

CrystalNets.jl: Identification and Classification of Crystal Topologies

Lionel Zoubritzky and François-Xavier Coudert*

Chimie ParisTech, PSL University, CNRS, Institut de Recherche de Chimie Paris, Paris, France

* fx.coudert@chimieparistech.psl.eu

February 28, 2022

Abstract

We present here an open-source Julia library for the topological identification of crystalline materials, with algorithmic and computational improvements over the previously available software in the field, resulting in a speed increase of one order of magnitude. This new algorithm and implementation can therefore be used at large scale in high-throughput screening methodologies. We have validated and benchmarked CrystalNets.jl against a diverse set of crystal databases, covering in particular metal–organic frameworks, aluminophosphates, zeolites, and other inorganic compounds.

Contents

1	Introduction	2
2	Algorithms and implementation	4
2.1	Periodic graph structure	4
2.2	Equilibrium placement	6
2.2.1	Definition	6
2.2.2	Computational aspects	7
2.3	Minimisation of the graph	8
2.4	Search space exploration	10
3	Reduction of the search space	14
3.1	The initial search space	14
3.2	Categorisation of vertices and edges	15
3.3	Reduction through symmetry	15
3.4	Ordering of keys	16
4	Application to materials databases	17
4.1	RCSR, EPINET and IZA-SC database	17
4.2	MOF structures	18
4.3	Database of aluminophosphates	20
4.4	Hypothetical zeolites	21
5	Conclusion and perspectives	22
	References	23

1 Introduction

In order to describe the average structure of crystalline materials at the nano-scale, reticular chemistry [1, 2] proposes to classify crystals, and in particular framework materials, in light of their topology. This approach consists in abstractly regrouping the atoms of a crystal into clusters, called *vertices*, bonded together by an abstract kind of linkers, called *edges* — the choice of grouping being informed by chemical sense. In simple compounds, such as dense oxides, the mapping can be trivial, with each atom as a separate vertex and each chemical bond as an edge. For zeolites and aluminosilicates, the tetrahedral atoms T are considered the vertices ($T = \text{Al, Si, P} \dots$) and the T–O–T linkages are the edges. In the case of metal–organic frameworks (MOFs), these vertices would be the inorganic *secondary building units* (SBUs) [3] and the edges, the organic SBUs for example. Any crystal can then be simplified into a topological structure, called a *net*, consisting of the sets of vertices and edges (see Figure 1).

Like the atomic structure it represents, the net is 3-dimensional and periodic. It contains all the information on the connectivity of the vertices, *i.e.* which pairs of vertices, called *neighbours*, are linked by an edge. Thus, it provides a convenient mathematical abstraction to classify crystal structures. However, it is important to note that this representation is dependent on the choice of vertices and therefore a crystal may be represented by different nets [8]. Moreover, reconstructing the full crystal from the net necessitates additional data, such as the unit cell parameters and the nature and exact geometric placement of the vertices and the edges [5]. A crystal may also be composed of multiple interpenetrating nets [9], which is common in ultraporous framework materials.

Several existing projects, such as the RCSR [7] or Epinet [10], aim at collecting the different nets underlying the structure of real or hypothetical crystals. This is key to allowing researchers to create new materials with a target topology, which is useful to obtain specific physical properties, given there is a strong link between materials topology and properties for a given chemical composition [11]. However, in order to check the accuracy of the collected nets with regard to existing crystals, one needs a way to determine the net from a given crystal structure.

There are two different approaches to solve this problem, which correspond to the two main softwares used in the field: TOPOS [12] and Systre [13]. TOPOS is a crystallographic tool that allows to retrieve an extensive set of topological properties from crystals. It is a proprietary software, but its topology detection tool can be accessed for free online [14]¹ for crystallographic files in CIF format [15]. Its strategy to identify the topology of a net consists in computing three specific properties for each vertex u :

- The coordination sequence up to distance 10. A vertex v is at distance k of u if there is a chain of $k + 1$ vertices $u = u_0, u_1, \dots, u_k = v$ such that $\forall i \in \llbracket 1, k \rrbracket, u_{i-1}$ and u_i are neighbours, and there is no smaller such chain. The k -th term of the coordination sequence is then defined as the number of vertices at distance k of u .
- The extended Schläfli symbol.
- The vertex symbol.

The latter two account for the shortest cycles and rings starting at vertex u [16]. It is found empirically that the combination of these three properties for all vertices forms an almost unique identifier of the topology. However, this approach is not completely sound as there can be different nets that do share these same three properties. A few other topological recognition tools specific to zeolites, like TOTOPOL [17], follow a method similar to that of TOPOS.

Systre is an open-source software²; its inputs are files representing the graph underlying the crystal structure, rather than the crystallographic files directly. It implements an algorithm

¹available online at <https://topcryst.com/>

²available online at <http://gavrog.org/>

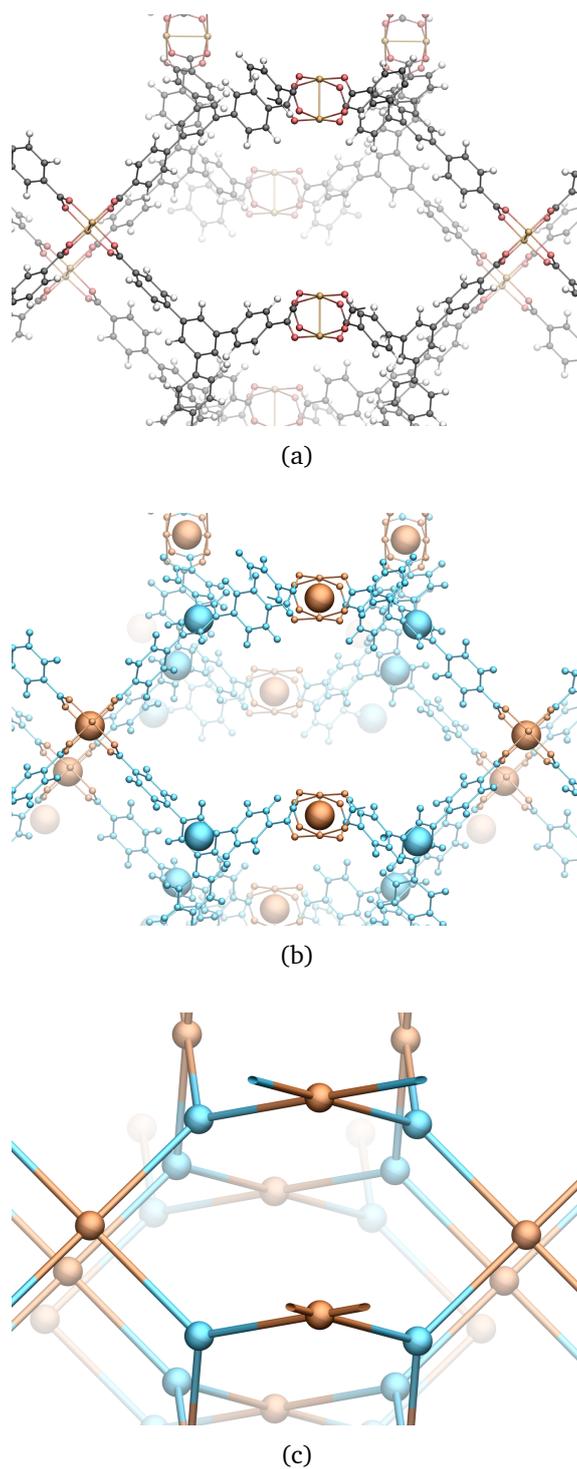


Figure 1: Decomposition of a crystalline framework into its underlying net:
(a) Crystalline structure of MOF-143 [4].
(b) Identification of the SBUs: Cu-based paddle-wheel [5] and benzene-tribenzoate (BTB) [6].
(c) The underlying net, called **pto** in the database of the RCSR [7].

that yields a so-called *topological genome* [18] of the net, which is a finite series of numbers. This genome is provably unique for each net, only depends on the net itself and not its representation, and can be computed in polynomial time of the size of the net. To do so, Systre starts by computing an *equilibrium placement* [19] of the vertices, which attributes to each vertex a position which is the barycentre of that of its neighbours. The net is then minimised, which amounts to finding a unit cell. Afterwards, Systre collects a number of *candidates*, which are frames of reference that do not depend on the current representation of the net. For each candidate, the vertices of the net are ordered according to their position in the frame of reference, and a key is determined which uniquely identifies the net with this ordering in this frame. The genome is then chosen as the lexicographically smallest key. Since the candidates do not depend on the initial representation of the graph, the keys do not either, hence the genome is indeed uniquely determined and only depends on the net itself.

In this context, we propose a new tool to efficiently determine the topology of crystal structures: the Julia package `CrystalNets.jl`. Julia [20] is a modern general-purpose and scientific programming language that offers both great performance and ease of programmability [21]. This, as well as some carefully chosen optimisations, allows `CrystalNets.jl` to outperform Systre while leveraging a very similar core algorithm to ensure soundness. Moreover, the package uses the Julia port of `chemfiles`³ which allows it to directly read common chemical file formats like CIF, making it easy to use for the analysis of single cases as well as large databases of crystal structures.

Delgado-Friedrichs and O’Keefe [13] explained the overall architecture of the Systre algorithm which `CrystalNets.jl` largely reuses: we will follow its structure, only briefly summarising the key concepts that are already explained elsewhere along the way, while focusing on the implementation choices and some proofs of correctness that do not appear in other publications. We will then focus on a series of new optimisations of the algorithm that allow to greatly improve its performance for most nets. Finally, we will demonstrate the applicability of the method by analysing several existing crystalline databases using `CrystalNets.jl`.

2 Algorithms and implementation

2.1 Periodic graph structure

Crystal nets are special instances of 3-periodic graphs. In general, a *graph* is a mathematical construct that contains objects, called *vertices*, some of which are linked together by *edges*. Two linked vertices are said to be *neighbours*. In a crystal, the pattern representing the atoms and their chemical bonds periodically repeats itself across three dimensions: for this reason, the graph representing a crystal is called 3-periodic.

