

# Predictive crystallography at scale: mapping, validating, and learning from 1000 crystal energy landscapes†‡

Christopher R. Taylor,  Patrick W. V. Butler  and Graeme M. Day \*

Received 17th May 2024, Accepted 22nd May 2024

DOI: 10.1039/d4fd00105b

Computational crystal structure prediction (CSP) is an increasingly powerful technique in materials discovery, due to its ability to reveal trends and permit insight across the possibility space of crystal structures of a candidate molecule, beyond simply the observed structure(s). In this work, we demonstrate the reliability and scalability of CSP methods for small, rigid organic molecules by performing in-depth CSP investigations for over 1000 such compounds, the largest survey of its kind to-date. We show that this highly-efficient force-field-based CSP approach is superbly predictive, locating 99.4% of observed experimental structures, and ranking a large majority of these (74%) as among the most stable possible structures (to within uncertainty due to thermal effects). We present two examples of insights such large predicted datasets can permit, examining the space group preferences of organic molecular crystals and rationalising empirical rules concerning the spontaneous resolution of chiral molecules. Finally, we exploit this large and diverse dataset for developing transferable machine-learned energy potentials for the organic solid state, training a neural network lattice energy correction to force field energies that offers substantial improvements to the already impressive energy rankings, and a MACE equivariant message-passing neural network for crystal structure re-optimisation. We conclude that the excellent performance and reliability of the CSP workflow enables the creation of very large datasets of broad utility and explanatory power in materials design.

## 1 Introduction

The discovery of new materials is important for addressing many critical societal needs, including energy production and storage, pollution remediation and healthcare. Research endeavours aimed at improving the success and efficiency of functional materials discovery, based on traditional efforts developing our

*School of Chemistry, University of Southampton, Southampton, SO17 1BJ, UK. E-mail: g.m.day@soton.ac.uk*

† Electronic supplementary information (ESI) available: Computational details of CSP and ML calculations, summary of molecule set including CSD refcodes, and summary of CSP matches to experiment. See DOI: <https://doi.org/10.1039/d4fd00105b>

‡ Complete CSP landscapes and ML models are available from <https://doi.org/10.5258/SOTON/D3094>.



understanding of the rules of crystal packing and, more recently, applications of machine learning, have benefited greatly from the availability of databases of stable crystal structures.

A major resource for such efforts has been the growth of accessible, curated databases of crystal structures and (some of) their properties. The most general and widely-used include the Cambridge Structural Database (CSD)<sup>1</sup> for organic and organometallic systems, and the Inorganic Crystal Structure Database.<sup>2</sup> These resources are unparalleled in their volume of experimental crystal structure information, but do not currently offer information about calculated or hypothetical structures. Historically, such data was limited to specialised areas (such as the Atlas of Prospective Zeolite Structures<sup>3</sup>), though recent developments have made remarkable progress in generalising this concept, including extensions to the Crystallography Open Database (COD)<sup>4–7</sup> and the Materials Project.<sup>8</sup> In the field of organic molecular crystals, however, much of our understanding of the rules of crystal packing derive from databases of experimentally observed structures.

Modern computational chemistry, employing both molecular and solid-state simulation techniques, can add significantly to the information that is available from experimentally determined crystal structures, and identify previously unverifiable trends. Computational studies of polymorphism<sup>9–13</sup> have evaluated the typical lattice energy differences between crystalline polymorphs of organic molecules, studies of conformations of flexible molecules<sup>11,14</sup> in their observed crystal structures have improved understanding of the limits of molecular strain in stable structures, and studies of the thermodynamics of co-crystallisation<sup>15–18</sup> have aided in rationalising a complex phenomenon with ramifications for experimental design.

A more complete and, crucially, predictive view of organic crystal packing can be obtained from crystal structure prediction (CSP).<sup>19,20</sup> A core concept in CSP is the crystal energy landscape (or CSP landscape) – the set of plausible crystal packings for a chemical species (or combination thereof), representing an exploration of the crystalline configuration space to identify candidate structures predicted to be stable (which ideally includes any observed structures), ranked in terms of their energetic stability. These CSP landscapes are of great utility in understanding and rationalising the thermodynamic and kinetic behaviour of crystal systems; multiple low-energy minima that are close-lying on a CSP landscape may be indicative of a significant risk of polymorphism, while dynamical simulations exploring these landscapes give insight into the kinetic trapping of metastable forms<sup>21</sup> and the observed absence of other predicted forms.<sup>22,23</sup>

Moreover, the energy landscape forms the foundation of the energy–structure–function maps that have in recent years demonstrated great power in materials discovery.<sup>24–27</sup> By associating computed properties (*e.g.* gas uptake, charge carrier mobility) with hypothetical crystal structures from the CSP landscape, it becomes possible to predict whether a molecule is a promising candidate for creating new functional materials, *i.e.* if it has one (or more) favourably-ranked crystal structures which are predicted to achieve the desired property.

The techniques and challenges in the field of organic CSP have been reviewed;<sup>19,20</sup> we provide a brief overview for the sake of context. CSP is typically considered a combination of two broad challenges: efficient and thorough



sampling of the configuration space of crystal packing, and accurate, cost-effective structural optimisation, ranking, and property calculation.

The sampling of hypothetical crystal packing arrangements is made extremely difficult due to the “curse of dimensionality”—the number of independent degrees of freedom to sample creates a vast configurational space. As a result, simple grid-based sampling approaches must be eschewed in favour of more sophisticated techniques, such as low-discrepancy quasi-random sampling<sup>28,29</sup> and genetic algorithm approaches.<sup>30</sup>

Meanwhile, the optimisation and ranking methods must be accurate enough to describe the fine balance of different intermolecular interactions (electrostatics, dispersion, hydrogen-bonding, *etc.*), resolving lattice energy differences often smaller than a  $\text{kJ mol}^{-1}$ , while sufficiently cost-effective to be applied to very large numbers ( $\gg 10^5$ ) of trial crystal structures. Historically, this has entailed the use of simple empirical force fields, but modern developments often employ tailor-made force fields<sup>18,31</sup> or machine-learned potentials derived from *ab initio* calculations. Still, empirical force fields retain their power even today due to their efficiency and broad transferability, often being the initial step of a hierarchy of increasingly accurate (and expensive) energy models employed in one CSP workflow.

Despite these two broad and ongoing challenges, organic CSP has demonstrated enormous success in diverse applications, crucially proving itself to be a truly predictive technique, guiding synthesis and discovery of novel forms of porous materials,<sup>24</sup> highly-flexible pharmaceutical molecules,<sup>31</sup> co-crystals,<sup>32</sup> simple molecules previously thought to be monomorphic<sup>33</sup> and templating of predicted metastable polymorphs.<sup>34</sup> A recurrent landmark in the field is the series of Blind Tests of CSP, showcasing the diversity of methods (and success rates thereof) employed within different CSP techniques to predict experimental crystal structures without any knowledge beyond the molecular chemical diagram.<sup>35</sup> Recent iterations of the Blind Test have demonstrated that CSP method development is successfully keeping pace with the complexity of molecules and crystal structures specifically selected to stress-test it.

Some of the most consequential recent developments in CSP have employed machine learning (ML) techniques in one form or another.<sup>36</sup> Among the most intuitive applications is the use of ML to learn relationships between the structure and lattice energy, either by learning the difference (*i.e.*  $\Delta$ -ML) between lattice energies computed at a lower level of theory to those at a higher level,<sup>37–40</sup> or to learn the relationship between the lattice energy and the ML descriptors directly through the training of ML potentials.<sup>41,42</sup> These approaches have achieved high accuracy predictions at a fraction of the cost of the full periodic density functional theory (DFT) reference calculations – particularly significant given the latter's ubiquity in recent Blind Test entries.

That said, ML techniques have further applications in CSP beyond improved optimisations and energy rankings. In particular, ML and related approaches applied to databases of experimental crystal structures and properties have demonstrated success in predicting NMR chemical shifts<sup>43</sup> and in formulating models of molecular hydrogen bond propensities within crystals.<sup>44</sup> Recent work by Cersonsky<sup>45</sup> has demonstrated machine-learning of the relationship between crystal lattice energies and the relative contributions to these by different chemical functional groups, paving the way for data-driven insights into new crystal engineering techniques. ML has also been shown to have potential to



enhance the analysis of crystal energy landscapes, by identifying structure–function relationships that might evade simple inspection but nevertheless offer explanatory and predictive power.<sup>46</sup>

Our aims in this work are threefold. Firstly, we seek to demonstrate the capability of our method for efficient, large-scale rigid-molecule organic CSP by presenting the results of the largest to-date CSP study, applying the methods to over 1000 molecules with observed crystal structures in the CSD. Secondly, we in turn assess the quality and reliability of our method by evaluating how often the experimental crystal structures of these molecules are reproduced in our CSPs, and how well they are ranked energetically compared to hypothetical structures on their CSP landscapes. Finally, we demonstrate example applications of this dataset, including assessing the distribution of CSP-derived space group preferences as a function of predicted lattice energy, spontaneous resolution of chiral molecules and the training of transferable machine-learned energetic models from a very large set of CSP landscapes.

## 2 Computational methods

### 2.1 Molecule selection

We used the CSD's ConQuest software and Python API<sup>1</sup> to search the CSD for crystal structures of rigid molecules to which our CSP methods could be applied on a very large scale. Restricting our search to solved crystal structures (*i.e.* with coordinates for all atoms in the asymmetric unit) of single chemical species (no co-crystals, solvates, inclusion compounds), we additionally filtered structures based on the following criteria: containing only elements from C, H, N, O, F;  $Z' \leq 1$ ; molecular weight less than 230; and importantly containing no rotatable bonds (as defined by the CSD's internal criteria).

