly found, we incrementally increase $\mu$ and check if a new layer of cells leads to any updates in the current $m \times k$ matrix containing distances from $m$ initial points to their $k$ nearest neighbors. If no updates happen for a point $p$, all $k$ minimum distances from $p$ are correctly found.

Table 1. *The Root Mean Square (RMS) deviation between the experimental T2-δ crystal and its closest simulated version with ID 14 from the T2 dataset in (Pulido* et al.*, 2017). The irregular dependence of RMS on a number of matched molecules makes this comparison unreliable. The computation over 35 molecules was about 1000 times longer than for 15.*

| matched molecules | 5 of 5 | 8 of 10 | 10 of 15 | 11 of 20 | 16 of 25 | 18 of 30 | 21 of 35 |
|---|---|---|---|---|---|---|---|
| RMS in Angstroms | 0.603 | 0.681 | 0.812 | 0.825 | 0.99 | 1.027 | 1.079 |

Table 2. *The point-wise distributions of distances (PDDs) and average minimum distances (AMDs) from Definitions 7 and 9 for* $S(r) = \{0, r, r+2, 4\} + 8\mathbb{Z}$,
$Q(r) = \{0, r+2, 4, r+4\} + 8\mathbb{Z}$, $0 < r \le 1$.

| PDD($S(r)$) | 1st distance | 2nd distance | 3rd distance |
|---|---|---|---|
| $p_1 = 0$ | $|0 - r| = r$ | $|0 - (2+r)| = 2+r$ | $|0 - 4| = 4$ |
| $p_2 = r$ | $|r - 0| = r$ | $|r - (2+r)| = 2$ | $|r - 4| = 4-r$ |
| $p_3 = 2+r$ | $|(2+r) - 4| = 2-r$ | $|(2+r) - r| = 2$ | $|(2+r) - 0| = 2+r$ |
| $p_4 = 4$ | $|4 - (2+r)| = 2-r$ | $|4 - r| = 4-r$ | $|4 - 0| = 4$ |
| AMD$_k(S(r))$ | AMD$_1 = 1$ | AMD$_2 = 2.5$ | AMD$_3 = 3.5$ |

| PDD($Q(r)$) | 1st distance | 2nd distance | 3rd distance |
|---|---|---|---|
| $p_1 = 0$ | $|0 - (2+r)| = 2+r$ | $|0 - (r+4-8)| = 4-r$ | $|0 - 4| = 4$ |
| $p_2 = 2+r$ | $|(2+r) - 4| = 2-r$ | $|(2+r) - (4+r)| = 2$ | $|(2+r) - 0| = 2+r$ |
| $p_3 = 4$ | $|4 - (4+r)| = r$ | $|4 - (2+r)| = 2-r$ | $|4 - 0| = 4$ |
| $p_4 = 4+r$ | $|(4+r) - 4| = r$ | $|(4+r) - (2+r)| = 2$ | $|(4+r) - 8| = 4-r$ |
| AMD$_k(Q(r))$ | AMD$_1 = 1 + 0.5r$ | AMD$_2 = 2.5 - 0.5r$ | AMD$_3 = 3.5$ |

Table 3. *The AMD distance below is the Euclidean distance between the vectors* $(\text{AMD}_1, \ldots, \text{AMD}_{100})$. *Crystal 15 from the dataset of 5679 T2 crystals (Pulido* et al.*, 2017) has a smaller AMD distance to the experimental T2-δ in Fig. 7 than the crystal 14, which was manually found by using Mercury calculating the RMS deviations of up to 15 molecules.*

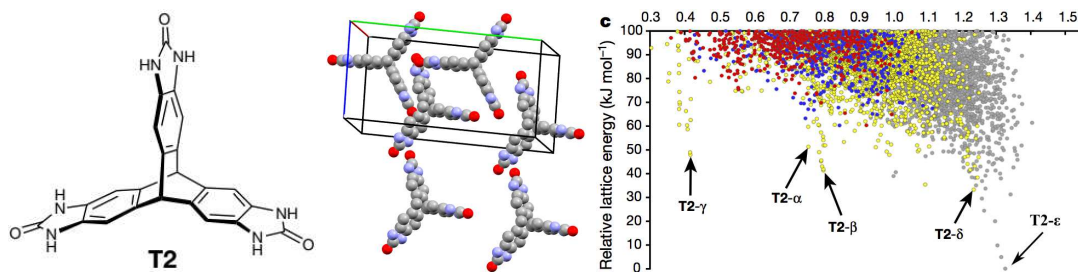| properties of crystals | experimental T2-δ | past matched crystal 14 | new matched crystal 15 |
|---|---|---|---|
| energy (kJ/mol) | unknown | -190.534 | -190.509 |
| density (g/cm$^3$) | 1.215 | 1.23606 | 1.23602 |
| RMS distance to T2-δ | 0 | 0.806 | 0.81 |
| AMD distance to T2-δ | 0 | 2.358 | 0.347 |