More generally an N -periodic graph (with $N \in \mathbb{N}^*$) can be represented using the vector model [22]. This model conceptually follows the description of a crystal as a lattice and a pattern. Any point in the lattice is denoted by an offset $o \in \mathbb{Z}^N$, whose i -th coordinate is the offset along the i -th axis of the lattice compared to a fixed origin. An N -periodic graph is then defined by a pair (V, E) where V is the set of vertex identifiers and $E \subset V \times V \times \mathbb{Z}^N$ is the set of edges. Any vertex in space is uniquely identified by a pair $(v, o) \in V \times \mathbb{Z}^N$ where $v \in V$ is its vertex identifier and $o \in \mathbb{Z}^N$ is the offset of the cell that contains the vertex. Let’s note $n = |V|$ the number of vertices per cell and identify V with the integer interval $\llbracket 1, n \rrbracket$ so that each vertex identifier can be considered to be a number in $\llbracket 1, n \rrbracket$. When working in a fixed cell, the vertices will simply be referred to by their vertex identifier. Finally, for most of the rest of this article N will be equal to 3 since nets are 3-dimensional, but the algorithms exposed in the

³available online at <https://chemfiles.org/>

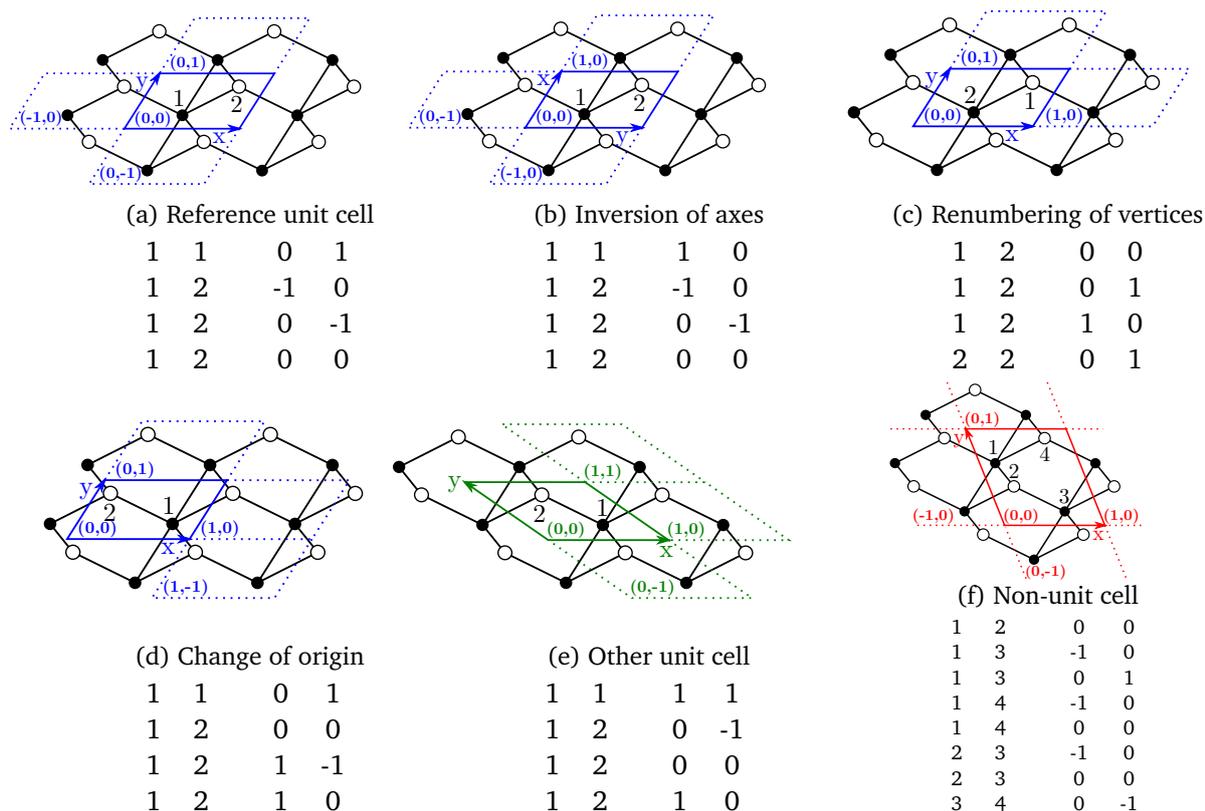


Figure 2: Six different representations of the same 2-periodic graph

rest of this section actually work for any dimension and CrystalNets.jl can identify 2-periodic and 1-periodic graphs as well.

An edge $e = (u, v, o)$ can be broken down into a source $u \in V$, and a destination $(v, o) \in V \times \mathbb{Z}^3$. The source can be seen as the vertex of identifier u and positioned in the cell at the origin of the lattice, and the destination is the vertex to which it is linked. When the source and the destination are in the same cell, the offset o is zero and noted $0_{\mathbb{Z}^3} = (0, 0, 0)$. Edges of the form $(u, u, 0_{\mathbb{Z}^3})$, called *loops*, are absent in nets because they do not convey any structural information.

For any edge $e = (u, v, o) \in E$, the reverse edge $(v, u, -o)$ is also in E : this property makes the graph *undirected*. All graphs used to represent crystals are undirected because chemical bonds are symmetric: if atom A is bonded to atom B , then atom B is bonded to atom A , and likewise for clusters of atoms.

As opposed to periodic graphs, *simple* graphs do not have a concept of offset: an edge of a simple graph is only a pair of vertices. Numerous computational representations of simple graphs exist in the literature, but they cannot be directly transposed to periodic graphs. In order to find a computational representation for crystals relevant to topological analyses, a useful observation is that nets share two distinctive features that come from chemistry:

- The number of edges is at most proportional to the number of vertices, because of the limited valence of atoms. Such graphs are generally called *sparse*. This remains true when vertices represent clusters instead of single atoms, although in that case the number of neighbours can be higher than the maximum valence of an atom — for instance, SBUs with up to 66 points of extension have been built. [3].
- The set of edges is small to medium-sized because the number of vertices per unit cell

is, in practice, between 1 and 10,000 in most real-world crystalline structures. This is small compared to the trillion of edges of some graphs in social media modelling for instance [23].

Finally, the most important operation on graphs that will subsequently be used is finding the set of neighbours of a given vertex. Driven by this, we chose to computationally represent a periodic graph by a list of size n whose u -th element is the set of neighbours of u : $N_u = \{(v, o) \in V \times \mathbb{Z}^3 \mid (u, v, o) \in E\}$. N_u is stored as a lexicographically sorted list.

Because of the previously mentioned undirectedness of the graph, if $(v, o) \in N_u$, then $(u, -o) \in N_v$. However, it is often convenient to be able to enumerate the edges of the graph without double counting an edge and its reverse. For this reason, the representation comprises another list of size n whose u -th element is the first index i in N_u such that $N_u[i] = (v, o)$ with either $v > u$ or $[v = u \text{ and } o > 0_{\mathbb{Z}^3}]$. Such an edge (u, v, o) will be called a *direct* edge. Since loops are forbidden, all edges are either a direct edge, or the reverse of a direct edge. Thus, since each list of neighbours is lexicographically ordered, reading the list of neighbours starting at the given indices i is guaranteed to yield all direct edges of the graph, hence each edge exactly once without its reverse. The lexicographically ordered sequence of direct edges is a unique identifier of a representation (V, E) of a periodic graph in the vector model: such a sequence will be called the *key* associated with the representation.

We implemented the graph structure detailed above in a companion Julia package, `PeriodicGraphs.jl`, which can be used independently of `CrystalNets.jl`.

Figure 2 shows several representations of the same 2-periodic graph. The graph is made of two kinds of vertices (represented by open and closed circles), but the choice of the vertex identifiers and of the cell influences the representation. Each choice of cell is characterised by its origin and its axes x and y ; some repeated cells are represented as well, along with their offset in the bottom-left corner. For each representation, the key is written below, each line referring to a direct edge. Recall that the Systre algorithm works by choosing the lexicographically smallest key from a set of representations, extracted from candidates that do not depend on the representation of the graph. Hence, if the set of representations was those shown in the figure, the topological genome of the graph would be the smallest key, *i.e.* that under representation (a).

Having a working representation for nets, let's now discuss our reimplementations of the Systre algorithm [13] to compute a unique identifier for each net topology.

2.2 Equilibrium placement

2.2.1 Definition

Given a periodic graph, the first step is to compute, for each vertex $(u, o') \in V \times \mathbb{Z}^3$, a position $p(u, o') \in \mathbb{R}^3$ such that

$$\forall u \in V, o' \in \mathbb{Z}^3, p(u, o') = \frac{1}{|N_u|} \sum_{(v, o) \in N_u} p(v, o' + o) \quad (1)$$

The set of positions p , called an *equilibrium placement*, is such that each vertex is at the barycentre of all its neighbours. p must also satisfy the additional constraint that

$$\forall u \in V, o \in \mathbb{Z}^3, p(u, o) = p(u, 0_{\mathbb{Z}^3}) + o \quad (2)$$

to enforce the periodicity of the equilibrium placement.

The equilibrium placement has been proven to be unique modulo isometries, because it is the solution of a full-rank linear problem [19]. **Equation 1** can be rewritten with **Equation 2**

in matrix form, $A \times P = O$, where $A \in \mathcal{M}_{n,n}(\mathbb{Z})$, $O \in \mathcal{M}_{n,3}(\mathbb{Z})$ and unknown $P \in \mathcal{M}_{n,3}(\mathbb{R})$ such that:

$$\begin{aligned} \bullet \quad \forall (u, v) \in V^2, A_{u,v} &= \begin{cases} |N_u| & \text{if } u = v \\ -1 & \text{if } u \text{ and } v \text{ are neighbours} \\ 0 & \text{otherwise} \end{cases} \\ \bullet \quad \forall u \in V, O_u &= \sum_{(v,o) \in N_u} o \end{aligned}$$

and from which we determine for all $u \in V$, $p(u, 0_{\mathbb{Z}^3}) = P_u$.

Since the coefficient of A and O are integers, the computed solution P is always in $M_{n,3}(\mathbb{Q})$ with the additional constraint that one of the vertices u be placed at the origin. The exact solution with coefficients in \mathbb{Q} is actually necessary for the subsequent computations and not an approximation, so the resolution of the system must be exact. This turns out to be a very costly operation: it requires storing integers of arbitrary size, which are computationally far more expensive than fixed-width machine integers.

2.2.2 Computational aspects

The GNU Multiple Precision Arithmetic Library (GMP) [24] provides access to heavily optimised structures that allow to efficiently perform exact integer and rational computations. Julia has native support for the arbitrarily-sized integer structure from GMP, through the `BigInt` type, but does not provide a full interface for GMP's rational numbers (although it is internally used for Julia's native `Rational{BigInt}` type). As part of this work, we provided this special support in a separate package, `BigRationals.jl`, which exports the `BigRational` type wrapping GMP's rational type.

In addition to the data structure used, a choice has to be made about the algorithm for solving this linear system. The most straightforward idea would be a general linear algebra technique such as LU decomposition followed by forward and backward substitution. Yet, matrix A is structured since it is both symmetric and sparse, calling for the use of more efficient algorithms.

Several algorithms have been discussed in the literature for the exact resolution of sparse and integer linear systems [25–28]. Following the analysis by Eberly *et al.* [25], we decided to implement Dixon's algorithm [27] which relies on a modulo-prime inversion of the matrix. In accordance with their conclusion, this algorithm largely outperforms a sparse LU decomposition algorithm (which we implemented as a reference and for benchmarking), both algorithms leveraging GMP instructions specific to the `BigRational` type for improved performance. Yet, Dixon's algorithm can fail if matrix A happens not to be invertible modulo the chosen prime: hence, we chose to successively invert it modulo three hard-coded large primes and stop at the first succeeding one, and finally resort to LU decomposition if all three fail. This implementation is up to orders of magnitude faster than that used in the original `Systre` program written in Java, especially for large graphs (with up to 920 vertices for the biggest net of the RCSR, called **tep**). This extra performance is crucial for the use of the topology detection algorithm as part of large-scale screening methodologies.