### 2.2 Molecular geometry optimisation

For each molecule, we began by extracting its in-crystal conformation from the corresponding CSD entry (where there is more than one entry, we select only the first listed in the database). This molecular conformation was optimised in DFT (as implemented in Gaussian<sup>47</sup>), using the B3LYP<sup>48,49</sup> exchange–correlation functional, a Pople-type 6-311G\*\* basis,<sup>50</sup> and Grimme's D3 dispersion correction<sup>51</sup> with Becke–Johnson damping.<sup>52</sup>

Distributed multipole analysis (DMA,<sup>53,54</sup> using the GDMA package) was performed on the resulting molecular conformers' charge densities to obtain atom-centred multipoles, as part of the model potential applied during lattice energy minimisation of trial crystal structures; multipoles up to hexadecapole were calculated for all atoms. The MUFFIT<sup>55,56</sup> software was used to fit atomic point charges for each molecule to best describe the molecular electrostatic potential. The resulting molecular conformers and their sets of multipoles are then used as the inputs to CSP, with each unique molecule represented by a single conformer and corresponding set of atomic multipoles and charges.

### 2.3 Crystal structure generation and optimisation

CSP was performed using the Global Lattice Energy Explorer (GLEE) package,<sup>28</sup> which uses quasi-random sampling of crystal packing variables to generate trial



crystal structures uniformly distributed across the lattice energy landscape, followed by rigid-molecule lattice energy minimisation using an anisotropic atom–atom intermolecular force field. All resulting local energy minima are treated as possible crystal structures of the molecule.

Space group symmetry is used to reduce the dimensionality of the search space, so that only the position and orientation of molecules in the asymmetric unit are sampled, with all other molecules in the unit cell generated by symmetry. In this study, we restrict ourselves to generating crystal structures with one independent molecule in the asymmetric unit ( $Z' = 1$ ). We sampled the 26 most commonly observed space groups for organic molecular crystals (listed in ESI†); these space groups cover over 99.4% of  $Z' \leq 1$  structures in the CSD. These space groups were sampled equally, irrespective of their observed frequency in the CSD: quasi-random structures are generated and lattice energy minimised until 10 000 successfully energy minimised crystal structures were found in each space group (260 000 structures per molecule). The CSP process is highly parallelisable, as each crystal structure is independent.

The trial crystal that passed geometric checks were lattice energy-minimised in three stages. Non-electrostatic interactions (principally intermolecular dispersion and exchange-repulsion) were described by the FIT  $exp - 6$  force field,<sup>57,58</sup> supplemented by fluorine parameters from Williams and Houpt.<sup>59</sup> At the final stage of optimisation, performed using DMACRYS,<sup>58</sup> the FIT potential was applied with atomic multipole electrostatics. Full details are provided in ESI.†

It is commonplace that multiple unique initial configurations optimise to the same local energy minimum. We remove these duplicates by fast comparison of simulated powder X-ray diffraction patterns, followed by structural comparisons using the COMPACK<sup>60</sup> algorithm as implemented in the CSD Python API.

**2.3.1 Locating experimental structures on the landscapes.** To assess our CSP workflow's performance, comparison of the known experimental crystal structures to the sets of predictions was automated using the COMPACK algorithm between the experimental crystal structures of these molecules and every unique crystal structure of that molecule in the CSP set.

## 2.4 Machine learned interatomic potentials

To investigate the potential of the CSP dataset to train data-derived models, a subset of the predicted crystal structures with energies within  $8.0 \text{ kJ mol}^{-1}$  of the global energy minimum on their CSP landscape, were selected for training a lattice energy correction to the FIT+DMA force field. The subset was determined by active learning *via* query-by-committee using a committee of 8 high-dimensional neural network potentials (NNPs), with selected structures evaluated by DFT+D single points at the B86bPBE+XDM level. From this, a dataset of the crystal structures and the corresponding energy correction between the FIT+DMA and B86bPBE+XDM lattice energies ( $\Delta E$ ) was created. B86bPBE+XDM lattice energies were calculated as the total energy less the energies of the isolated molecules from the unit cell calculated with the same basis set and tolerances.

The initial dataset (before the active learning iterations) was generated by randomly selecting up to 10 low energy predicted crystal structures for each compound, resulting in a dataset of 10 249 structures approximately evenly distributed across the compounds. To evaluate transferability, the total dataset



was partitioned into a training dataset consisting of the CSP structures for a randomly selected *ca.* 85% of the compounds and an extrapolation test set consisting of CSP structures for the remaining compounds. A further in-domain test set was formed by randomly extracting one structure per compound from the training set. The NNPs were then trained on the remaining training dataset to yield a  $\Delta$ -ML model capable of predicting the lattice energy correction. The standard deviation between the ensemble of NNPs was used to estimate the uncertainty of predictions, and was exploited in the active learning iterations to add high uncertainty candidates from the remaining low energy predicted structures of the training compounds, which overall added a further 1000 structures to the training set. Corrected CSP landscapes were calculated using the final model by adding predicted lattice energy corrections to the FIT+DMA energies (FIT+DMA+ $\Delta$ -ML). Additionally, for performing unconstrained geometry optimisations MACE equivariant message-passing neural network (MPNN) models were trained using a dataset derived from that of the NNP correction model. Further details of the datasets and machine learning models are provided in the ESI.<sup>†</sup>

## 3 Results

### 3.1 Diversity of survey set

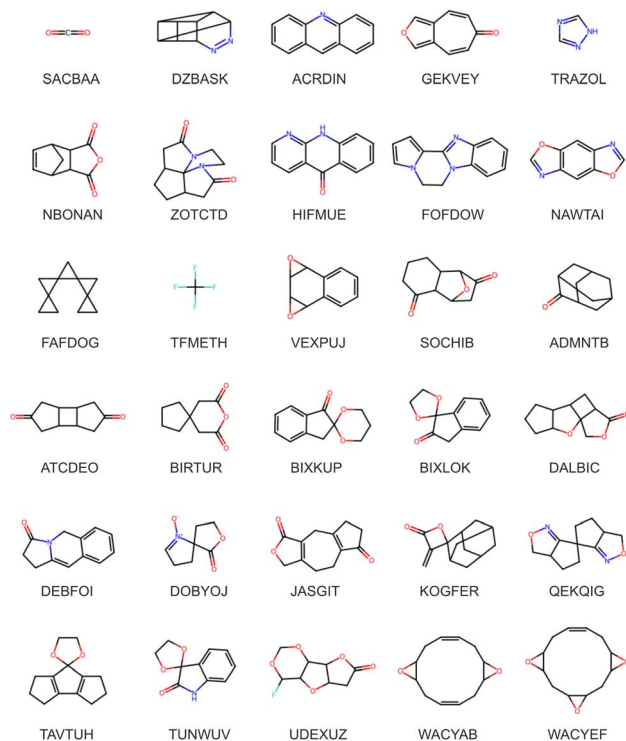
Our aforementioned search criteria yielded 1007 distinct molecules crystallising in 1040 crystal structures observed in the CSD. The constraint of no rotatable bonds necessarily limits chemical diversity, but despite this, our candidate set still displays a variety of chemical functionalities and molecular structures, as seen in Fig. 1 and Table 1. A complete list of molecules, including formulae, CSD refcode identifiers, SMILES strings, and systematic and common names is available in the ESI.<sup>†</sup>

**3.1.1 Deviation between experimental and gas-phase conformations.** Despite restricting our set to molecules containing no rotatable bonds, this does not preclude molecular flexibility entirely. More complex collective motions cannot be described in terms of a single torsional angle about a covalent bond, and so molecules displaying such conformational flexibility are present in this candidate set – the prototypical example is a ring “flip” or buckling, such as the boat-chair interconversion of cyclohexane rings.

As a measure of the typical deviation in molecular geometry between the observed crystal structures and the DFT optimised molecules used in CSP, we present in Fig. 2 the histogram of all-atom root-mean-squared deviations (RMSD) in atomic positions between the gas-phase conformers used in CSP and their in-crystal initial conformations as extracted from the CSD.

As might be expected from such rigid molecules, the average RMSD after gas-phase optimisation is very small – approximately 0.11 Å, which corresponds to *e.g.* adjusting the C–H bond lengths in fluorobenzene by 0.16 Å. While the distribution is skewed towards small conformational changes, the outliers with larger RMSD values demonstrate the limitations of defining molecular flexibility in terms of rotatable bonds alone. The largest RMSD values correspond to systems where changes in ring conformation cause large overall molecular changes – the largest observed RMSD of 0.60 Å occurs in 7-oxa-1-azaspiro(4.4)non-1-en-6-one 1-oxide (CSD reference code: DOBYOJ), in which a 5-membered saturated ring can twist about a spiro carbon centre. The highest molecular RMSD values are largely





**Fig. 1** Molecular diagrams and crystal structure CSD reference codes for (top three rows) a random selection of the 1007 molecules included in the large-scale CSP study. The bottom three rows show molecules in the set with the largest differences between in-crystal and optimised molecular geometries (as measured by all-atom RMSD); the CSP landscapes for these molecules were re-optimised using the transferrable MACE model (final section).

**Table 1** Total counts of selected functional groups across the full set of molecules, as assessed by RDKit<sup>61</sup> from SMILES strings

Functional group	Count	Functional group	Count	Functional group	Count
Benzene	418	Ether	572	Ester	134
Pyridine	76	Amide	297	Carbonyl	702
Imidazole	35	Imide	43	Imine	13
Furan	22	Lactone	91	Secondary amine	405
Piperidine	51	Epoxide	73	Tertiary amine	537
Urea	22	Ketone	269	Halogen (F)	240

associated with molecules containing spiro carbon atoms (denoted by the red portion of the bars in Fig. 2).