Fig. 1. **1st**: T2 molecule (triptycenetrisbenzimidazolon). **2nd**: millions of initial almost random arrangements are iteratively perturbed to minimize an energy. **3rd**: a CSP software can output thousands of simulated crystals often visualized as an energy landscape in Pulido *et al.* (2017, Fig. 2d), where every crystal has two coordinates (density,energy). This landscape 'hints' at deep minima in downward spikes. The past manual matching of five synthesized crystals to simulated predictions is now automated by the new invariant-based algorithm in section 7.



Fig. 2. **Left**: a current energy landscape is a list of simulated crystals. **Middle**: isometry invariants will 'join the dots' and sample a crystal space to find energy barriers. **Right**: a 'mapped' energy landscape (the energy function over a space of crystals) with highlighted deep minima (most stable crystals in red), energy barriers (blue) and other approximations to local minima (black).



Fig. 3. Isometry classes of crystals can be distinguished only by isometry invariants.

Fig. 4. **First two**: the square lattice and its perturbation with a rectangular primitive cell, so the volume of a primitive cell is unstable. **Third**: the experimental crystal T2-$\delta$ overlaid with its closest simulated version was deposited in the CSD with id SEMDIA only new invariants have shown big differences with deposited crystals. **Fourth**: these non-isometric sets can not be distinguished by pairwise distances.
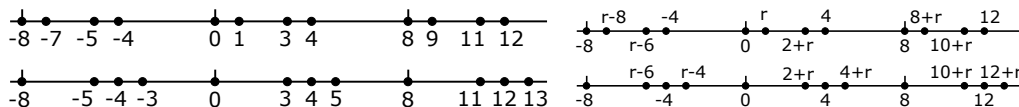


Fig. 5. Non-isometric homometric sets, see Definition 5. **Top**: $S(r) = \{0, r, r + 2, 4\} + 8\mathbb{Z}$. **Bottom**: $Q(r) = \{0, r + 2, 4, r + 4\} + 8\mathbb{Z}$, $0 < r \leq 1$ is a parameter. The simpler versions on the left correspond to $r = 1$. The circular versions are in Fig. 6.