The equilibrium placement naturally embeds the graph in a Euclidean space, namely \mathbb{Q}^3 . The first vertex of the graph is placed in position $(0, 0, 0) = 0_{\mathbb{Q}^3}$, which is allowed because any translation of the equilibrium placement still results in an equilibrium placement. Moreover, because of Equation 2, for any vertex identifier $u \in V$ there is an offset o such that $p(u, o) \in [0, 1[$ ³. In order to simplify the representation, this offset is set to $0_{\mathbb{Z}^3}$ by adjusting the offsets of the edges of the graph accordingly. In this new representation, equations 1 and 2 still hold, as well as the new property:

$$\forall u \in V, p(u) \in (\mathbb{Q} \cap [0, 1[)^3 \quad \text{where } p(u) = p(u, 0_{\mathbb{Z}^3})$$

In order to accelerate the subsequent algorithms, the vertices of the graph are lexicographically sorted according to their new position. The first vertex, being at $0_{\mathbb{Q}^3}$, remains first.

Nets where at least two vertices have the same position are called *unstable* and, similarly to Systre, are not handled by our implementation. This is generally not problematic since their probability of occurrence is low [29]. Our implementation can however handle unstable nets where each collision site comprises at most four vertices if there is no edge between two different collision sites and no collision site is bonded to two representatives of the same vertex. This accounts for roughly half of the encountered unstable nets in the databases we explored. We will not detail the case of these unstable nets here but focus on stable ones. We can thus assume that no two vertices have the same position.

Figure 3 represents the equilibrium placement for the net of diamond (called **dia** in the RCSR). Note that the natural embedding in \mathbb{Q}^3 imposes a cubic unit cell for all nets, which is not representative of the actual unit cell of diamond; however, the geometry of the unit cell is irrelevant here because only the topology of the net matters. As a consequence, the bond lengths are not representative of that of diamond either, which explains why not all have the same lengths in this representation. Nonetheless, each vertex is indeed at the barycentre of the positions of its neighbours.

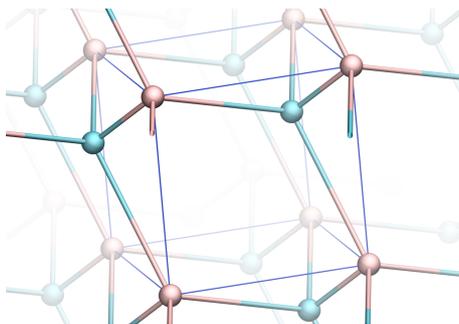


Figure 3: Visualisation of the natural Euclidean embedding of the equilibrium placement for the **dia** net.

2.3 Minimisation of the graph

Once the equilibrium placement is found, the next step of the algorithm is to find a unit cell. In general, a *cell* is a parallelepiped whose three axes $(\vec{a}, \vec{b}, \vec{c})$ are linearly independent translations that constitute symmetries of the graph. Such translations will be called *valid translations*. Another representation for this cell is the matrix A whose columns are respectively \vec{a} , \vec{b} and \vec{c} . Its volume is $\det(A) = (\vec{a} \wedge \vec{b}) \cdot \vec{c}$. It is a *unit cell* if it contains the smallest possible number of vertices, which is equivalent to having a minimal volume.

At this point, it is important to recall that the Euclidean space in which the graph is embedded through the equilibrium placement depends on its initial representation, so there is no simple way of finding a unique unit cell independently of the original representation. In particular, common cell reduction techniques such as Delaunay [30] or Niggli [31, 32], which yield a uniquely defined cell, do not solve this problem. Unfortunately, existing libraries implementing these techniques (like `spglib` [33] or `cctbx` [34]) cannot be used either for the simpler problem of finding any unit cell because they do not handle arbitrary precision.

First, note that, by construction of the equilibrium placement, the initial cell is the canonical basis of \mathbb{Q}^3 , so its matrix is the identity I . Let's now assume that I is not a unit cell and there is a unit cell $A \in \mathbb{Q}^{3 \times 3}$ of volume $\det(A) \in]0, 1[$.

Let's define $f : \mathbb{Q} \rightarrow \mathbb{Q} \cap [0, 1[$, $\frac{p}{q} \mapsto \frac{p \% q}{q}$ where $p \% q$ is the remainder of the euclidean division of p by q (hence $p \% q \in \llbracket 0, q - 1 \rrbracket$). Let $B = (f(A_{i,j}))_{1 \leq i, j \leq 3} \in \mathbb{Q}^{3 \times 3}$. By definition of f , the columns of B are axes within the initial cell; moreover, since $\forall x \in \mathbb{Q}$, $x - f(x) \in \mathbb{Z}$, each axis of B can be seen as the corresponding axis of A plus an offset $o \in \mathbb{Z}^3$. But since I is a repeating unit, such an offset is a valid translation, so each axis of B is a valid translation. Finally, since $\det(A) \notin \mathbb{Z}$, not all coefficients of A can be integers, so A has at least one coefficient x such that $f(x) \neq 0$, so at least one axis of B is not $0_{\mathbb{Q}^3}$.

To summarise, the previous observations show that if the initial cell was not a unit cell, then there exists a non-zero valid translation $t \in (\mathbb{Q} \cap [0, 1[)^3 \setminus \{0_{\mathbb{Q}^3}\}$.

Let's use this fact to find a unit cell by enumerating all such translations. Symmetries map vertices to vertices, so the first vertex (lexicographically sorted by position) is mapped to another. But since the first vertex was put in position $0_{\mathbb{Q}^3}$, the translation is directly given by the position of the other vertex. Hence, only $n - 1$ translations need to be checked, and if none is valid then the initial cell is a unit cell.

Checking whether a translation is valid can be done with $O(n \log n)$ time complexity. Vertices are already sorted by position, which allows to efficiently check that all the vertices are mapped to vertices. Doing so first allows to collect the map from initial vertices to their new counterparts and return early if the translation is invalid. Afterwards, it remains to be checked that all edges are also mapped to edges, which can also be done efficiently since the lists of neighbours N_u where $u \in V$ are also lexicographically sorted.

If I is not a unit cell, then a valid translation $t = (t_1, t_2, t_3) \in (\mathbb{Q} \cap [0, 1[)^3 \setminus \{0_{\mathbb{Q}^3}\}$ must have been found. Let $k \in \llbracket 1, 3 \rrbracket$ be such that t_k is the smallest non-zero coefficient of t . Then a new smaller cell is given by matrix $A^{(t)}$ where

$$A^{(t)}_{i,j} = \begin{cases} t_i & \text{if } j = k \\ 1 & \text{if } j \neq k \text{ and } i = j \\ 0 & \text{otherwise} \end{cases}$$

The volume of $A^{(t)}$ is $\det(A^{(t)}) = t_k \in]0, 1[$ and by constructions its axes are all valid translations, so A is indeed a new smaller cell. This process can then be iterated in the new Euclidean space created by mapping each point $x \in \mathbb{Q}^3$ to $(A^{(t)})^{-1}x$, until no valid translation is left.

There are several ways to improve the computational complexity of this procedure. The first one consists in observing that if $t \in \mathbb{Q}^3$ is a valid translation, then so is $k \times t + o$ where $k \in \mathbb{Z}$ and $o \in \mathbb{Z}^3$. This allows to minimise the volume of the new cell $A^{(t)}$ thanks to algebraic considerations. Let's illustrate this point in two dimensions for simplicity, although it can be generalised to any number of dimensions.

Let $t = (\frac{p_1}{q_1}, \frac{p_2}{q_2})$ be a valid translation. Let's assume that $p_1 \neq 0$ and $q_1 \geq q_2$ (the opposite reasoning can be done in the other case). Let $k \in \mathbb{Z}$ and $k' \in \mathbb{Z}$ be Bézout coefficients for p_1 and q_1 , i.e. such that $kp_1 + k'q_1 = 1$, and let $o = (k', 0) \in \mathbb{Z}^2$. Then $k \times t + o = (\frac{1}{q_1}, \frac{kp_2}{q_2})$.

Hence, matrix $A^{(kt+o)}$ can replace $A^{(t)}$ to reduce the cell to a volume $\frac{1}{q_1}$ instead of $\frac{p_1}{q_1}$.

Another potential improvement consists in using all the valid translations instead of stopping at the first one encountered. This way, the translation that minimises the volume of the new cell can be chosen instead of any valid translation. However the cost of collecting all valid translations can be up to n times higher than finding the first valid translation.

Furthermore, $A^{(t)}$ is just a special case of the set of possible new cells. In general, any invertible matrix whose columns are valid translations (including columns of the identity matrix)

is the matrix of a valid cell. The best new cell is the one with the smallest volume, but combinatorial exploration of all possible such matrices that lead to smaller cells has a worst-case complexity of $O(n^3)$. In practice, many candidate matrices can be eliminated before computing their determinant so using this method is actually quite efficient.

2.4 Search space exploration

At this stage in the workflow, the graph has a unit cell and a corresponding equilibrium placement. Let's now look for the topological genome. The strategy follows that implemented in Systre, which consists in taking the lexicographically smallest key among those computed for a number of candidates. Both the computation of the key from a given candidate and the selection of the candidates should not depend on the current representation of the graph.

A candidate is defined as a triplet of linearly independent edges, and represented by a pair $\{\text{vertex, basis of } \mathbb{Q}^3\}$ where the vertex is the source of the first edge and the columns of the basis are the translations induced by the edges. The process of selecting candidates is explained in [section 3](#). This subsection focuses on the computation of a key from a given candidate.

This precise algorithm is explained in [\[13\]](#) and analysed as Algorithm 5 in [\[19\]](#). It is detailed in [Figure 4](#). The main point is that the algorithm creates a new representation of the graph which is independent of the previous one by redefining both the unit cell and the numbering of vertices thanks to a unique euclidean space derived from the candidate. The first element of the candidate, the vertex, is the origin of the new space and the columns of the basis are its three axes.

Step **(b)** is a simple affine transformation of the positions and steps **(c)** to **(h)** are straightforward to implement, they only require keeping track of the positions of the numbered vertices. Step **(j)** is also straightforward, it consists in multiplying each translation by the inverse of the matrix of the new unit cell to obtain an integer offset. This allows to retrieve the set of edges E of the graph, and since $V = \llbracket 1, n \rrbracket$, this directly yields the representation of a periodic graph in the vector model.

Step **(i)** deserves special consideration. It consists in building a new unit cell from the set of positive translations t of the edges, positive in the sense of being lexicographically greater than $0_{\mathbb{Q}^3}$. Negative translations are negated to only keep positive ones. By construction, this set does not depend on the previous representation of the graph, it only depends on the candidate. As a consequence, any deterministic algorithm that transforms a set of translations into a unit cell can be used for this purpose. Delgado-Friedrichs and O'Keefe [\[13\]](#) mention this point and suggest triangulating the translations, without detailing their algorithm. Before presenting the one used here, let's sketch the general idea and prove a lemma.