Regardless, the overall inflexibility displayed by this set serves as a strong indicator that our assumption of a single, near-crystalline conformer of each molecule is reasonable, and unlikely to be a systematic source of error in predicting many known structures of molecules in this set.





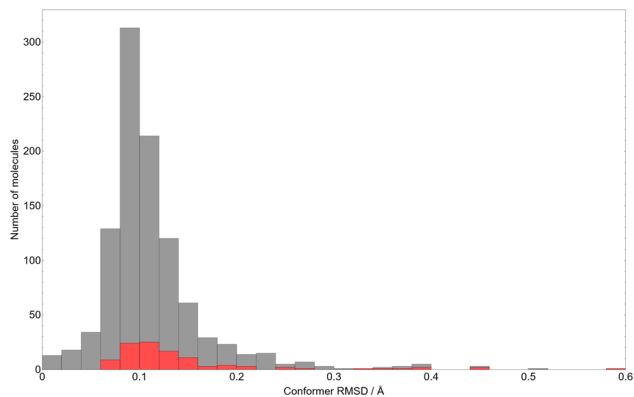


Fig. 2 A histogram of the RMSD between all atomic positions in the gas-phase optimised CSP candidate molecules, relative to their initial in-crystal conformations. The red portions of each bar indicate the molecules in that bin with spiro carbon centres.

**3.1.2 Crystalline diversity.** Fig. 3 shows a comparison of the distribution of crystallographic space groups for most small molecules in the CSD to the rigid-molecule subset selected for our CSP survey. There is no significant difference in the relative occurrence of different crystalline symmetries between the two sets, though the relative ordering of space groups by frequency varies slightly. For example, space group 15 ( $C2/c$ ) is observed slightly more frequently in our rigid set and space group 2 ( $P\bar{1}$ ) slightly less so. Regardless, our rigid molecule set is reasonably representative of the range of crystal packing symmetries observed in the CSD.

In contrast, our rigid molecule set displays a diminished frequency of hydrogen bonding (H-bonding) in the observed crystal structures compared to the

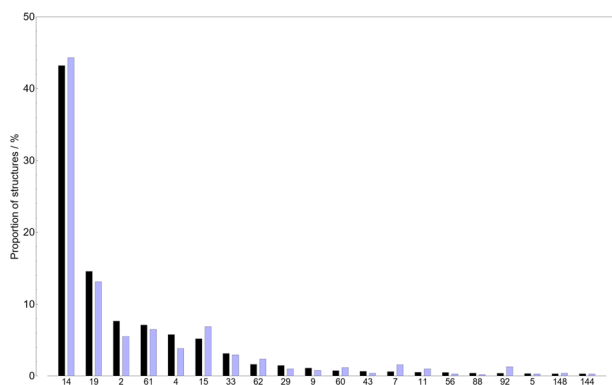


Fig. 3 The relative frequency of space groups observed for crystal structures (where  $Z' \leq 2$ ) of molecules in the CSD with a molecular weight under 230 (black) and the subset (blue) of these molecules with no rotatable bonds selected for our CSP surveys. The distributions are presented for only the 20 most common space groups (of the general molecule case) for clarity.





CSD more generally, as might be expected by excluding rotatable bonds. 26.7% (277) of crystal structures of molecules in our subset contain at least one intermolecular H-bond. This compares to 62.3% of structures of molecules of similar size without rotatable bond restrictions, demonstrating the bias introduced by the omission of common H-bonding groups due to their flexibility, including alcohols ( $-\text{OH}$ ), carboxylic acids ( $-\text{COOH}$ ), and primary amines ( $-\text{NH}_2$ ). However, the proportion of H-bonded systems is still significant enough that meaningful H-bond chemistry is incorporated in our survey, albeit underrepresented. Consequently, our set in turn over-represents chemistry such as  $\pi$ - $\pi$  stacking and “weak” (*i.e.* more isotropic, less localised) interactions, and assessments of our CSP energy model's performance must be made with these biases in mind.

### 3.2 Quality assessment of the CSP results

We propose that the dataset of predicted crystal structure landscapes across a large, diverse set of molecules is valuable for the development of future predictive models. To evaluate the quality of the dataset, we assess three aspects: completeness of the landscapes; how well, geometrically, the CSP calculations reproduce the known crystal structures of these molecules; and the quality of the relative energies of the predicted structures.

The 26 space groups included in our standard search include those with an observed frequency above 0.05% in the CSD. Naturally, in a very large survey of molecules, we include some that crystallise in less frequent space groups; for these 8 cases, we added the space group of the observed crystal structure to the search. In 3 additional cases, the symmetry of the observed structure means that it could not have been sampled without performing CSP with multiple independent molecules ( $Z' > 1$ ); in these cases, we include the datasets in our study, knowing that the observed crystal structures could not have been located. Of those that could have been located, the searches find matches for 1034 of the 1040 observed crystal structures: reasons for the 6 missed matches are discussed below.

A main assumption in CSP is that the observed crystal structures of a molecule correspond to the lowest energy possible structures. We use this assumption to assess the quality of the calculated energies. Fig. 5 summarises the energetic ranking of experimentally observed crystal structures within the CSP landscapes *via* the distribution of  $\Delta E$ , the energy difference between the prediction matching the experimental crystal structure and the global lattice energy minimum crystal structure on their parent molecule's CSP landscape; an example landscape is shown in Fig. 4. In the case of molecules with chiral centers and only one stereoisomer present,  $\Delta E$  is calculated only among Sohncke space groups (those whose symmetry elements contain only translations, rotations and roto translations).

Of the 1034 experimentally determined crystal structures where a matching structure was identified on the CSP landscape, 424 (41%) correspond to the global lattice energy minimum from CSP. Those that do not correspond to global energy minima could be due to inaccuracy of the model potential (FIT+DMA), neglect of other contributions to the lattice free energy or where the kinetics of crystallisation favour a metastable structure. Nyman and Day found that, for crystals of rigid molecules, lattice vibrational contributions to room temperature free energy



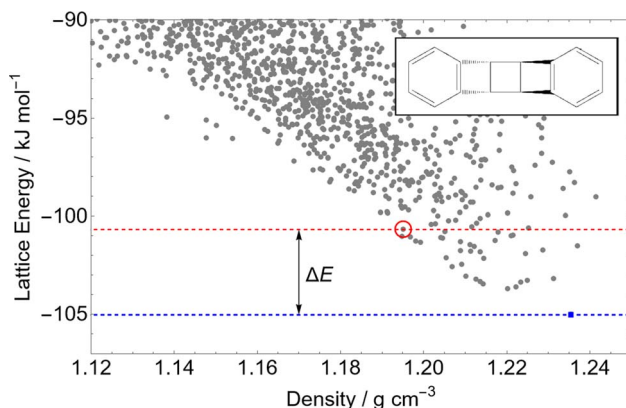


Fig. 4 An example CSP landscape for a single molecule 3,4,7,8-dibenzo-tricyclo(4.2.0.0<sup>2.5</sup>)octa-3,7-diene – the set of predicted lattice energy minima from our CSP workflow. Each point represents a local predicted minimum, and thus stable hypothetical crystal packing. The blue square point is the global energy minimum, which in the absence of experimental information is taken to be the most likely crystal structure. If the red circled point is a match to the experimentally-observed crystal structure,  $\Delta E$  is the difference in energy between it and the global minimum, the energy rank used in this work as a measure of the quality of the calculated energies.

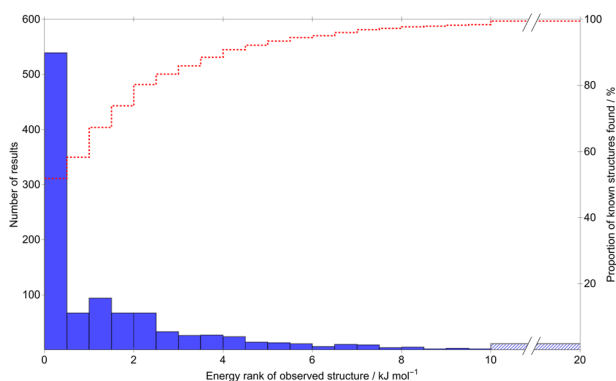


Fig. 5 A histogram of the frequency (blue bars) with which our CSP workflow achieves a match to the experimentally-known structures of our molecule set, grouped by the relative energy of that match compared to the CSP global minimum (0.5 kJ mol<sup>-1</sup> bins). The red dashed line relates to the secondary (right) y-axis, the proportion of known structures located successfully as a function of the relative energy at that bin and below. Note the broken x-axis; the highest-energy bin (blue hatching) encompasses all matches with relative energy greater than 10 kJ mol<sup>-1</sup>.

differences between polymorphs rarely exceed 2 kJ mol<sup>-1</sup>;<sup>9,10</sup> 767 (74%) of the observed crystal structures are found within 2 kJ mol<sup>-1</sup> of the global energy minimum – the estimated error introduced by neglecting lattice vibrations and thermal expansion in the CSP calculations.

Furthermore, the known structure(s) almost always lie within 8 kJ mol<sup>-1</sup> of the global lattice energy minimum (1011, 97.8%, of the observed crystal structures),



consistent with the observation<sup>10</sup> that known polymorphic pairs of small, rigid molecules are rarely separated by more than this energy difference. Thus, the energy model used here, combining empirically parameterised repulsion–dispersion with atomic multipole electrostatics, provides energy ranking of crystal structures that is consistent with observed polymorphism, within the limits of temperature-free lattice energy-based predictions.