Fig. 6. **First two**: circular versions of the homometric sets $S, Q$ in Fig. 5. Each circle splits into 8 equal arcs. The distances between points (shown outside the disk) are arc lengths (shown inside the disk). **Last two**: homometric sets $S(r) = \{0, r, r + 2, 4\} + 8\mathbb{Z}$, $Q(r) = \{0, r + 2, 4, r + 4\} + 8\mathbb{Z}$, $0 < r < 2$. The distances between points (shown outside the disk) are arc lengths (inside the disk).
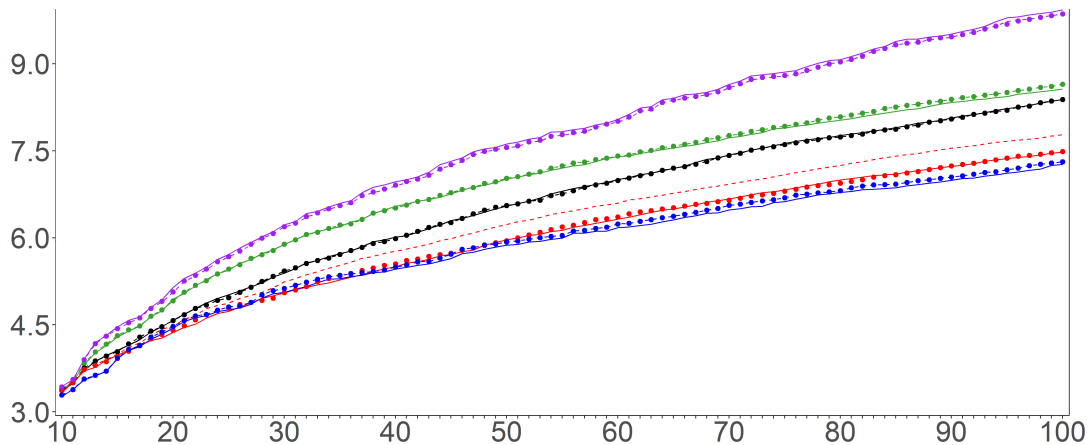
Fig. 7. The average minimum distances $\mathrm{AMD}_k$ in Angstroms with $k = 10, \ldots, 100$. Solid curves are for experimental crystals. Dashed curves are for the past simulated matches reported in (Pulido *et al.*, 2017). Dotted curves are for new matches found by smallest Euclidean distances between AMD vectors with $k = 100$. From top to bottom: purple T2-$\gamma$, green T2-$\alpha$, black T2-$\beta$, red T2-$\delta$ (new match 15 is much closer by AMDs to the experimental crystal than the old match 14), blue T2-$\varepsilon$.
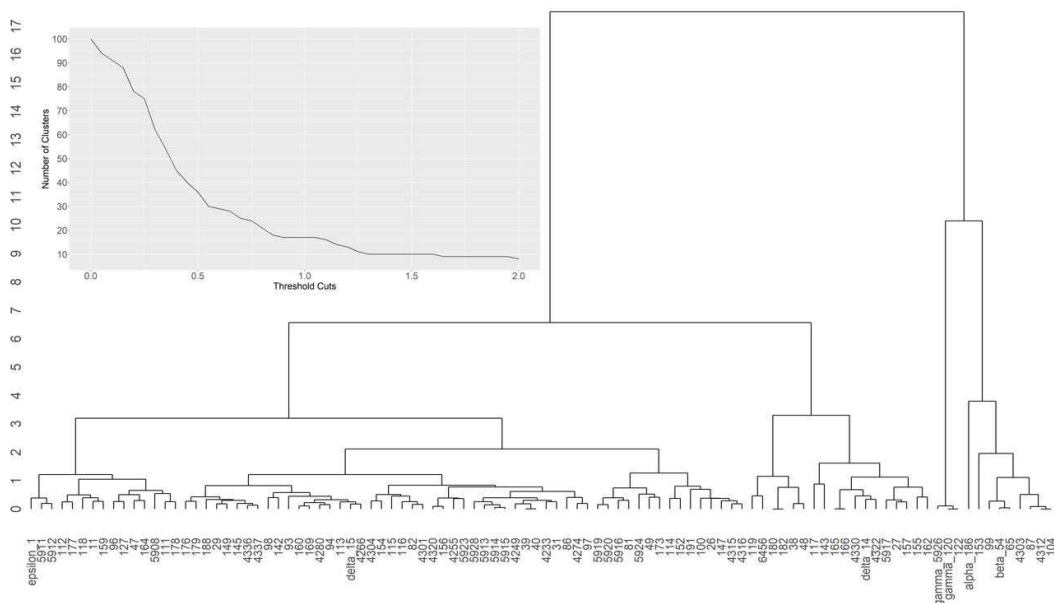


Fig. 8. Complete-linkage clustering by $(\mathrm{AMD}_1, \ldots, \mathrm{AMD}_{100})$ of 100 crystals with lowest energies from the dataset of 5679 simulated crystals (Pulido *et al.*, 2017). The inset image shows how the number of clusters decreases when a threshold grows.