Given their origin, translations of the set should correspond to integer offsets between vertices and their representatives in the new unit cell to be found. Thus, the matrix $B \in \mathbb{Q}^{3 \times 3}$ representing the unit cell must be such that each positive translation $t \in \mathbb{Q}^3$ of the set verifies:

$$B^{-1}t = o \in \mathbb{Z}^3 \quad (3)$$

Not all such basis B will do, since it should also represent a unit cell, hence the columns of B must be valid translations. Since all positive translations of the set map some vertex to one of its representatives, they are all valid translations, so it is tempting to look among the integer combinations of these translations to find the columns of B . However, it remains to be proven that such combinations indeed span the entire group of valid translations.

To show this, let's take the Euclidean space of step **(b)** and consider a path between the origin and any of its representatives. This path consists in a series of edges $(u, v, t) \in V \times V \times \mathbb{Q}^3$

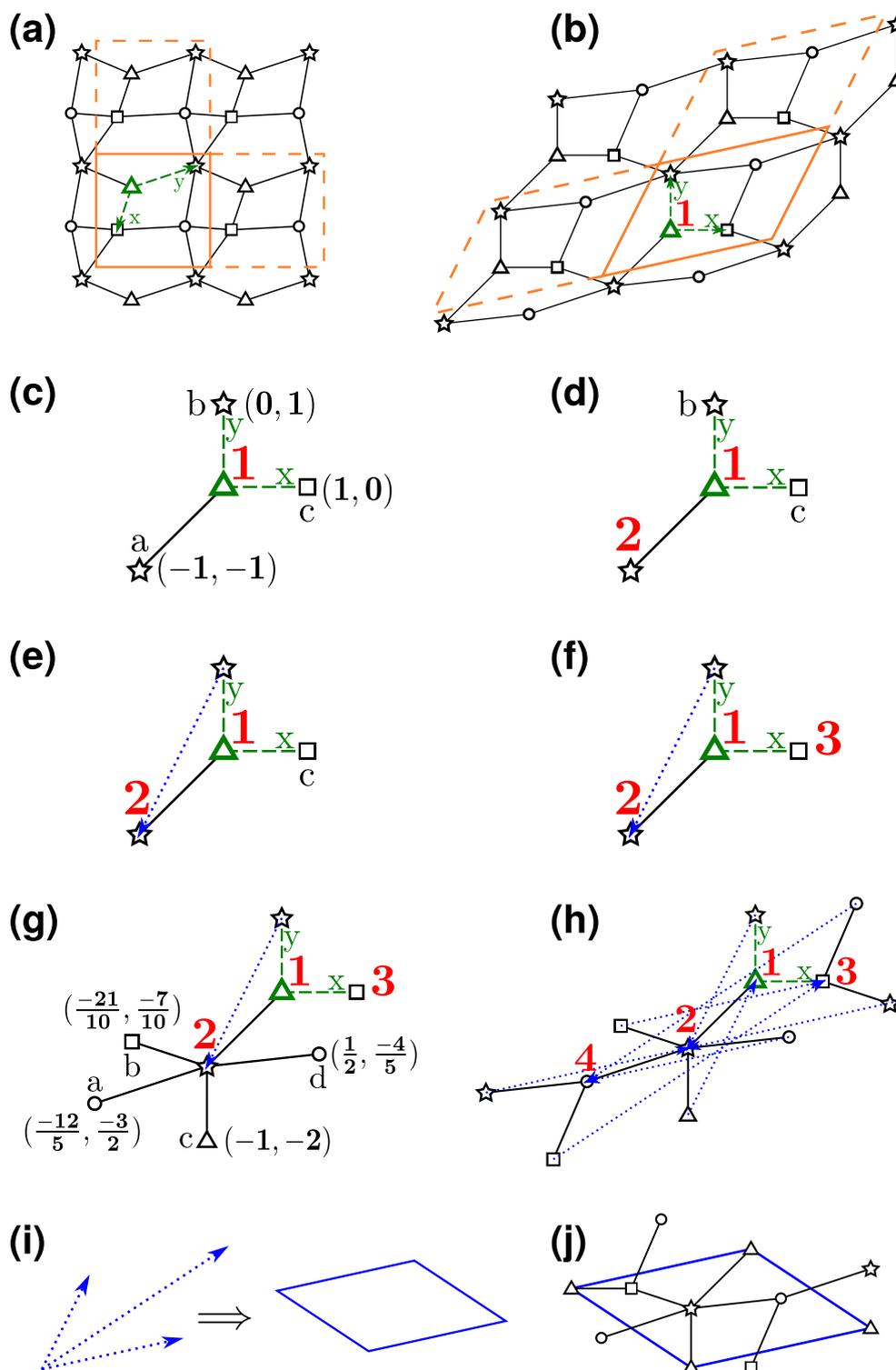


Figure 4: Illustration of the algorithm used to extract a key from a candidate on a 2-dimensional graph

- (a) Start from a candidate (in green), whose origin is Δ and axes are x and y . Initialise an empty queue of vertices.
- (b) Represent the graph in the basis given by the candidate (the former unit cell and two of its neighbours are represented in orange). Note that axes x and y are now orthogonal. Assign number 1 to the origin Δ and push it to the queue.
- (c) Take the first vertex u of the queue and sort its neighbours according to the lexicographical order of their position. Here, $u = 1$ and the sorted neighbours are a, b and c.
- (d) For each neighbour v' , in order, if it is not the representative of an already encountered vertex, assign it the first available number between 1 and n and push it to the queue. In this case, neighbour a is assigned number 2.
- (e) Otherwise, store the translation between the newly encountered representative v' and that which has been recorded v (the blue arrow). Here, v' is neighbour b and v is vertex 2: both are representatives of the same vertex.
- (f) Continue until all the neighbours have been processed. In this case, only neighbour c remains and it is assigned number 3.
- (g) Repeat the steps from (c) until the queue is empty. Here, the neighbours of 2 are being processed.
- (h) When the queue is empty, all numbers between 1 and n have been assigned and a set of translations has been recorded.
- (i) Triangulate the set of positive translations to obtain a new unit cell (see section 2.4).
- (j) Convert the translations to integer offsets using the unit cell to find a new representation of the graph, and extract the key.

Figure 4: *Continued:* Illustration of the algorithm used to extract a key from a candidate on a 2-dimensional graph

where t is the translation between vertex v and the vertex v' which is the actual neighbour of u . After any number of steps, let's note $(u, \tau) \in V \times \mathbb{Q}^3$ the pair where u is identifier of the vertex on which the path has just landed and τ is the sum of the translations t encountered so far. The initial pair thus is $(1, 0_{\mathbb{Q}^3})$, since 1 is the number of the origin. An induction on the number of steps shows that the position of the exact vertex on which the path lands when reaching pair (u, τ) is that of vertex u translated by τ . As a consequence, when reaching the final vertex of the path, the pair is $(u_{\text{end}}, \tau_{\text{end}}) = (1, \tau_{\text{end}})$. Thus, by construction τ_{end} is the sum of translations t stored in the edges during the algorithm, so it is an integer combination of the positive translations of the set. But since this reasoning was valid for the translation between the origin and any of its representatives, this proves that any valid translation can be written as an integer combination of positive translations of the set.

Let's now proceed to the algorithm that converts the set of positive translations into a matrix B satisfying Equation 3.

First, note that most positive offsets present in the usual representations of periodic graphs are one of either $(1, 0, 0)$, $(0, 1, 0)$ or $(0, 0, 1)$. As a consequence, in the most common cases, there are only three elements in the set of positive translations, which can simply be mapped to the three canonical offsets. This fact serves as a heuristic optimisation for the more general algorithm, by representing all the positive translations in the basis $A \in \mathbb{Q}^{3 \times 3}$ whose columns are the three first linearly independent positive translations. The reverse transformation will be applied at the end of the algorithm. Overall, this operation is not crucial for the algorithm to work and the input remains a set of translations, but it tends to reduce the denominators of the coefficients of the translations, which is computationally helpful.

Given a set of translations, the goal now is to find three combinations such that they form the columns of a matrix B satisfying Equation 3. In other words, these combinations must span the group of valid translations. This problem on rationals is lifted to integers by storing the least common multiple of the denominators of the coefficients of the translations columnwise, and pointwise multiplying the translations by this triplet of integers in order to obtain a set of integer translations.

From that point, several algorithms exist in order to solve the problem on this integer form [25–28]. The standard procedure is to build the rectangular matrix whose lines are the triplets representing integer translations and put it in Hermite normal form [35], which is precisely such that it is built from an integer combination of the initial triplets and still spans the same group as that of the triplets.

Determining the exact Hermite normal form is useless in this case however, since there is no need for the complexity brought by its other properties, and the existing solvers are specialised to have their best performance on big square matrices, not thin rectangular matrices one of its dimensions being only 3. Thus, we implemented a custom solver for this particular problem which performs the following steps, for column j from 3 to 1:

- Find Bézout coefficients of the j -th element of the integer translations, as well as their greatest common divisor $d \in \mathbb{N}^*$. To do so, we implemented an efficient algorithm from Havas *et al.* [35].
- Set the j -th column of B as the sum of the integer translations weighted by their Bézout coefficient.
- For each integer translation $t \in \mathbb{Z}^3$, update it to zero its j -th element by subtracting the new column of B times the quotient of the j -th element of t by d . By construction, after this step all translations have at most $(j - 1)$ non-zero elements.

Finally, columnwise divide B by the triplet of least common multiples that was found before to get back to rationals, and left-multiply it by A to go back to the Euclidean space specified by

the candidate.

The resulting matrix B has all the required properties for being a unit cell. Thus, each translation $t \in \mathbb{Q}^3$ of an edge now corresponds to an offset $o = B^{-1}t \in \mathbb{Z}^3$ with respect to the unit cell, so the lexicographically ordered set of edges $(u, v, o) \in V \times V \times \mathbb{Z}^3$ is a representation of the initial graph that does not depend on its initial representation, but only on the candidate.

Let's note that this method transforms a set of translations into a unit cell, and this is precisely the problem encountered when minimising the cell from the set of all valid translations in [subsection 2.3](#). However, in the current case, the procedure cannot rely on a starting cell to minimise, which is why it is more complex.

To conclude, let's observe that this representation might be impossible to embed, in the sense it may be impossible to put all the vertices of the graph, as they were determined by the algorithm, within the same unit cell given the constraints on their edges. This can be circumvented by taking appropriate representatives for each vertex given a fixed unit cell, but even that is unnecessary since the goal is only to extract a key, and it does not matter whether the key itself can be embedded.

3 Reduction of the search space

The only aspect of the algorithm that remains to be explored is the selection of the candidates, which is a crucial point for the overall complexity of the algorithm. After explaining how the implementation of Systre does this, we will present a series of novel reductions of the search space, *i.e.* of the number of candidates, that yield a notable performance gain.