We quantify the geometric quality of the predictions using an all-atom RMSD within 30-molecule clusters (RMSD<sub>30</sub>) from experimentally-determined crystal structures and their corresponding match within the CSP sets. A histogram of RMSD<sub>30</sub> (Fig. 6) shows that geometric agreement is generally very good: RMSD<sub>30</sub> is

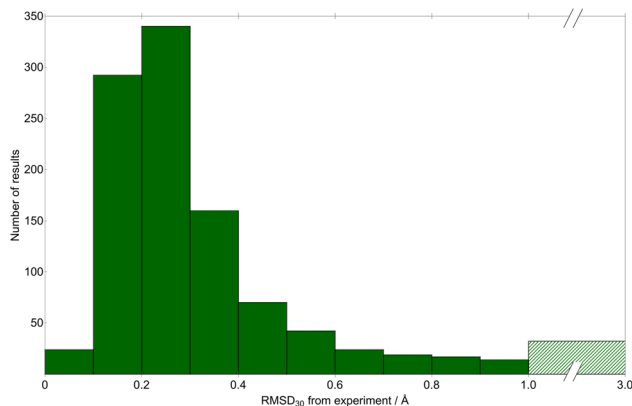


Fig. 6 A histogram of the geometric deviation between experimentally-determined crystal structures and the corresponding matching structures from CSP. The deviation is measured as the RMSD in atomic positions within 30-molecule clusters from experimental and CSP structures. Note the broken x-axis – the largest-deviation bin (green hatching) includes all matches with RMSD<sub>30</sub> greater than 1.0 Å.

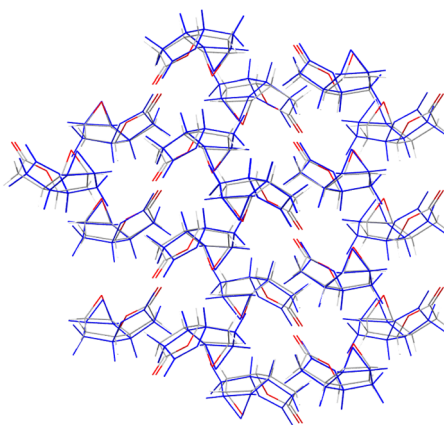


Fig. 7 Overlay of 30-molecule clusters from the X-ray determined crystal structure (atoms coloured by element, CSD reference code SIBJIX) with the matching prediction (blue) – the global energy minimum structure for (1aR,2aS,5aS,5bS)-perhydro-4*H*-oxir-eno(3,4)cyclopenta(1,2-*b*)furan-4-one. For this match, RMSD<sub>30</sub> = 0.393 Å.



below 0.4 Å in 78.9% (816) of matches. As a visualisation of this level of agreement, Fig. 7 shows an overlay of the X-ray determined crystal structure of (1a*R*,2a*S*,5a*S*,5b*S*)-perhydro-4*H*-oxireno(3,4)cyclopenta(1,2-*b*)furan-4-one<sup>62</sup> (CSD reference code SIBJIX) and the predicted global lattice energy minimum, with  $\text{RMSD}_{30} = 0.393$  Å. As a reference for these values, consider the  $\text{RMSD}_{30}$  between structural determinations of the same crystal structure at ambient and low temperature:  $\text{RMSD}_{30} = 0.204$  Å between neutron diffraction crystal structures of naphthalene at 5 K and 295 K, and  $\text{RMSD}_{30} = 0.160$  Å between 20 K and 330 K crystal structures of form I paracetamol.<sup>63</sup> 327 matches have  $\text{RMSD}_{30}$  below 0.204 Å, *i.e.* have geometric deviations that are of a magnitude that can be explained by the temperature-free nature of structural optimisations used in CSP. While known crystal structures are reproduced very well by CSP in most cases, there are a small number where agreement is less satisfactory: in 32 cases (3% of structures),  $\text{RMSD}_{30} > 1$  Å. Despite what might be assumed, we find no significant correlation between the  $\text{RMSD}_{30}$  of the experimental match and the molecular RMSD of the parent gas-phase conformer used for CSP. Assuming an experimental match is found, the geometric quality of the match is only weakly sensitive to the difference in the molecular conformation used for CSP; even the most extreme outlier in molecular conformational change (the aforementioned DOBYOJ) achieves a reasonable geometric match to the experimental crystal structure, with  $\text{RMSD}_{30} = 0.654$  Å.

It is evident that an overwhelming proportion of such rigid molecules can successfully be treated with the CSP workflow implemented here. Our sampling procedure followed by a cost-effective, approximate minimisation method successfully locates the vast majority of observed crystal structures with excellent geometric agreement with the experimental structure, and routinely ranks them as among the most stable structures on the landscape.

**3.2.1 Analysis of individual outliers.** The cases where no match whatsoever is found to the experimental structure are limited – approximately 0.5% of the experimental structures considered.

One missed match was for the molecule 4,8*b*-dihydropyrrolo[3,4-*b*]indole-1,3(2*H*,3a*H*)-dione (CSD reference code BUVGAC), which displays a large deviation in molecular geometry between the observed crystal and the gas phase optimised geometry used in CSP. Flexibility in the fused ring system allows the molecule to “fold up” in the gas-phase optimisation, sufficient to change its packing behaviour. As a result, the experimental crystal structure is not a minimum on a CSP landscape derived from the gas-phase conformer.

In the case of DNNEPH10 (1,8-dinitroso-naphthalene), we fail to find the experimental structure despite an apparently rigid, planar molecule that changes very little in the gas-phase optimisation. However, even optimising the known crystal structure using our FIT+DMA energy model results in a final structure that does not match, *i.e.* the experimental structure appears to be unstable at our level of theory. This may indicate a failure of our energy model for a case of somewhat unusual chemistry.

Two of the missed matches are unusual cases in which a refcode “stem” (the initial six letters, typically shared by a “family” of multiple CSD entries of the same species) is used by crystal structures containing distinct chemical bonding arrangements. Our procedure takes a single representative of a CSD refcode family as the source for the molecular connectivity, which was assumed to stay



fixed. In the case of refcodes XUGHUD/XUHGUD01, there is a tautomeric difference, which led to no match with XUHGUD01. For IHEPUG/IHEPG02, the subject molecule is a diastereomeric fused ring system, which exists in an *anti* configuration in IHEPUG, but a *syn* configuration in IHEPG02. It is arguable that these crystal structures should not be considered part of the same “family”, as these distinct bonding arrangements are not interconvertible. CSP was performed with the isomer found in IHEPUG, so no crystal structure matching IHEPG02 was located. These two systems demonstrate a shortcoming of our approach, in that we assumed that a given CSD refcode “stem” always denotes the same molecular connectivity, including protonation states. Fortunately, these are the only instances of this assumption failing in the entire set.

Our final missed match occurs for QIBCEK, benzo(*f*)phthalazin-4(3*H*)-one, another rigid planar molecule. However, upon inspection, we posit that there are flaws in the experimental determination of this structure as held in the CSD – there are extremely close hydrogen contacts ( $<1.4$  Å) and an unsatisfied potential hydrogen-bonding arrangement despite 1 : 1 donor–acceptor availability. While we retain it as a missed match for the purposes of conservatively assessing our CSP method’s performance, we also emphasise that such a large-scale, unbiased workflow has potentially identified an incorrect experimental structure simply through its absence from the CSP landscape.

### 3.3 Revisiting empirical rules

Large databases of experimentally determined crystal structures have been analysed to uncover trends in the packing preferences of organic molecules. The availability of high quality crystal energy landscapes should allow organic solid state researchers to gain deeper insight into these preferences and, we hope, to discover new rules that will benefit the field of crystal engineering. We give two examples here: the unequal frequency with which molecules occupy the possible space groups and the spontaneous resolution of chiral molecules.

**3.3.1 Space group preferences.** There are strong space group preferences for experimentally observed crystal structures: over 80% of molecular crystals occupy 5 of the 230 three-dimensional space groups. Having applied an approach in generating trial crystal structures that is unbiased in how the 26 space groups included in our searches are sampled, we analyse the results of CSP to investigate space group preferences within the low energy structures on the set of crystal energy landscapes.

We preface our analysis by emphasising that the space group frequencies presented in Fig. 8 are those for the CSP landscapes, *i.e.* only those of structures generated using an asymmetric unit containing a complete molecule (in these single-species systems,  $Z' = 1$ ) and without detecting and assigning additional symmetry after the minimisation. In contrast, those of the observed crystal structures in the CSD are the full, maximal-symmetry space group (allowing fractions of molecules in the asymmetric unit, *i.e.*  $Z' \leq 1$ ). Hence the space groups enumerated for the CSP set can be thought of as “lower bounds” on the symmetry – higher-symmetry space groups may be assignable if the molecules present internal symmetry.

The global lattice energy minimum for each molecule is the energetically preferred packing; the distribution of space groups among the global minima



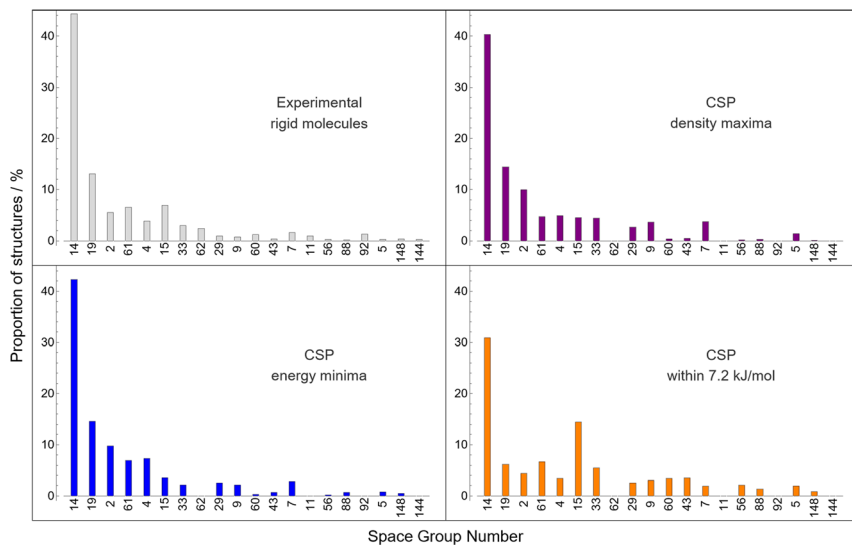


Fig. 8 The relative frequencies of space groups of crystal structures. Grey (top left) are those of rigid molecules in the CSD (as in Fig. 3), while the rest are obtained from CSP landscapes at the global density maximum (purple, top right), at the global energy minimum (blue, bottom left), and within 7.2 kJ mol<sup>-1</sup> of the global energy minimum (orange, bottom right). The ordering of space groups on the x-axis is chosen to match that in Fig. 3. Only the 20 most common space groups for rigid molecules are presented for clarity.

structures are highly consistent with the observed statistics – perhaps as expected, as CSP has been demonstrated in this work to often identify observed structures of these molecules as global minima.