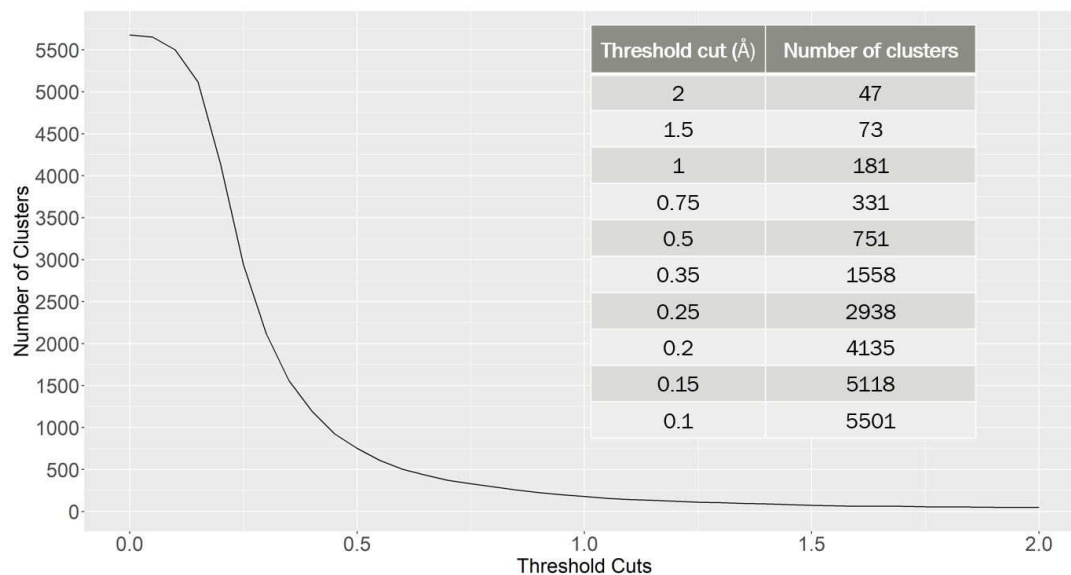| Threshold cut (Å) | Number of clusters |
|---|---|
| 2 | 47 |
| 1.5 | 73 |
| 1 | 181 |
| 0.75 | 331 |
| 0.5 | 751 |
| 0.35 | 1558 |
| 0.25 | 2938 |
| 0.2 | 4135 |
| 0.15 | 5118 |
| 0.1 | 5501 |

Fig. 9. The number of clusters by complete-linkage clustering is continuously decreasing when a threshold grows. Such continuity was impossible by past tools.
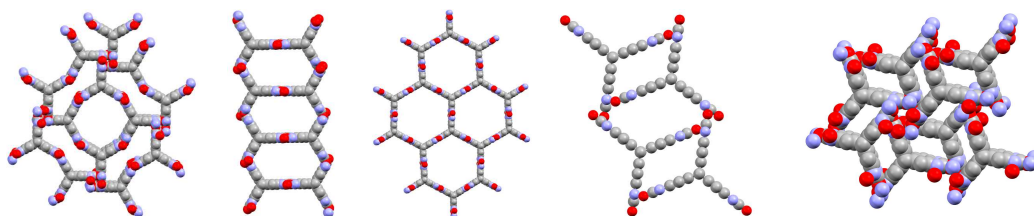


Fig. 10. The crystals T2-$\alpha$, T2-$\beta$, T2-$\gamma$, T2-$\delta$, T2-$\varepsilon$ based on the T2 molecule were synthesized for methane capture following the CSP (Pulido *et al.*, 2017) in Fig. 1.
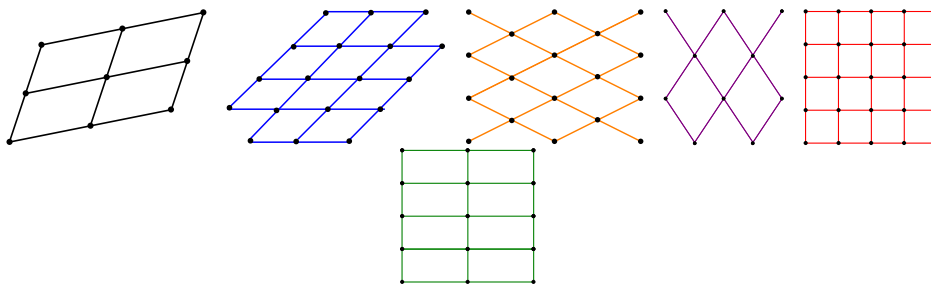
Fig. 11. 2D lattices whose AMD curves are in Fig. 12.    **1st**: a generic black lattice with the basis $(1.25, 0.25), (0.25, 0.75)$. **2nd**: the blue hexagonal lattice with the basis $(1, 0), (1/2, \sqrt{3}/2)$. **3rd**: the orange rhombic lattice with the basis $(1, 0.5), (1, -0.5)$. **4th**: the purple rhombic lattice with the basis $(1, 1.5), (1, -1.5)$. **5th**: the red square lattice. **6th**: the green rectangular lattice with the basis $(2, 0), (0, 1)$.
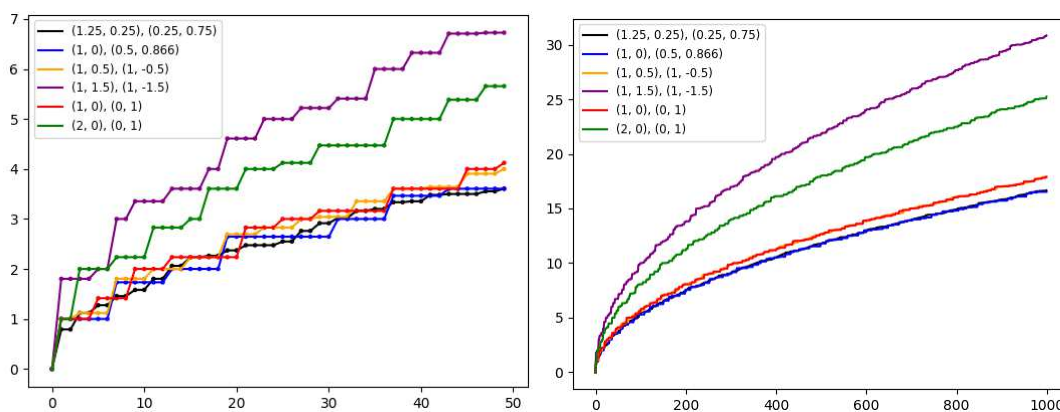