3.1 The initial search space

Let's recall that a candidate is defined as a triplet of edges, which is used to create a new representation for the graph as explained in [Figure 4](#). The topological genome is chosen as the lexicographically smallest key extracted from the candidates, so the set of candidates must be chosen in a fashion that is independent of the initial representation of the graph, lest the selection biases the search space. As a consequence, Delgado-Friedrichs and O'Keefe [[13](#)] proposed the following selection pattern:

- If there exists at least one vertex such that its outgoing edges are not all coplanar, then the set of candidates is the set of triplets of linearly independent edges that all start at the same vertex.
- Otherwise, the candidates are all triplets of linearly independent edges such that the two first start from a common vertex and the third starts from another vertex.

In both cases, the set of valid candidates is not empty (otherwise it would mean that the graph is either not connected or not 3-periodic) and independent of the initial representation of the graph. It can be further reduced by removing duplicate candidates that yield the same key, using symmetry considerations as will be shown in [subsection 3.3](#).

For many graphs, this search space remains rather large compared to the set of all possible triplets of edges. Our goal in this project was for our software to be able to tackle very big structures very efficiently, for two different types of applications: high-throughput screening of hypothetical materials, where the topological detection is performed on hundreds of thousands (and sometimes several millions) of structures [[36](#)], and application to very large periodic structures, such as models of amorphous framework materials [[37](#)]. Therefore, reducing this search space is crucial for performance since it directly tackles one of the bottlenecks in the original implementation of Systre.

3.2 Categorisation of vertices and edges

Since the search space must remain independent of the representation of the graph, it can only depend on some properties that are common to all its representations. Such a property is called a *topological invariant*. The only invariant used so far is whether there exists a vertex whose outgoing edges are not all coplanar, but many other invariants exist, such as the three used by TOPOS [12]. These properties can be used to considerably reduce the search space by separating vertices according to them.

The core idea of vertex categorisation consists in trying to separate vertices into as many categories as possible, each of these categories being identifiable independently of the representation of the graph. For instance, in this work, vertices are separated according to their coordination sequence, so that two vertices are in the same category if and only if they share the same coordination sequence.

Vertex categories are then sorted by the total number of outgoing edges from vertices in that category and, in case of equality, by the invariant they were defined from. For example, ties are broken in our case by lexicographically sorting according to the coordination sequences.

With these sorted categories in mind, let's look for candidates and try to minimise their number. The following strategy is used in our improved algorithm:

- Let c be the first vertex category such that there is a triplet of non-coplanar edges starting from the same vertex in c (if no such category exists, skip to the alternative below). Then, for all such triplets, let c_1 , c_2 and c_3 be the sorted categories of the vertices to which the initial vertex is bonded with these edges. There are four subsets of such triplets, characterised by whether they satisfy:
 1. $c_1 = c_2 = c_3$.
 2. $c_1 < c_2 = c_3$.
 3. $c_1 = c_2 < c_3$.
 4. $c_1 < c_2 < c_3$.

Then the candidates are all the triplets of edges whose categories are in the last non-empty such subset with (c_1, c_2, c_3) lexicographically minimal.

- Alternatively, if all vertices have only coplanar outgoing edges, let c be the first category such that there exists another category c' such that there is a non-coplanar triplet of edges, the two first of which start from a vertex of c and the last from c' . c' is chosen as the first such category. Let c_1 , and c_2 be the sorted categories of the destinations of the two first vertices, and c_3 be the destination of the third. Then the candidates are all the triplets of edges, the two first of which start in c and the last from c' such that (c_1, c_2, c_3) is lexicographically minimal.

The main difference with the strategy of Systre is that the categorisation of vertices bounds the number of candidates with respect to the maximal size of a category, hence the goal of having as many small categories as possible. However, the two strategies are exactly the same in the worst-case scenario where there is only one category. Note that, although rather complex, the constraints on the retained candidates described above can actually be implemented in an efficient way, which is necessary to prevent the search space reduction from taking more time than the actual search space exploration. We provide a multithreaded implementation of the first case to slightly further reduce the time taken by this procedure.

3.3 Reduction through symmetry

In addition to the individual properties of vertices which are used to separate them into categories, symmetry considerations on the equilibrium placement can be useful in reducing the

search space. Indeed, if two candidates are symmetry-related, then the representations extracted from both will be exactly identical. This optimisation is already used in Systre, and in the current work, it is useful to perform a symmetry analysis on the equilibrium placement of the unit cell before even categorising vertices.

Unfortunately, we do not know of any library that can find symmetries for an input given in arbitrary precision; all existing libraries take approximate positions, and thus return only approximate symmetries. We use the *spglib* library [33] to this purpose. *spglib* is a C library, which makes it easy to interface with Julia without sacrificing speed. Since all its inputs and outputs are approximate, the returned symmetries need to be checked against the exact equilibrium placement of the graph to rule out false positives. Because of the imprecision of the input, there might also be some false negatives – that is, some symmetries may be missing – but they do not matter because symmetries can only help to eliminate candidates that are symmetric images of other ones already encountered. Hence, the use of symmetries is purely an optimisation, but it has no influence on the result of the algorithm.

Symmetries are useful to avoid needlessly recomputing coordination sequences when dividing the vertices into categories. Indeed, the search space is reduced by singling out one vertex per symmetry orbit, and using the convention that only these can serve as the source of the first edge of a candidate. No such constraint can be placed on the destinations of the edges unfortunately, as that would be unsound.

3.4 Ordering of keys

Finally, the algorithm can be made faster by changing the mode of comparison of the candidates. In Systre, the topological genome was defined as the smallest key extracted from a candidate. However, we observe that, arrived at step **(h)** of the algorithm in Figure 4, both the embedding of the graph as well as the ordering of its vertices only depend on the candidate. Hence, we can define at that point a pseudo-key consisting in the sorted list of triplets $(s, d, t) \in V \times V \times \mathbb{Q}^3$ if there is an edge between the vertex defined by its number s and the representative of vertex number d which is translated from the latter by t . This pseudo-key only depends on the intermediate representation of the graph obtained in step **(h)**, which means it is exclusively determined by the candidate, so we can define the genome as the key extracted from the candidate with the smallest pseudo-key.

Instead of comparing the keys extracted after a costly lifting of the positive translations to a proper representation of the graph, it is thus sufficient to compare the pseudo-keys. Moreover, this pseudo-key can be computed incrementally during the main loop of the algorithm. To do so, at step **(c)**, the neighbours of vertex u should be sorted in a different manner (represented on Figure 5): all representatives of the same vertex should be contiguous and lexicographically sorted according to their position within each category; the categories corresponding to vertices which have a number (assigned in step **(d)**) are placed first and sorted according to this number, while the others are placed afterwards and sorted according to the position of the first vertex in the category. This makes it so that each edge of the graph is added to the pseudo-key in the order in which it is studied (in steps **(d)** – **(f)**), while preserving the invariant that the pseudo-key is lexicographically sorted.

Since the pseudo-key is constructed on-the-fly, it can thus be compared at each step with the previous smallest pseudo-key. If a triplet of the current pseudo-key is lexicographically greater than the triplet at the same position in the previous best pseudo-key, the computation for the candidate is aborted; if it is lower then this comparison can be skipped for the rest of the computation and the pseudo-key will be the new smallest. Finally, once all the candidates have been checked, the minimal pseudo-key is converted to an actual key through steps **(i)** and **(j)**, yielding the topological genome of the graph.

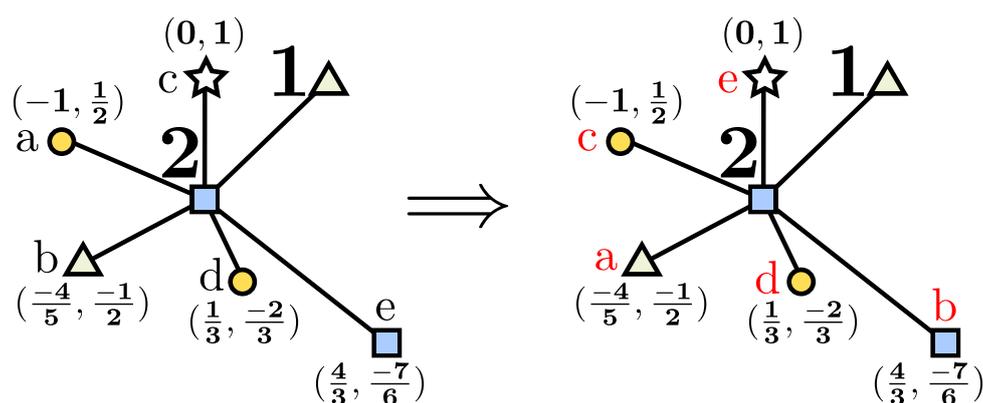


Figure 5: Difference between previous and new ordering of neighbours at step (c) on an example. Representatives of the same vertex have the same symbol.

4 Application to materials databases

CrystalNets.jl is a Julia package which can be intuitively used from an interactive Julia shell or as an executable to determine the structure of single crystallographic files in the variety of file formats recognised by chemfiles⁴. Since chemical bonds are not systematically present in such files, it uses a custom bond detection algorithm inspired from that of chemfiles, and can export both the input with the detected bonds as well as the extracted underlying net to allow visually checking the correspondence with the input.

By combining the bond detection heuristic with topology determination, CrystalNets.jl can also be used to survey large databases of crystalline structures, without the need for user input or relying on other chemical information (such as choice of SBUs) on the structures themselves. In this section, we demonstrate the code by exploring a series of available materials databases.

4.1 RCSR, EPINET and IZA-SC database

The Reticular Chemistry Structure Resource (RCSR) [7]⁵ is a repository of 3,132 structures with 2 or 3-periodicity (excluding weaving and 0- and 1-periodic structures), each given a canonical name. It offers a Systre archive containing 2,928 nets that we reprocessed to obtain the 2,928 topological keys corresponding to each structure. These keys are not the same as those of the archive because of our algorithmic differences explained before. Such a few thousand periodic graphs are processed in two minutes on a single thread on a laptop.

The RCSR partially overlaps with the larger EPINET database [10]⁶ which contains, among other data, 14,532 stable 3D nets, which we processed similarly. Finally, both databases partially overlap that of Structure Commission of the International Zeolite Association (IZA-SC), which provides names for 255 zeolitic framework structures and corresponding CIF files.⁷ These structures have unambiguous vertices, defined by the tetrahedral atoms, and edges between two such atoms bonded by an oxygen, thus a uniquely-defined underlying net.

The results for all three databases are stored within CrystalNets.jl, so that they may be used to recognise the topology of structures processed by the code. In the following, we call *unknown* nets those that are not part of any of these databases — although it should be noted that other complementary databases of nets do exist (in particular the proprietary databases

⁴full list available at <https://chemfiles.org/chemfiles/latest/formats.html#list-of-supported-formats>

⁵available online at <http://rcsr.anu.edu.au/nets> (3D) and <http://rcsr.anu.edu.au/layers> (2D)

⁶available online at <http://epinet.anu.edu.au>

⁷available online at <http://www.iza-structure.org/databases/>

used by TOPOS).

Finally, we note that this processing of the databases exposed a few mismatches between the names of the structures given in the three databases, which were all reported. Namely, the RCSR nets **ucn**, **rad**, **jst**, **fu**v and **fv**b have as IZA-SC names SBN, ATS, JST, NPT and PTT, which did not appear on the RCSR. Additionally, we noticed that 137 structures from EPINET were missing their RCSR name, while 6 RCSR structures mentioned an incorrect EPINET counterpart. This highlights the need for systematic and automated checks of material and topology databases for consistency.