Space group frequencies are often explained by close packing arguments: the commonly observed space groups have combinations of symmetry elements that facilitate close packing of irregular shapes.<sup>64</sup> Examining the space group distribution amongst the densest predicted crystal structure for each molecule gives a similar distribution to that from global energy minima; minimising energy and maximising density lead to similar space group preferences. However, there are differences after the three most popular space groups: the next few have almost equal frequencies among high density structures, suggesting that they are equally good at promoting close packing and that observed differences between these space groups relate to subtler influences of symmetry on lattice energy.

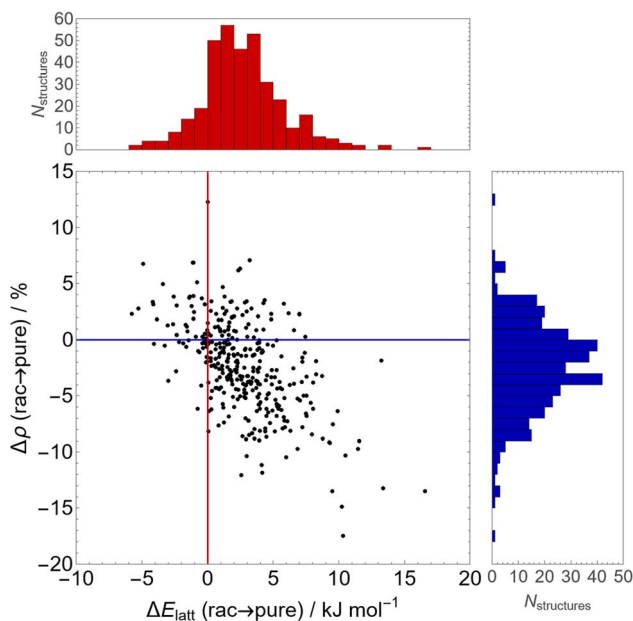
Considering hypothetical structures higher on the CSP energy landscape (up to the usual energy limit of polymorphism, 7.2 kJ mol<sup>-1</sup>, in the bottom right panel of Fig. 8), we see a further flattening of the distribution, and an over-representation of space group 15 (*C2/c*) compared to global energy minima structures. These changes in distribution with energy, which we do not examine in deeper detail here, are only available from access to complete energy landscapes and are relevant in the discovery of high energy, metastable materials, which have sometimes been observed to have attractive properties.<sup>24,65</sup> Thus, we feel that large datasets of CSP landscapes hold potential for generalising our understanding of symmetry preferences in molecular crystals.



**3.3.2 Spontaneous resolution of chiral molecules.** As a second example of insight that can be gained from large numbers of crystal energy landscapes, we examine the tendency for spontaneous resolution of chiral molecules. It is generally accepted that crystallisation from a racemic solution of a chiral molecule more frequently yields racemic crystals rather than undergoing spontaneous resolution into a mixture of crystals, each containing a single stereoisomer.<sup>66</sup> However, information in structural databases alone is limited: knowing whether crystals were grown from a racemic or enantiomerically pure solution is necessary to interpret the incidence rate of spontaneous resolution and the molecular characteristics that influence this behaviour. Furthermore, where enantiomers separate upon crystallisation, it is not possible to grow racemic crystals, so that comparison of racemic vs. enantiomerically pure crystal structures is not possible. Computed crystal energy landscapes make it possible to compare the structures and relative energies of the alternative crystallisation outcomes.

In Fig. 9, we show the difference in stability and density between the energy minimum across all Sohncke (*i.e.* enantiopure) space groups and the minimum across enantiogenic (*i.e.* racemic) space groups for the 356 molecules in our set containing at least one chiral centre. This energy difference represents the propensity for spontaneous resolution of a racemic solution compared to a racemic crystal of both enantiomers.

In general, there is a slight but consistent lattice energy penalty to enantiopurity – the average difference in lattice energy between the enantiopure global



**Fig. 9** The difference in lattice energy ( $\Delta E_{\text{latt}}$ , top red histogram) and relative density ( $\Delta \rho$ , right blue histogram) between the energy minimum in Sohncke (enantiopure) space groups and that in racemic space groups, for all molecules in our set containing at least one chiral centre. The scatter plot (center) displays the relationship between these values for each comparison (red and blue lines indicate the origin, *i.e.* no change in either quantity).





minimum and the racemic global minimum is  $2.7 \text{ kJ mol}^{-1}$ , favouring the racemate. Of the 356 chiral molecules, racemic crystallisation is preferred in 86% (305) of molecules and spontaneous resolution is predicted to occur for approximately 14% (51) of molecules, although those with small lattice energy differences could be influenced by thermal contributions that are not included here. Moreover, this is accompanied (or perhaps driven) by improved packing in the racemate – on average, optimal enantiopure structures are 2.3% less dense than optimal racemic structures, consistent with the empirical Wallach's rule.<sup>67,68</sup>

### 3.4 Machine learned interatomic potentials

One of the more straightforward applications of machine learning in CSP is for the prediction of high quality lattice energies, to reduce the cost of geometry optimisations or of the final energy ranking of structures.<sup>69</sup> Such approaches have been demonstrated in molecular organic CSP by training models on the landscapes of individual molecules.<sup>37,39,40</sup> The CSP dataset developed in this work has significant potential for training data-derived models for organic crystals that could be applied more broadly.

**3.4.1 Transferable  $\Delta$ -ML lattice energy corrections.** To illustrate the use of CSP to train transferable machine-learned energy models, we trained a committee neural network potential using atom-centred symmetry functions<sup>70,71</sup> for lattice energy corrections to the force field used in this work (FIT+DMA), correcting the lattice energies to the B86bPBE+XDM level. An initial model was trained on 7950 selected crystal structures from *ca.* 85% of the CSP landscapes (up to 9 crystal structures per landscape), randomly selected from within  $8 \text{ kJ mol}^{-1}$  of the global energy minimum of each landscape. This corresponds to just under 5% of the crystal structures (166 395) within this energy range for these landscapes. One crystal structure per landscape was withheld as an in-domain test set, while 10 crystal structures per landscape from the remaining CSP landscapes are used as an extrapolation test set. Following initial training, active learning was applied to identify crystal structures from the training landscapes with highest uncertainty in the lattice energy correction. After two iterations (adding 1000 training structures), a slight decrease in errors was observed in the test set, but no improvement in the extrapolation set (see ESI†), so training was halted. We also tried a third iteration of active learning with a wider energy window on each landscape, potentially including more diverse structures, but this did not lower the errors on the test or extrapolation sets. Consequently, we decided to proceed with the model after two iterations of active learning.

The performance of the resulting model on the test set shows remarkably low error (Fig. 10), returning an MAE of just  $0.93 \text{ kJ mol}^{-1}$ . Moreover, a similarly low MAE of  $1.57 \text{ kJ mol}^{-1}$  is achieved on the extrapolation test set, which contained crystal structures of compounds not included in the training of the correction. Compared to the errors for the baseline FIT+DMA, which returned MAEs of 7.80 and  $7.95 \text{ kJ mol}^{-1}$  on the test set and extrapolation set respectively, the correction offers a marked improvement in accuracy. The fact that the errors for the correction are slightly higher on the extrapolation set shows that there are limitations to the transferability of the correction. We expect that better transferability can be achieved as we increase the number of CSP landscapes available for training; the generation of additional landscapes can be targeted to weaknesses in



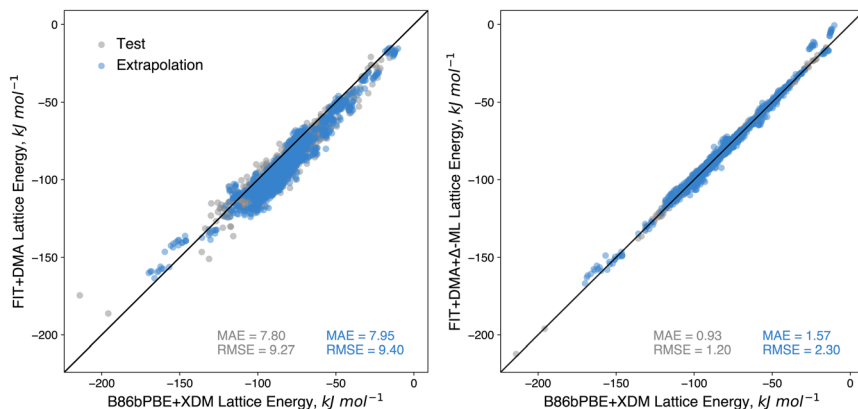


Fig. 10 Correlation of (left) force field (FIT+DMA) lattice energies and (right) force field with machine learned correction vs. DFT (B86bPBE+XDM) lattice energies. The test set (grey data points) consists of unseen crystal structures from CSP landscapes of molecules that were seen during training. The extrapolation set consists of predicted crystal structures from molecules withheld from training.

the underlying force field and molecular types where the current machine learned correction has large errors.