Fig. 12. **Left**: $AMD_k$, $k \in [0, 50]$, for the 2D lattices with given bases. **Right**: extended $AMD_k$ up to $k = 1000$. The orange and red graphs are very close as well as the blue and black graphs, see clearer differences for smaller $k$ on the left.

### References

Andrews, L., Bernstein, H. & Pelletier, G. (1980). *Acta Crystallographica A*, **36**(2), 248–252.

Andrews, L. C., Bernstein, H. J. & Sauter, N. K. (2019). *Acta Crystallographica Section A: Foundations and Advances*, **75**(3), 593–599.

Boutin, M. & Kemper, G. (2004). *Advances in Applied Mathematics*, **32**(4), 709–735.

Brown, R. A. (2015). *J. Computer Graphics Techniques)*, **4**(1), 50–68.

Chisholm, J. & Motherwell, S. (2005). *J. Applied Crystallography*, **38**(1), 228–231.

Flor, G., Orobengoa, D., Tasci, E., Perez-Mato, J. & Aroyo, M. (2016). *J. Applied Crystallography*, **49**(2), 653–664.

Hahn, T., Shmueli, U. & Arthur, J. W. (1983). *International tables for crystallography*, vol. 1.

Himanen, L., Jäger, M. O., Morooka, E. V., Canova, F. F., Ranawat, Y. S., Gao, D. Z., Rinke, P. & Foster, A. S. (2020). *Computer Physics Communications*, **247**, 106949.

Kerber, M., Morozov, D. & Nigmetov, A. (2017). *J Experimental Algorithmics*, **22**, 1–20.

Kitaigorodsky, A. (2012). *Molecular crystals and molecules.* Elsevier.

Lai, R. & Zhao, H. (2014). *arXiv:1406.3758.*

Liberti, L. & Lavor, C. (2017). *Euclidean distance geometry: an introduction.* Springer.

Morissette, S. L., Soukasene, S., Levinson, D., Cima, M. J. & Almarsson, Ö. (2003). *Proceedings of the National Academy of Sciences*, **100**(5), 2180–2184.

Mosca, M. & Kurlin, V. (2020). *Crystal Research and Technology*, **55**(5), 1900197.

Oliynyk, A. O., Antono, E., Sparks, T. D., Ghadbeigi, L., Gaultois, M. W., Meredig, B. & Mar, A. (2016). *Chemistry of Materials*, **28**(20), 7324–7331.

Patterson, A. (1939). *Nature*, **143**, 939–940.

Patterson, A. (1944). *Physical Review*, **65**, 195.

Pele, O. & Werman, M. (2008). In *European Conference on Computer Vision*, pp. 495–508.

Price, S. L. (2018). *Faraday Discussions*, **211**, 9–30.

Pulido, A., Chen, L., Kaczorowski, T., Holden, D., Little, M., Chong, S., Slater, B., McMahon, D., Bonillo, B., Stackhouse, C., Stephenson, A., Kane, C., Clowes, R., Hasell, T., Cooper, A. & Day, G. (2017). *Nature*, **543**, 657–664.

Rubner, Y., Tomasi, C. & Guibas, L. (2000). *Intern. Journal of Computer Vision*, **40**(2), 99–121.

---

## Synopsis

The Average Minimum Distances form an infinite sequence of real-valued crystal descriptors that are invariant under rigid motions and are provably stable under atomic vibrations.