4.2 MOF structures

In order to test CrystalNets.jl on a realistic workload, we surveyed the CoRE MOF 2019 [38] database of metal-organic frameworks (version 1.1.3). To do so, we implemented a custom clustering algorithm to regroup atoms into SBUs. Since the definition of MOF SBUs is not universal [3] and IUPAC recommends providing complementary topological descriptions when necessary [39], we implemented several clustering options to account for different strategies used in the literature.

We provide five main clustering options: “all nodes“, “single nodes“, “standard“, “PE” and “PEM“. Except for the “standard” representation which is uniquely defined by its implementation in TOPOS [40, 41], the other clustering options are only described but not uniquely determined by an explicit algorithm [8, 39, 40, 42, 43]. We also chose not to follow the original implementation of the “standard” representation [41], primarily because of the required computational cost of determining the minimal ring around each bond. Instead we use a custom algorithm whose main steps are as follows:

- First, organic cycles are replaced by a new virtual carbon atom. Afterwards, all metallic atoms are assigned to the “inorganic” class and all carbon atoms to the “organic” class. Then, in a loop, each remaining atom that is bonded to an atom in the “inorganic” class is itself classified as “inorganic”, until no such atom remains. All remaining unassigned atoms are then classified as “organic”. P and S atoms are treated differently and considered organic if directly bonded to a carbon, metallic if directly bonded to a metal, and part of a third class otherwise.
- Afterwards, SBUs are formed by regrouping together all adjacent atoms belonging to the same class. Paddle-wheel patterns are identified separately and regrouped into a single inorganic SBU even when the two opposite metal centres are not bonded together. The *points of extension*, which are organic atoms bonded to an inorganic atom, form their own SBU each. The algorithm stops here if there is no periodic SBU, *i.e.* if there is no infinite set of adjacent atoms belonging to the same class.

Otherwise, if there are such periodic SBUs, they are split by virtually removing the atoms with the highest connectivity and regrouping the other atoms of the periodic SBU into the remaining connected components. A few other heuristics are used to ensure that all final SBUs are finite, and to tackle the various possible configurations encountered in framework structures.

This algorithm yields our “all nodes” representation [39]. It can be converted to “PE” (standing for “points of extension”) [44] by removing all metallic atoms and bonding their corresponding points of extension to describe the coordination polyhedra around the removed metals. The “all nodes” representation is also convertible to “single nodes” [39] simply by merging the organic points of extension with their neighbouring organic SBU. It can also be transformed to yield the “PEM” (for “points of extension and metals”) representation [44] by

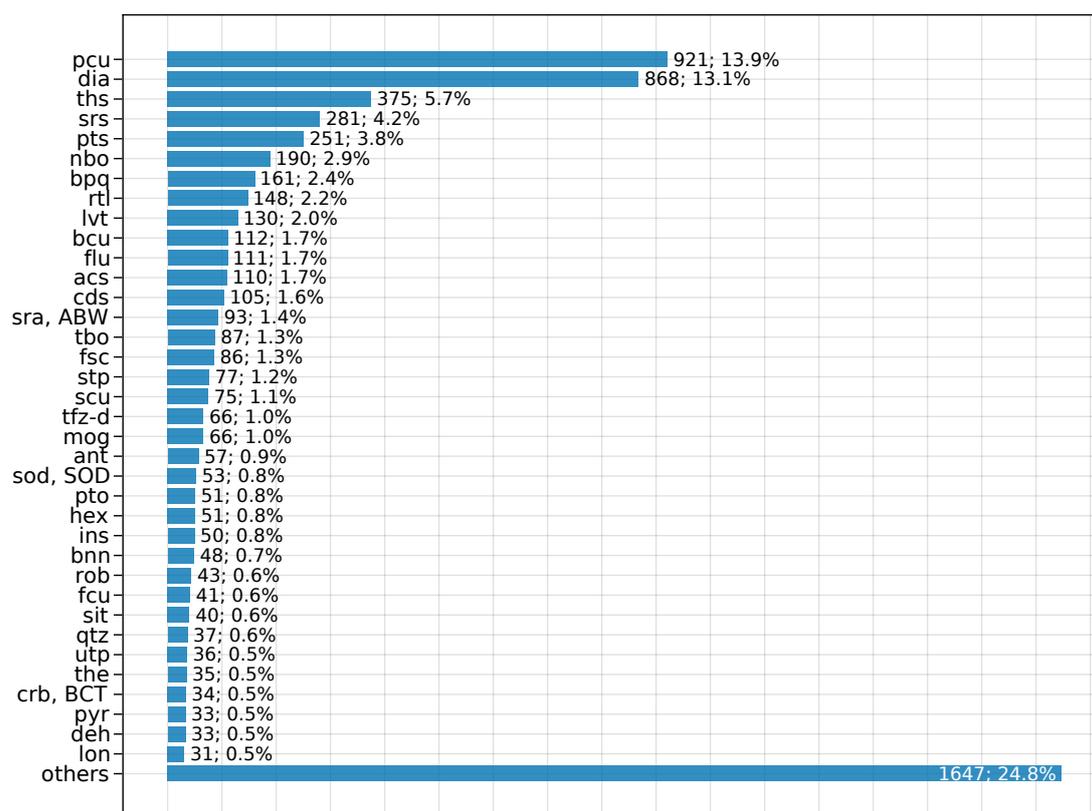


Figure 6: Number of structures per known net in the ordered structures of ASR.

splitting inorganic clusters to keep only one metal atom per SBU. Combining the algorithm for “PEM” and “single nodes” yields our “standard” representation [40]. For simple cases, we also provide an “each vertex” clustering option which bypasses the entire clustering step and keeps each atom as its own SBU, as well as an “input” option which determines clusters directly from the Residues of the chemfiles-parsed input file, so that the user may provide their own SBUs.

Once the SBUs are determined, they are converted to vertices and bonded together according to the atomic bonds. Finally, as with all nets in CrystalNets.jl, if a vertex has one or zero neighbour, it is removed, and if it has exactly two neighbours it is converted into an edge between its neighbours, and this iteratively until convergence. Once vertices and edges are defined, interpenetrating nets are separated, each of them accounting for a single structure.

To simplify the presentation, we will only discuss the results when using the “single nodes” approach, as well as the “MOF” structure option which widens metallic bond radii. CrystalNets.jl recognises 60% of all 21,691 resulting structures of CoRE MOF 2019 as one of the nets in either the RCSR, the IZA-SC or the EPINET databases (only 33% structures are recognised with “PEM” and between 45% and 55% with the other clustering options). Among the unrecognised structures, less than 2% correspond to unstable nets which CrystalNets.jl does not handle. These proportions are similar for each of the four subcategories of the database, divided whether only free solvent was removed (FSR) or all solvent molecules (ASR), and within each, whether the structure is disordered or not. The closeness of these numbers is expected since most structures from FSR also belong to ASR with only molecules of solvent as difference, which should not affect the topology. This highlights the robustness of the algorithm.

Next, we focus on the 11,189 ordered structures of the ASR subset to avoid duplicates. By studying the distribution of the number of structures per corresponding net, we observe that

this distribution drops faster for unknown nets than for known ones. This is not surprising since we expect known nets to correspond to the most common occurring in nature, whereas unknown nets should only correspond to at most a few corresponding structures. Indeed, among the latter, 1,489 unknown nets only correspond to one MOF structure each, and 302 to two structures. Yet, 5 of them still correspond to thirty MOF structures or more each: we manually checked through the TopCryst interface [14] that these 6 nets were part of the TOPOS database, but we could not do so for all the unknown nets since these databases are not open. In comparison, for known nets, 110 of them only correspond to a single MOF structure, 54 to two and 38 to thirty or more each. These most represented nets are detailed in Figure 6: **pcu** (primitive cubic lattice), **dia** (diamond) are by far the most common, together accounting for 27% of know nets, followed by **ths**, **srs** and **pts**. For zeolitic nets, ABW, SOD (sodalite) and BCT are the most commonly encountered.

The high number of nets with known topologies is a good indicator of the quality of our clustering algorithm. Yet, since it partly relies on heuristics, it should not be considered fail-proof. To assess its validity, CrystalNets.jl exports both the input and the detected net in VTF format, which can be superposed in a molecular visualisation software to check the correspondence between the vertices of the nets and the expected clusters. Moreover, we note that since our analysis relies solely on CrystalNets.jl, which is independent from the tools used to build and clean the CoRE MOF database, it allowed us to identify about 80 improper structures in the database, which were reported.

4.3 Database of aluminophosphates

Additionally, we surveyed the topologies of Zheng, Li and Yu’s database of aluminophosphates [45]. These 385 experimental structures are given as CIF files with little to no cleaning from the publications they stem from. In particular, solvent is included and the structures can contain multiple overlapping atoms with partial occupancy, while atoms of the main structure are often partially substituted. CrystalNets.jl automatically removes 0-dimensional residues like solvent and keeps for each site only the atom with the highest occupancy. Using the “Zeolite” structure option of CrystalNets.jl, we recognise 64% of the 205 3D nets in the database. This number rises to 73% by adding a heuristic to remove P–O–P and Al–O–Al linkages, which are usually not accounted for when computing the topologies of AlPOs.

The 31 known nets represented on more than one structure in the database are represented on Figure 7. With the exception of **fel** (originating from feldspars) and **cag** (occurring in variscite), they are all referenced in the database of zeolite frameworks of IZA-SC. This is expected since AlPOs are structurally close to zeolites, with the replacement of Si atoms by Al and P. Conversely, it shows that **fel** and **cag** may be nets underlying the structure of unreported zeolites, especially since both only have 4-coordinates vertices, corresponding to the T-atoms, and can be embedded to exhibit pores.

Most known structures have a net represented only once (45 nets *i.e.* 30%) or twice (19 nets *i.e.* 25%). In comparison, among the unknown nets, 1 occurs on three structures, 5 on two and the 32 others appear only once each. Both distribution tails and comparison of known and unknown nets follow a pattern similar to the CoRE MOF database.