Although the performance of the correction on the test sets is encouraging, an important question is whether the improved energies are significant enough to produce improved stability rankings of organic crystals. To evaluate the influence on ranking, we applied the correction to all of the CSP landscapes, re-ranking all the predicted structures based on the corrected energies. Thereafter, we compared the ranking of the matches to the experimental structures in terms of both their ranking and their energy above the global minimum with that found using FIT+DMA. Structures with predicted uncertainties greater than  $25 \text{ kJ mol}^{-1}$  were omitted from this analysis. This amounted to 257 structures (out of 3.9 million total) with examination indicating the structures were predominately high energy structures containing voids.

The resulting distributions (Fig. 11) over both the training and extrapolation compounds illustrate clearly that the corrected energies in general improve the rankings of the experimental structures. For instance, the correction results in an increase from 424 (FIT+DMA) to 501 (FIT+DMA+ $\Delta$ -ML) of the observed structures ranked as global energy minima, and an increase from 767 (FIT+DMA) to 839 (FIT+DMA+ $\Delta$ -ML) within  $2 \text{ kJ mol}^{-1}$  of the global energy minimum (the approximate limit of vibrational contributions to free energy differences). These improvements are observed in the CSP landscapes of the extrapolation and training molecules (see ESI†). Importantly, the correction also greatly improves the worst ranked experimental matches with the proportion ranked above 35 in energy ranking decreasing by over a third, from 69 (FIT+DMA) to 45 (FIT+DMA+ $\Delta$ -ML) (see ESI†).

**3.4.2 Transferable MACE total energy model.** While the  $\Delta$ -ML approach successfully improves the quality of the CSP energy rankings, some of the CSP matches to experimental crystal structures are still relatively high on their energy landscape, even after applying the lattice energy correction. Considering the demonstrated accuracy of the correction, the high relative energies could be the



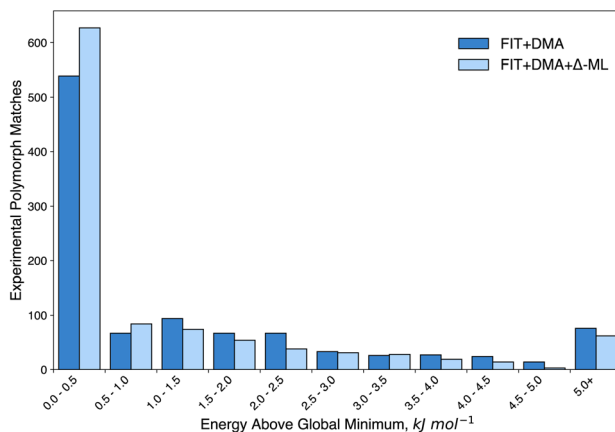


Fig. 11 Histogram of the relative energies of matches to the experimentally-determined crystal structures over all CSPs before (FIT+DMA) and after (FIT+DMA+ $\Delta$ -ML) applying the data-derived B86bPBE+XDM lattice energy correction. Separate distributions for molecules from the training and extrapolation sets are shown in the ESI.†

result of limitations with the rigid-body lattice energy minimised geometries, which an energy correction is unable to remedy. Indeed, many of the structures with high relative energies are for compounds which had large geometric deviations between the experimental in-crystal molecular conformations and the gas-phase optimised molecular geometries used (and kept rigid) during CSP. Improving the performance of CSP for these more flexible molecules in our study likely requires re-optimising the predicted crystal structures and relaxing the rigid-molecule constraint.

The FIT+DMA+ $\Delta$ -ML model is not suitable for this task because it is only a correction to the intermolecular contribution to the lattice energy. Therefore, using a dataset derived from the CSP structures selected for the energy correction model, with perturbed atomic coordinates to sample conformational degrees of freedom, we trained a total energy MACE (higher order equivariant message passing neural network) model. Full details are provided in the ESI.† The trained MACE model was then applied to geometry optimise the predicted structures of 15 compounds with large differences in molecular geometry between the observed crystal structure and the DFT-optimised molecule (see Fig. 1, bottom three rows).

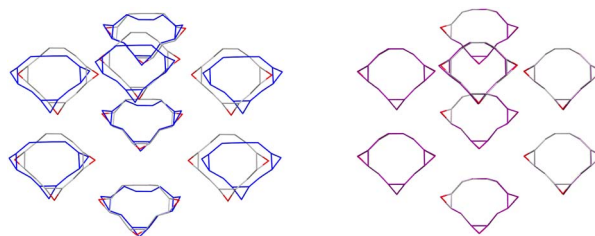
Re-optimisation with the trained MACE model yielded considerable improvements in the geometric agreement of predicted structures with experiment and of their energy ranking on the CSP landscapes (Table 2). RMSD<sub>30</sub> between experimental and predicted structures decreased upon re-optimisation for 14 of the 15 compounds, by up to 1.4 Å, and moved the match to the observed structure closer to the global energy minimum in all but 2 cases, with 7 becoming the global energy minimum structure. Fig. 12 illustrates the improved geometric agreement for one of these molecules.

Re-optimisation of CSP structures was also run for 4,8*b*-dihydropyrrolo[3,4-*b*]indole-1,3(2*H*,3*aH*)-dione, where no match to the experimental crystal structure (CSD reference code BUVGAC) was identified in the CSP. However, no match was identified after re-optimisation with the MACE model; in this case, molecular



**Table 2** Crystal structure RMSD<sub>30</sub> and energetic ranking of matches to the experimentally-determined crystal structures within CSP landscapes for 15 molecules re-optimised with the transferable MACE model. Changes between the rigid-molecule force field results (FIT+DMA) and MACE are shown in parentheses. Bold entries highlight where MACE re-optimisation leads to improvement

Crystal structure	FIT+DMA		MACE	
	RMSD <sub>30</sub> (Å)	$\Delta E$ (kJ mol <sup>-1</sup> )	RMSD <sub>30</sub> (Å)	$\Delta E$ (kJ mol <sup>-1</sup> )
ATCDEO	0.341	7.26	<b>0.125 (-0.216)</b>	<b>6.67 (-0.59)</b>
BIRTUR	0.396	6.06	<b>0.391 (-0.005)</b>	<b>0.02 (-6.03)</b>
BIXKUP	0.813	16.32	<b>0.208 (-0.605)</b>	<b>0.00 (-16.32)</b>
BIXLOK	1.048	4.95	<b>0.265 (-0.783)</b>	<b>0.00 (-4.95)</b>
DALBIC	0.343	1.45	<b>0.206 (-0.137)</b>	3.70 (+2.25)
DEBFOI	1.208	4.18	<b>0.110 (-1.098)</b>	<b>0.00 (-4.18)</b>
DOBYOJ	0.654	1.27	<b>0.640 (-0.014)</b>	3.55 (+2.28)
JASGIT	1.837	15.77	<b>1.776 (-0.061)</b>	<b>9.72 (-6.05)</b>
KOGFER	0.247	0.37	0.248 (+0.001)	<b>0.00 (-0.37)</b>
QEKQIG	0.824	4.84	<b>0.139 (-0.684)</b>	<b>0.00 (-4.84)</b>
TAVTUH	1.751	12.03	<b>0.325 (-1.426)</b>	<b>1.87 (-10.16)</b>
TUNWUV	1.103	9.55	<b>0.399 (-0.704)</b>	<b>0.00 (-9.55)</b>
UDEXUZ	0.725	7.41	<b>0.216 (-0.509)</b>	<b>0.75 (-6.66)</b>
WACYAB	0.721	9.81	<b>0.216 (-0.505)</b>	<b>0.92 (-8.89)</b>
WACYEF	0.975	14.02	<b>0.200 (-0.775)</b>	<b>0.00 (-14.02)</b>



**Fig. 12** Overlay of the experimentally determined crystal structure (atoms coloured by element) of (exo,exo,exo)-1,2:4,5:7,8-triepoxy cyclododec-10-ene (CSD reference code WACYEF) with: (left) the matching structure from the force field (FIT+DMA) CSP (blue); and (right) the matching CSP structure after re-optimisation with the transferable MACE model (purple). Hydrogen atoms are hidden for clarity. The large structural deviation in the FIT+DMA structure is driven by deviation in the molecular geometry.

flexibility would have been required during structure generation, rather than post-CSP re-optimisation. Now that a transferable MACE model has been trained, it could potentially be implemented earlier in the CSP workflow.

It should be noted though that the MACE model is a semi-local model and so neglects long-range interactions that can be important for highly accurate relative energies. Overall, the results of the unconstrained optimisations confirm that the high relative energies of these crystal structure matches estimated by FIT+DMA+ $\Delta$ -ML are largely a result of limitations with the rigid-body lattice energy minimised geometries. It is worth reiterating that the limitations of the CSP geometries only affected a small number of structures; in the vast majority of



cases, as shown earlier, the predicted structures at the FIT+DMA level achieved high quality matches to the experimentally-determined crystal structures. Undoubtedly, the success of the FIT+DMA+ $\Delta$ -ML energy correction is a reflection of the performance of the baseline empirical force field.

## 4 Conclusions

We have presented (to our knowledge) the largest CSP dataset produced to-date, serving as a computational survey of crystal packing in the organic solid state for small, rigid molecules. Using established, well-characterised CSP methods, we produced crystal structure energy landscapes for over 1000 such molecules, all of which have at least one known, solved crystal structure available in the CSD. In total, our CSP landscapes contain over 4 million unique crystal structures, each with an associated lattice energy at a consistent level of theory.

We have assessed the quality of the dataset by evaluating its reliability at predicting the known crystal structures of these molecules, both in terms of the quality of the geometric match of the crystal structures and the resulting energy ranking of the experimental forms on their respective landscapes compared to other hypothetical structures. Our CSP approach is overwhelmingly successful at predicting crystal structures of these simple molecules – over 99% of all experimental structures have a match located in our CSP searches. 41% of experimental structures are predicted to be the global energy minimum on their landscapes, and 74% are found within 2 kJ mol<sup>-1</sup> of this minimum, a margin equivalent to the estimated error introduced by ignoring thermal effects and ranking based solely on static, 0 K lattice energies. Geometrically, the typical discrepancy between experimental structures and their closest predicted matches is comparable to the thermal fluctuations in experimental crystal structure solutions of the same solid form obtained at low temperatures *versus* ambient conditions. Such remarkable performance demonstrates the consistency and accuracy of our chosen methods for optimising and ranking these structures.

Such a large dataset of many possible crystal packings should prove a valuable resource for identifying a variety of crystal packing trends, and we make this data available to the community as part of this work. Herein, we studied space group distributions among low-energy hypothetical structures compared to those observed in the CSD, and find substantial overlap, particularly among the most common space groups. We also demonstrated the potential of this varied dataset to explore chirality in the organic solid state, finding very good agreement with established empirical rules concerning the propensity of racemic mixtures to crystallise in racemic crystal structures as opposed to separate enantiopure crystals.

Additionally, we have shown the power of such large-scale CSP to train transferable machine learned potentials for organic solid-state systems. A committee neural network potential trained on single-point periodic DFT lattice energies achieved excellent accuracy in correcting our force field energy landscapes to the DFT level, reducing energy errors by approximately 8 fold. While the NNP performance is reduced on molecules reserved as an extrapolation set compared to those seen in training, the potential still demonstrates improvements to the quality of the resulting CSP rankings increasing the number of experimental structures ranked as the global energy minimum by 18% overall. We further demonstrated the development of a transferable MACE potential using



structures derived from the CSP landscapes to allow re-optimisation of crystal structures, testing it successfully on those molecules in our set where the molecular geometry distortion between the conformation in the known crystal structure and the gas-phase-optimised one was largest. The results showed improved structural agreement with experimental structures in almost all cases and much improved energy rankings, moving several poorly-ranked observed structures to the global energy minimum.

Using these mature, well-tested CSP methods alongside modern machine learning approaches, we have demonstrated the ability of CSP to create very large, diverse datasets of hypothetical crystal structures, and the utility of this information in both understanding broad trends in organic crystal structures and in training more advanced energetic models for refinement and transferability. It is our hope that the variety and quantity of CSP data presented here, alongside our demonstrations of possible applications, enables the greater organic solid-state computational community to develop even more sophisticated models and techniques in pursuit of truly predictive, computational data-driven discovery of novel materials.

## Author contributions

Christopher R. Taylor: conceptualisation, methodology, software, investigation, validation, formal analysis, data curation, writing – original draft, writing – review & editing. Led on production, curation and analysis of large-scale CSP datasets. Patrick W. V. Butler: methodology, software, investigation, validation, formal analysis, writing – original draft, writing – review & editing. Training of ML energy models and their evaluation. Graeme M. Day: conceptualisation, validation, formal analysis, resources, writing – original draft, writing – review & editing, supervision, funding acquisition.

## Conflicts of interest

There are no conflicts to declare.

## Acknowledgements

This project has received funding from the European Research Council under the European Union's Horizon 2020 Research and Innovation Program (grant agreement no. 856405, ERC-ADAM) and from the AI for Chemistry: AICChem hub (EPSRC grant EP/Y028759/1). P. W. V. B. is funded by a University of Southampton Presidential Scholarship. The authors acknowledge the use of the IRIDIS High Performance Computing Facility, and associated support services at the University of Southampton, in the completion of this work. *Via* our membership of the UK's HEC Materials Chemistry Consortium, which is funded by EPSRC (EP/X035859), this work used the ARCHER2 UK National Supercomputing Service (<https://www.archer2.ac.uk>).

## Notes and references

- 1 C. R. Groom, I. J. Bruno, M. P. Lightfoot and S. C. Ward, *Acta Crystallogr., Sect. B: Struct. Sci., Cryst. Eng. Mater.*, 2016, **72**, 171–179.



- 2 D. Zagorac, H. Müller, S. Ruehl, J. Zagorac and S. Rehme, *J. Appl. Crystallogr.*, 2019, **52**, 918–925.
- 3 M. M. J. Treacy, I. Rivin, E. Balkovsky, K. H. Randall and M. D. Foster, *Microporous Mesoporous Mater.*, 2004, **74**, 121–132.
- 4 A. Le Bail, *J. Appl. Crystallogr.*, 2005, **38**, 389–395.
- 5 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, *J. Appl. Crystallogr.*, 2009, **42**, 726–729.
- 6 A. Vaitkus, A. Merkys and S. Gražulis, *J. Appl. Crystallogr.*, 2021, **54**, 661–672.
- 7 A. Vaitkus, A. Merkys, T. Sander, M. Quirós, P. A. Thiessen, E. E. Bolton and S. Gražulis, *J. Cheminf.*, 2023, **15**, 123.
- 8 A. Jain, S. P. Ong, G. Hautier, W. Chen, W. D. Richards, S. Dacek, S. Cholia, D. Gunter, D. Skinner, G. Ceder and K. A. Persson, *APL Mater.*, 2013, **1**, 011002.
- 9 J. Nyman and G. M. Day, *Phys. Chem. Chem. Phys.*, 2016, **18**, 31132–31143.
- 10 J. Nyman and G. M. Day, *CrystEngComm*, 2015, **17**, 5154–5165.
- 11 A. J. Cruz-Cabeza and J. Bernstein, *Chem. Rev.*, 2014, **114**, 2170–2191.
- 12 A. J. Cruz-Cabeza, S. M. Reutzel-Edens and J. Bernstein, *Chem. Soc. Rev.*, 2015, **44**, 8619–8635.
- 13 P. Sacchi, P. Neoptolemos, R. J. Davey, S. M. Reutzel-Edens and A. J. Cruz-Cabeza, *Chem. Sci.*, 2023, **14**, 11775–11789.
- 14 H. P. G. Thompson and G. M. Day, *Chem. Sci.*, 2014, **5**, 3173–3182.
- 15 H. C. S. Chan, J. Kendrick, M. A. Neumann and F. J. J. Leusen, *CrystEngComm*, 2013, **15**, 3799–3807.
- 16 C. R. Taylor and G. M. Day, *Cryst. Growth Des.*, 2018, **18**, 892–904.
- 17 L. M. LeBlanc, S. G. Dale, C. R. Taylor, A. D. Becke, G. M. Day and E. R. Johnson, *Angew. Chem., Int. Ed.*, 2018, **57**, 14906–14910.
- 18 J. Hoja and A. Tkatchenko, *Faraday Discuss.*, 2018, **211**, 253–274.
- 19 G. M. Day, *Crystallogr. Rev.*, 2011, **17**, 3–52.
- 20 G. J. O. Beran, *Chem. Sci.*, 2023, **14**, 13290–13312.
- 21 S. Yang and G. M. Day, *Commun. Chem.*, 2022, **5**, 1–13.
- 22 N. F. Francia, L. S. Price, J. Nyman, S. L. Price and M. Salvalaglio, *Cryst. Growth Des.*, 2020, **20**, 6847–6862.
- 23 P. W. V. Butler and G. M. Day, *Proc. Natl. Acad. Sci. U. S. A.*, 2023, **120**, e2300516120.
- 24 A. Pulido, L. Chen, T. Kaczorowski, D. Holden, M. A. Little, S. Y. Chong, B. J. Slater, D. P. McMahon, B. Bonillo, C. J. Stackhouse, A. Stephenson, C. M. Kane, R. Clowes, T. Hasell, A. I. Cooper and G. M. Day, *Nature*, 2017, **543**, 657–664.
- 25 J. E. Campbell, J. Yang and G. M. Day, *J. Mater. Chem. C*, 2017, **5**, 7574–7584.
- 26 G. M. Day and A. I. Cooper, *Adv. Mater.*, 2018, **30**, 1704944.
- 27 C. Zhao, L. Chen, Y. Che, Z. Pang, X. Wu, Y. Lu, H. Liu, G. M. Day and A. I. Cooper, *Nat. Commun.*, 2021, **12**, 817.
- 28 D. H. Case, J. E. Campbell, P. J. Bygrave and G. M. Day, *J. Chem. Theory Comput.*, 2016, **12**, 910–924.
- 29 P. G. Karamertzanis and C. C. Pantelides, *J. Comput. Chem.*, 2005, **26**, 304–324.
- 30 F. Curtis, X. Li, T. Rose, A. Vásquez-Mayagoitia, S. Bhattacharya, L. M. Ghiringhelli and N. Marom, *J. Chem. Theory Comput.*, 2018, **14**, 2246–2264.