For the structures for which the underlying net has been recognised as one of the IZA-SC, there is little doubt on the validity of the detected net. For the others however, there is no systematic way of determining whether the detected net is indeed that representing the structure or not except from visually checking the adequate superposition of the detected net and the initial framework for each file. This limitation comes from the lack of cleaning of the input CIF files, which can lead CrystalNets.jl to interpret solvent residues as part of the framework for instance, or duplicate atoms as two different atoms instead of one. This issue

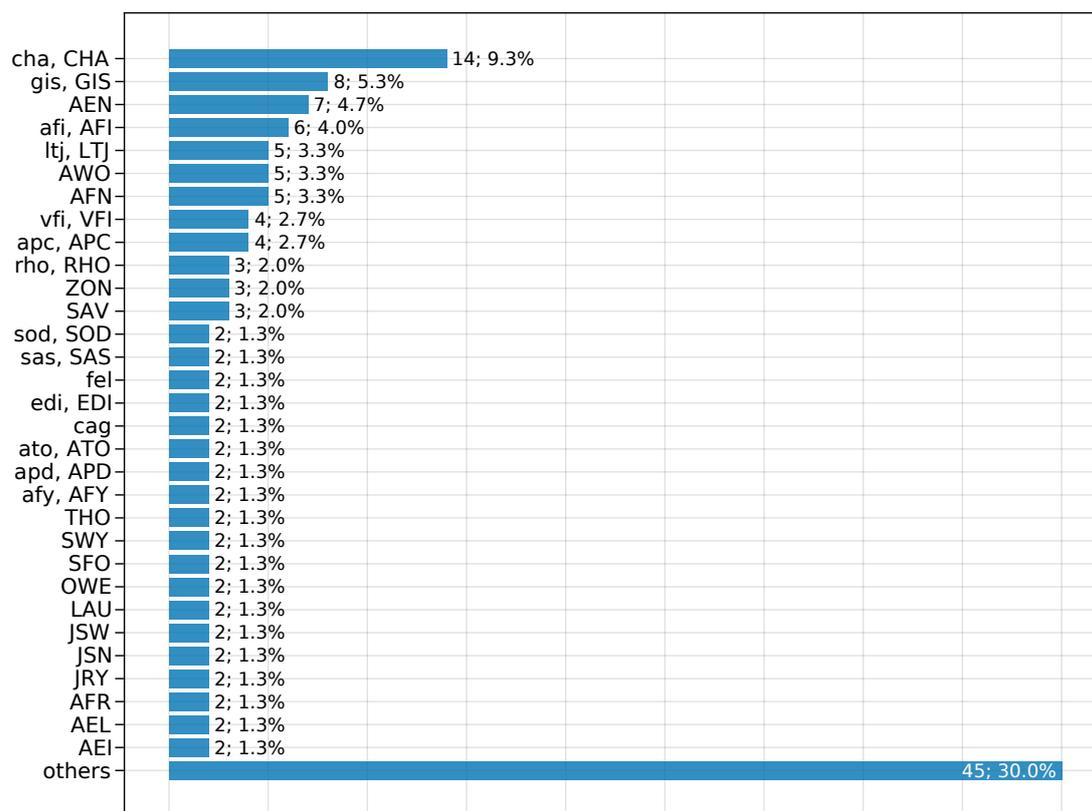


Figure 7: Number of structures per known net in the ALPO database.

can be circumvented by either cleaning the CIF inputs, as in the CoRE MOF database, or by manually checking the correspondence if working on individual files. Obviously, this limitation is inherent to the quality of the input files and affects all topology determination codes.

4.4 Hypothetical zeolites

The study of framework topology has historically been extensively applied to zeolites [46]. One reason is that the structure of zeolites, made only from tetrahedral ($T = \text{Si}, \text{Al}, \dots$) atoms bridged by $T\text{-O-T}$ linkages, is simple enough to show an almost one-to-one correspondence between each framework and its underlying net. Using this fact, large databases of hypothetical zeolites have been devised: the core idea consists in embedding 4-connected nets in real space, either by enumeration of the nets [47] or by Monte-Carlo sampling of vertex positions [48]. The vertices correspond to Si atoms while O atoms are placed along the edges. Each obtained structure is subsequently refined by minimizing its energy through molecular dynamics and with simulated annealing.

Since these structures are built from their underlying net, we tried to perform the opposite operation on the Predicted Crystallography Open Database (PCOD) [48] to assess our implementation. This database contains 331,172 CIF files representing zeolites with energies within 30 kJ/mol (per SiO_2 unit) of α -quartz, which were analyzed by CrystalNets.jl in 45 minutes using 20 cores (Intel Xeon Silver 4210R CPUs, 2.4 GHz). Bond guessing forced O atoms to have at most 2 neighbours, which is the default heuristic for O atoms not bound to metals: this incidentally ensures that all vertices of the net correspond to Si atoms.

As the SiO_4 tetrahedra are somewhat distorted during the refinement steps, the bonds guessed from the resulting structure do not always correspond to those from which it was

built. As a result, we found 11 frameworks whose T-atoms were not all 4-connected, and 75 other frameworks whose coordination sequences did not match that recorded in the database. Two structures were found to have the same underlying net: 8013876 and 8035701, the latter having different detected coordination sequences than indicated. Barring these 86 mismatched structures, 18 nets only have a 2D periodicity and 31 are unstable nets not handled by CrystalNets.jl. All those “problematic” nets represent less than 0.04% of all nets in the database. For the rest, no two structures had the same underlying net.

Indeed, by construction of the database, structures whose topology is already part of the database are removed. The invariants used to identify the topology are the coordination sequences and the ring sizes: this strategy does ensure that no two structures with the same net are present, but it is too restrictive since two distinct nets could share these invariants. In the future, the use of a unique topological identifier through CrystalNets.jl could thus potentially lead to the generation of more structures.

5 Conclusion and perspectives

We have presented CrystalNets.jl, a Julia library implementing topology identification on crystalline materials. The core algorithm for topology detection is strongly inspired from that used in Systre [13] which computes a topological genome, that is a sequence of numbers uniquely identifying the net underlying the given structure. The algorithmic and implementation improvements of CrystalNets.jl still preserve the soundness and polynomial time complexity of the initial algorithm, while offering a much better performance across all kinds of nets. On the **tep** net for instance, which is the largest of the RCSR, CrystalNets.jl computes a topological genome in 11 s, compared to 25 minutes for the original Systre program implemented in Java, and 192 s for its stripped-down implementation in Javascript. For the much smaller **pcu**, CrystalNets.jl takes 10 ms, which is half the time taken by the Javascript implementation of Systre. Moreover, our software directly takes crystallographic files as input and can automatically detect bonds if they are absent; it also includes a general clustering algorithm with different target options for framework materials, and heuristics adapted to metal–organic frameworks specifically.

Overall, CrystalNets.jl was successfully tested against a large diversity of crystalline structures: MOFs, aluminophosphates, zeolites and many other crystals. Even on large databases, our implementation shows great performance and can reliably recognise any structure in the Systre archive of the RCSR, in the IZA-SC database of zeolite structures and in EPINET. It successfully attributes known topologies to 60% of the structures in CoRE MOF after automatically clustering the vertices according to the “single nodes” policy. It was also able to recognise 73% aluminophosphate structures from Zheng *et al.*'s database [45], automatically removing solvent residues and colliding atoms when possible.

We believe CrystalNets.jl can be used as a faster and more convenient alternative to Systre and TOPOS for the identification of nets underlying crystalline structures. Its programmatic interface can be integrated to any computational pipeline written in Julia, as well as some other languages through bindings. Looking forward, we plan on providing a website interface similar to TopCryst [14], which could be useful to compare any structure to one of those surveyed in the previously exposed databases. We also plan on including results from the topological survey of the COD [49] and CSD [50] databases. The latter could then serve as comparison for our bonding and clustering algorithms against those of TOPOS [41].

Acknowledgments

We acknowledge financial support from the Agence Nationale de la Recherche under project “MATAREB” (ANR-18-CE29-0009-01) and from ANRT Cifre project 2021/1151 in association with Air Liquide.

References

- [1] N. W. Ockwig, O. Delgado-Friedrichs, M. O’Keeffe and O. M. Yaghi, *Reticular chemistry: Occurrence and taxonomy of nets and grammar for the design of frameworks*, ChemInform **36**(24) (2005), doi:[10/c5pqkx](https://doi.org/10/c5pqkx).
- [2] O. M. Yaghi, M. O’Keeffe, N. W. Ockwig, H. K. Chae, M. Eddaoudi and J. Kim, *Reticular synthesis and the design of new materials*, Nature **423**(6941), 705 (2003), doi:[10/bbfkw4](https://doi.org/10/bbfkw4).
- [3] D. J. Tranchemontagne, J. L. Mendoza-Cortés, M. O’Keeffe and O. M. Yaghi, *Secondary building units, nets and bonding in the chemistry of Metal–Organic frameworks*, Chemical Society Reviews **38**(5), 1257 (2009), doi:[10/br3pwj](https://doi.org/10/br3pwj).
- [4] H. Furukawa, Y. B. Go, N. Ko, Y. K. Park, F. J. Uribe-Romo, J. Kim, M. O’Keeffe and O. M. Yaghi, *Isorecticular expansion of Metal–Organic frameworks with triangular and square building units and the lowest calculated density for porous crystals*, Inorganic chemistry **50**(18), 9147 (2011), doi:[10/cmdwhm](https://doi.org/10/cmdwhm).
- [5] S. Bureekaew, V. Balwani, S. Amirjalayer and R. Schmid, *Isorecticular isomerism in 4,4-connected paddle-wheel metal–organic frameworks: Structural prediction by the reverse topological approach*, CrystEngComm **17**(2), 344 (2015), doi:[10/gnx9cq](https://doi.org/10/gnx9cq).
- [6] J. Kim, B. Chen, T. M. Reineke, H. Li, M. Eddaoudi, D. B. Moler, M. O’Keeffe and O. M. Yaghi, *Assembly of Metal–Organic frameworks from large organic and inorganic secondary building units: New examples and simplifying principles for complex structures*, Journal of the American Chemical Society **123**(34), 8239 (2001), doi:[10/bqxdrj](https://doi.org/10/bqxdrj).
- [7] M. O’Keeffe, M. A. Peskov, S. J. Ramsden and O. M. Yaghi, *The reticular chemistry structure resource (RCSR) database of, and symbols for, crystal nets*, Accounts of chemical research **41**(12), 1782 (2008), doi:[10/fwg963](https://doi.org/10/fwg963).
- [8] M. O’Keeffe and O. M. Yaghi, *Deconstructing the crystal structures of Metal–Organic frameworks and related materials into their underlying nets*, Chemical reviews **112**(2), 675 (2012), doi:[10/fkhc93](https://doi.org/10/fkhc93).
- [9] S. R. Batten, *Topology and interpenetration*, In *Metal-Organic Frameworks: Design and Application*, pp. 91–130. John Wiley & Sons, MacGillivray, L. R. (2010).
- [10] S. Hyde, O. D. Friedrichs, S. Ramsden and V. Robins, *Towards enumeration of crystalline frameworks: The 2D hyperbolic approach*, Solid State Sciences **8**(7), 740 (2006), doi:[10/d8mvj6](https://doi.org/10/d8mvj6).
- [11] J. C. Tan, T. D. Bennett and A. K. Cheetham, *Chemical structure, network topology, and porosity effects on the mechanical properties of Zeolitic Imidazolate Frameworks*, Proceedings of the National Academy of Sciences **107**(22), 9938 (2010), doi:[10/d7pf3n](https://doi.org/10/d7pf3n).
- [12] V. A. Blatov, *Multipurpose crystallochemical analysis with the program package TOPOS*, IUCr CompComm Newsletter **7**(4), 4 (2006).