- 31 M. A. Neumann, J. v. d. Streek, F. P. A. Fabbiani, P. Hidber and O. Grassmann, *Nat. Commun.*, 2015, **6**, 1–7.
- 32 D.-K. Bučar, G. M. Day, I. Halasz, G. G. Z. Zhang, J. R. G. Sander, D. G. Reid, L. R. MacGillivray, M. J. Duer and W. Jones, *Chem. Sci.*, 2013, **4**, 4417–4425.
- 33 C. R. Taylor, M. T. Mulvey, D. S. Perenyi, M. R. Probert, G. M. Day and J. W. Steed, *J. Am. Chem. Soc.*, 2020, **142**, 16668–16680.
- 34 D. H. Case, V. K. Srirambhatla, R. Guo, R. E. Watson, L. S. Price, H. Polyzois, J. K. Cockcroft, A. J. Florence, D. A. Tocher and S. L. Price, *Cryst. Growth Des.*, 2018, **18**, 5322–5331.
- 35 A. M. Reilly, R. I. Cooper, C. S. Adjiman, S. Bhattacharya, A. D. Boese, J. G. Brandenburg, P. J. Bygrave, R. Bylsma, J. E. Campbell, R. Car, D. H. Case, R. Chadha, J. C. Cole, K. Cosburn, H. M. Cuppen, F. Curtis, G. M. Day, R. A. DiStasio Jr, A. Dzyabchenko, B. P. van Eijck, D. M. Elking, J. A. van den Ende, J. C. Facelli, M. B. Ferraro, L. Fusti-Molnar, C.-A. Gatsiou, T. S. Gee, R. de Gelder, L. M. Ghiringhelli, H. Goto, S. Grimme, R. Guo, D. W. M. Hofmann, J. Hoja, R. K. Hylton, L. Iuzzolino, W. Jankiewicz, D. T. de Jong, J. Kendrick, N. J. J. de Klerk, H.-Y. Ko, L. N. Kuleshova, X. Li, S. Lohani, F. J. J. Leusen, A. M. Lund, J. Lv, Y. Ma, N. Marom, A. E. Masunov, P. McCabe, D. P. McMahon, H. Meekes, M. P. Metz, A. J. Misquitta, S. Mohamed, B. Monserrat, R. J. Needs, M. A. Neumann, J. Nyman, S. Obata, H. Oberhofer, A. R. Oganov, A. M. Orendt, G. I. Pagola, C. C. Pantelides, C. J. Pickard, R. Podeszwa, L. S. Price, S. L. Price, A. Pulido, M. G. Read, K. Reuter, E. Schneider, C. Schober, G. P. Shields, P. Singh, I. J. Sugden, K. Szalewicz, C. R. Taylor, A. Tkatchenko, M. E. Tuckerman, F. Vacarro, M. Vasileiadis, A. Vazquez-Mayagoitia, L. Vogt, Y. Wang, R. E. Watson, G. A. de Wijs, J. Yang, Q. Zhu and C. R. Groom, *Acta Crystallogr., Sect. B: Struct. Sci., Cryst. Eng. Mater.*, 2016, **72**, 439–459.
- 36 J. Behler and G. Csányi, *Eur. Phys. J. B*, 2021, **94**, 142.
- 37 F. Musil, S. De, J. Yang, J. E. Campbell, G. M. Day and M. Ceriotti, *Chem. Sci.*, 2018, **9**, 1289–1300.
- 38 D. McDonagh, C.-K. Skylaris and G. M. Day, *J. Chem. Theory Comput.*, 2019, **15**, 2743–2758.
- 39 O. Egorova, R. Hafizi, D. C. Woods and G. M. Day, *J. Phys. Chem. A*, 2020, **124**, 8065–8078.
- 40 P. W. V. Butler, R. Hafizi and G. M. Day, *J. Phys. Chem. A*, 2024, **128**, 945–957.
- 41 V. Kapil and E. A. Engel, *Proc. Natl. Acad. Sci. U. S. A.*, 2022, **119**, e2111769119.
- 42 D. P. Kovács, J. H. Moore, N. J. Browning, I. Batatia, J. T. Horton, V. Kapil, W. C. Witt, I.-B. Magdău, D. J. Cole and G. Csányi, MACE-OFF23: Transferable Machine Learning Force Fields for Organic Molecules, *arXiv*, 2023, preprint, arXiv:2312.15211[physics], DOI: [10.48550/arXiv.2312.15211](https://doi.org/10.48550/arXiv.2312.15211).
- 43 F. M. Paruzzo, A. Hofstetter, F. Musil, S. De, M. Ceriotti and L. Emsley, *Nat. Commun.*, 2018, **9**, 4501.
- 44 P. T. A. Galek, L. Fábrián, W. D. S. Motherwell, F. H. Allen and N. Feeder, *Acta Crystallogr., Sect. B: Struct. Sci.*, 2007, **63**, 768–782.
- 45 R. K. Cersonsky, M. Pakhnova, E. A. Engel and M. Ceriotti, *Chem. Sci.*, 2023, **14**, 1272–1285.



- 46 J. Yang, S. De, J. E. Campbell, S. Li, M. Ceriotti and G. M. Day, *Chem. Mater.*, 2018, **30**, 4361–4371.
- 47 M. J. Frisch, G. W. Trucks, H. B. Schlegel, G. E. Scuseria, M. A. Robb, J. R. Cheeseman, G. Scalmani, V. Barone, B. Mennucci, G. A. Petersson, H. Nakatsuji, M. Caricato, X. Li, H. P. Hratchian, A. F. Izmaylov, J. Bloino, G. Zheng, J. L. Sonnenberg, M. Hada, M. Ehara, K. Toyota, R. Fukuda, J. Hasegawa, M. Ishida, T. Nakajima, Y. Honda, O. Kitao, H. Nakai, T. Vreven, J. A. Montgomery Jr, J. E. Peralta, F. Ogliaro, M. Bearpark, J. J. Heyd, E. Brothers, K. N. Kudin, V. N. Staroverov, R. Kobayashi, J. Normand, K. Raghavachari, A. Rendell, J. C. Burant, S. S. Iyengar, J. Tomasi, M. Cossi, N. Rega, J. M. Millam, M. Klene, J. E. Knox, J. B. Cross, V. Bakken, C. Adamo, J. Jaramillo, R. Gomperts, R. E. Stratmann, O. Yazyev, A. J. Austin, R. Cammi, C. Pomelli, J. W. Ochterski, R. L. Martin, K. Morokuma, V. G. Zakrzewski, G. A. Voth, P. Salvador, J. J. Dannenberg, S. Dapprich, A. D. Daniels, Ö. Farkas, J. B. Foresman, J. V. Ortiz, J. Cioslowski and D. J. Fox, *Gaussian09 Revision D.01*.
- 48 C. Lee, W. Yang and R. G. Parr, *Phys. Rev. B: Condens. Matter Mater. Phys.*, 1988, **37**, 785–789.
- 49 A. D. Becke, *J. Chem. Phys.*, 1993, **98**, 5648–5652.
- 50 K. Raghavachari, J. S. Binkley, R. Seeger and J. A. Pople, *J. Chem. Phys.*, 1980, **72**, 650–654.
- 51 S. Grimme, J. Antony, S. Ehrlich and H. Krieg, *J. Chem. Phys.*, 2010, **132**, 154104.
- 52 A. D. Becke and E. R. Johnson, *J. Chem. Phys.*, 2005, **123**, 154101.
- 53 A. J. Stone and M. Alderton, *Mol. Phys.*, 2002, **100**, 221–233.
- 54 A. J. Stone, *J. Chem. Theory Comput.*, 2005, **1**, 1128–1132.
- 55 P. J. Winn, G. G. Ferenczy and C. A. Reynolds, *J. Phys. Chem. A*, 1997, **101**, 5437–5445.
- 56 G. G. Ferenczy, P. J. Winn and C. A. Reynolds, *J. Phys. Chem. A*, 1997, **101**, 5446–5455.
- 57 D. S. Coombes, S. L. Price, D. J. Willock and M. Leslie, *J. Phys. Chem.*, 1996, **100**, 7352–7360.
- 58 S. L. Price, M. Leslie, G. W. A. Welch, M. Habgood, L. S. Price, P. G. Karamertzanis and G. M. Day, *Phys. Chem. Chem. Phys.*, 2010, **12**, 8478–8490.
- 59 D. E. Williams and D. J. Houpt, *Acta Crystallogr., Sect. B: Struct. Sci.*, 1986, **42**, 286–295.
- 60 J. A. Chisholm and S. Motherwell, *J. Appl. Crystallogr.*, 2005, **38**, 228–231.
- 61 RDKit: Open-Source Cheminformatics, <https://www.rdkit.org>.
- 62 C. Peifer, M. Tschertsche, D. Schollmeyer and S. Laufer, *Acta Crystallogr., Sect. E: Struct. Rep. Online*, 2007, **63**, o1359–o1360.
- 63 C. Wilson, *Z. Kristallogr. – Cryst. Mater.*, 2000, **215**, 693–701.
- 64 M. J. Buerger, *Science*, 1962, **135**, 912.
- 65 C. M. Aitchison, C. M. Kane, D. P. McMahon, P. R. Spackman, A. Pulido, X. Wang, L. Wilbraham, L. Chen, R. Clowes, M. A. Zwijnenburg, R. S. Sprick, M. A. Little, G. M. Day and A. I. Cooper, *J. Mater. Chem. A*, 2020, **8**, 7158–7170.
- 66 C. P. Brock and J. D. Dunitz, *Chem. Mater.*, 1994, **6**, 1118–1127.
- 67 O. Wallach, *Justus Liebigs Ann. Chem.*, 1895, **286**, 90–118.



- 68 C. P. Brock, W. B. Schweizer and J. D. Dunitz, *J. Am. Chem. Soc.*, 1991, **113**, 9811–9820.
- 69 R. J. Clements, J. Dickman, J. Johal, J. Martin, J. Glover and G. M. Day, *MRS Bull.*, 2022, **47**, 1054–1062.
- 70 J. Behler, *J. Chem. Phys.*, 2011, **134**, 074106.
- 71 M. Gastegger, L. Schwiedrzik, M. Bittermann, F. Berzsenyi and P. Marquetand, *J. Chem. Phys.*, 2018, **148**, 241709.