- [13] O. Delgado-Friedrichs and M. O’Keeffe, *Identification of and symmetry computation for crystal nets*, Acta Crystallographica Section A: Foundations of Crystallography **59**(4), 351 (2003), doi:[10/dbg89q](https://doi.org/10/dbg89q).
- [14] E. V. Alexandrov, A. P. Shevchenko and V. A. Blatov, *Topological databases: Why do we need them for design of coordination polymers?*, Crystal Growth & Design **19**(5), 2604 (2019), doi:[10/gpjds](https://doi.org/10/gpjds).
- [15] S. R. Hall, F. H. Allen and I. D. Brown, *The crystallographic information file (CIF): A new standard archive file for crystallography*, Acta Crystallographica Section A: Foundations of Crystallography **47**(6), 655 (1991), doi:[10/fw2589](https://doi.org/10/fw2589).
- [16] O. Delgado-Friedrichs and M. O’Keeffe, *Crystal nets as graphs: Terminology and definitions*, Journal of Solid State Chemistry **178**(8), 2480 (2005), doi:[10/c8t2sh](https://doi.org/10/c8t2sh).
- [17] M. Treacy, M. Foster and K. Randall, *An efficient method for determining zeolite vertex symbols*, Microporous and Mesoporous Materials **87**(3), 255 (2006), doi:[10.1016/j.micromeso.2005.08.024](https://doi.org/10.1016/j.micromeso.2005.08.024).
- [18] O. Delgado-Friedrichs, S. T. Hyde, M. O’Keeffe and O. M. Yaghi, *Crystal structures as periodic graphs: The topological genome and graph databases*, Structural Chemistry **28**(1), 39 (2017), doi:[10/gnx9cp](https://doi.org/10/gnx9cp).
- [19] O. Delgado-Friedrichs, *Barycentric drawings of periodic graphs*, In *International Symposium on Graph Drawing*, pp. 178–189. Springer (2003).
- [20] J. Bezanson, A. Edelman, S. Karpinski and V. B. Shah, *Julia: A fresh approach to numerical computing*, SIAM review **59**(1), 65 (2017), doi:[10/f9wkpj](https://doi.org/10/f9wkpj).
- [21] J. Bezanson, J. Chen, B. Chung, S. Karpinski, V. B. Shah, J. Vitek and L. Zoubritzky, *Julia: Dynamism and performance reconciled by design*, Proceedings of the ACM on Programming Languages **2**(OOPSLA), 1 (2018), doi:[10/gfgr7m](https://doi.org/10/gfgr7m).
- [22] S. J. Chung, T. Hahn and W. Klee, *Nomenclature and generation of three-periodic nets: The vector method*, Acta Crystallographica Section A: Foundations of Crystallography **40**(1), 42 (1984), doi:[10/b39rn4](https://doi.org/10/b39rn4).
- [23] A. Ching, S. Edunov, M. Kabiljo, D. Logothetis and S. Muthukrishnan, *One trillion edges: Graph processing at Facebook-scale*, Proceedings of the VLDB Endowment **8**(12), 1804 (2015), doi:[10/gfx98n](https://doi.org/10/gfx98n).
- [24] T. Granlund et. al., *The GNU Multiple-Precision (GMP) arithmetic library*.
- [25] W. Eberly, M. Giesbrecht, P. Giorgi, A. Storjohann and G. Villard, *Solving sparse rational linear systems*, In *Proceedings of the 2006 International Symposium on Symbolic and Algebraic Computation*, pp. 63–70, doi:[10/dh2gph](https://doi.org/10/dh2gph) (2006).
- [26] W. Eberly, M. Giesbrecht, P. Giorgi, A. Storjohann and G. Villard, *Faster inversion and other black box matrix computations using efficient block projections*, In *Proceedings of the 2007 International Symposium on Symbolic and Algebraic Computation*, pp. 143–150, doi:[10/bfm7qm](https://doi.org/10/bfm7qm) (2007).
- [27] J. D. Dixon, *Exact solution of linear equations using P-adic expansions*, Numerische Mathematik **40**(1), 137 (1982), doi:[10/fs5tss](https://doi.org/10/fs5tss).

- [28] E. Kaltofen and B. D. Saunders, *On Wiedemann's method of solving sparse linear systems*, In *International Symposium on Applied Algebra, Algebraic Algorithms, and Error-Correcting Codes*, pp. 29–38. Springer, doi:[10/cg8p5p](https://doi.org/10/cg8p5p) (1991).
- [29] O. Delgado-Friedrichs, S. T. Hyde, S.-W. Mun, M. O'Keeffe and D. M. Proserpio, *Nets with collisions (unstable nets) and crystal chemistry*, *Acta Crystallographica Section A: Foundations of Crystallography* **69**(6), 535 (2013), doi:[10/gnx9cr](https://doi.org/10/gnx9cr).
- [30] B. Delaunay, *Neue darstellung der geometrischen kristallographie*, *Zeitschrift für Kristallographie-Crystalline Materials* **84**(1-6), 109 (1933), doi:[10/gnx9qg](https://doi.org/10/gnx9qg).
- [31] P. Niggli, W. Wien and F. Harms, *Handbuch Der Experimentalphysik*, Akademische Verlagsgesellschaft, Leipzig (1928).
- [32] A. Santoro and A. Mighell, *Determination of reduced cells*, *Acta Crystallographica Section A: Crystal Physics, Diffraction, Theoretical and General Crystallography* **26**(1), 124 (1970), doi:[10/fmzmtq](https://doi.org/10/fmzmtq).
- [33] A. Togo and I. Tanaka, *Spglib: A software library for crystal symmetry search*, arXiv preprint arXiv:1808.01590 (2018), [1808.01590](https://arxiv.org/abs/1808.01590).
- [34] R. W. Grosse-Kunstleve, N. K. Sauter, N. W. Moriarty and P. D. Adams, *The Computational Crystallography Toolbox: Crystallographic algorithms in a reusable software framework*, *Journal of Applied Crystallography* **35**(1), 126 (2002), doi:[10/ddpwq8](https://doi.org/10/ddpwq8).
- [35] G. Havas, B. S. Majewski and K. R. Matthews, *Extended gcd and Hermite normal form algorithms via lattice basis reduction*, *Experimental Mathematics* **7**(2), 125 (1998), doi:[10/gnx9cs](https://doi.org/10/gnx9cs).
- [36] E. Ren, P. Guilbaud and F.-X. Coudert, *High-throughput computational screening of nanoporous materials in targeted applications* (2022), [2202.09886](https://arxiv.org/abs/2202.09886).
- [37] N. Castel and F.-X. Coudert, *Atomistic models of amorphous metal-organic frameworks* (2022), [2202.07059](https://arxiv.org/abs/2202.07059).
- [38] Y. G. Chung, E. Haldoupis, B. J. Bucior, M. Haranczyk, S. Lee, H. Zhang, K. D. Vogiatzis, M. Milisavljevic, S. Ling, J. S. Camp, B. Slater, J. I. Siepmann *et al.*, *Advances, Updates, and Analytics for the Computation-Ready, Experimental Metal–Organic Framework Database: CoRE MOF 2019*, *Journal of Chemical & Engineering Data* **64**(12), 5985 (2019), doi:[10/gpjdtr](https://doi.org/10/gpjdtr).
- [39] C. Bonneau, M. O'Keeffe, D. M. Proserpio, V. A. Blatov, S. R. Batten, S. A. Bourne, M. S. Lah, J.-G. Eon, S. T. Hyde and S. B. Wiggin, *Deconstruction of crystalline networks into underlying nets: Relevance for terminology guidelines and crystallographic databases*, *Crystal Growth & Design* **18**(6), 3411 (2018), doi:[10/gdsxqk](https://doi.org/10/gdsxqk).
- [40] A. P. Shevchenko and V. A. Blatov, *Simplify to understand: How to elucidate crystal structures?*, *Structural Chemistry* **32**(2), 507 (2021), doi:[10/gpjt7b](https://doi.org/10/gpjt7b).
- [41] E. Alexandrov, V. Blatov, A. Kochetkov and D. Proserpio, *Underlying nets in three-periodic coordination polymers: Topology, taxonomy and prediction from a computer-aided analysis of the Cambridge Structural Database*, *CrystEngComm* **13**(12), 3947 (2011), doi:[10/fr67j8](https://doi.org/10/fr67j8).

- [42] M. Li, D. Li, M. O’Keeffe and O. M. Yaghi, *Topological analysis of metal–organic frameworks with polytopic linkers and/or multiple building units and the minimal transitivity principle*, Chemical Reviews **114**(2), 1343 (2014), doi:[10/gnb5gk](https://doi.org/10/gnb5gk).
- [43] B. J. Bucior, A. S. Rosen, M. Haranczyk, Z. Yao, M. E. Ziebel, O. K. Farha, J. T. Hupp, J. I. Siepmann, A. Aspuru-Guzik and R. Q. Snurr, *Identification schemes for metal–organic frameworks to enable rapid search and cheminformatics analysis*, Crystal Growth & Design **19**(11), 6682 (2019), doi:[10/gf9g2q](https://doi.org/10/gf9g2q).
- [44] F. M. A. Noa, M. Abrahamsson, E. Ahlberg, O. Cheung, C. R. Göb, C. J. McKenzie and L. Öhrström, *A unified topology approach to dot-, rod-, and sheet-MOFs*, Chem **7**(9), 2491 (2021), doi:[10/gmwkzz](https://doi.org/10/gmwkzz).
- [45] C. Zheng, Y. Li and J. Yu, *Database of open-framework aluminophosphate structures*, Scientific Data **7**(1) (2020), doi:[10/gpjdtq](https://doi.org/10/gpjdtq).
- [46] J. V. Smith, *Topochemistry of zeolites and related materials. 1. Topology and geometry*, Chemical Reviews **88**(1), 149 (1988), doi:[10/b4qp7h](https://doi.org/10/b4qp7h).
- [47] M. Treacy, I. Rivin, E. Balkovsky, K. Randall and M. Foster, *Enumeration of periodic tetrahedral frameworks. II. Polynodal graphs*, Microporous and Mesoporous Materials **74**(1-3), 121 (2004), doi:[10/d4rpgc](https://doi.org/10/d4rpgc).
- [48] R. Pophale, P. A. Cheeseman and M. W. Deem, *A database of new zeolite-like materials*, Physical Chemistry Chemical Physics **13**(27), 12407 (2011), doi:[10/d3sjr2](https://doi.org/10/d3sjr2).
- [49] S. Gražulis, D. Chateigner, R. T. Downs, A. F. T. Yokochi, M. Quirós, L. Lutterotti, E. Manakova, J. Butkus, P. Moeck and A. Le Bail, *Crystallography Open Database – an open-access collection of crystal structures*, Journal of Applied Crystallography **42**(4), 726 (2009), doi:[10.1107/S0021889809016690](https://doi.org/10.1107/S0021889809016690).
- [50] C. R. Groom, I. J. Bruno, M. P. Lightfoot and S. C. Ward, *The Cambridge structural database*, Acta Crystallographica Section B: Structural Science, Crystal Engineering and Materials **72**(2), 171 (2016), doi:[10.1107/S2052520616003954](https://doi.org/10.1107/S2052520616003954).