

PNAS



1

2 **Supporting Information for**

3 **Reducing Overprediction of Molecular Crystal Structures *via* Threshold Clustering**

4 **Patrick W. V. Butler, Graeme M. Day**

5 **Graeme M. Day**

6 **E-mail: g.m.day@soton.ac.uk**

7 **This PDF file includes:**

8 Supporting text

9 Figs. S1 to S7

10 Table S1

11 SI References

12 Supporting Information Text

13 1. Crystal Structure Prediction

14 The initial CSP landscapes were generated for benzene, acrylic acid, and resorcinol using our GLEE(1) program. For the CSPs
15 of benzene and acrylic acid, we followed our previously described methodology based on rigid-body lattice optimisations using an
16 empirically parametrised intermolecular atom-atom exp-6 potential combined with atomic multipole electrostatics (FIT+DMA),
17 the parameters sourced from the FIT(2) force field. The molecular geometries were optimised at the B3LYP/6-311G(d,p) level
18 using Gaussian09(3) and held fixed through the rest of the search. Distributed, atom-centered multipoles up to hexadecapole
19 were derived from the resulting electron density by a distributed multipole analysis and partial charges were fitted to the
20 multipoles.(4, 5) A quasi-random search of the lattice packing space with one molecule in the asymmetric unit was then
21 conducted in selected space groups. For benzene and acrylic acid the 25 most common space groups were searched. Valid
22 structures were lattice energy minimised using the software packages PMIN(6) and DMACRYST(7) in a 3-stage protocol
23 consisting of: PMIN at ambient pressure with partial charge electrostatics, FIT+DMA at 0.1 GPa with multipole electrostatics,
24 and lastly the FIT+DMA once more at ambient pressure with multipole electrostatics. The search was terminated for both
25 systems once 250,000 valid structures were generated. In each case, the structure set was clustered iteratively by comparison of
26 simulated powder X-ray diffraction (pXRD) patterns generated by PLATON(8) followed by structural overlay comparisons
27 using the CSD API with a molecular cluster size of 15 molecules and an RMSD cutoff in atomic positions of 0.3 Å.(9) The
28 unique structures were then ranked by the lattice energy of the final optimization stage to yield the CSP landscapes.

29 In the case of resorcinol, the CSP landscape was generated by applying our recently developed flexible-molecule CSP
30 protocol.* This protocol is largely similar to that described for the rigid systems, however, rather than searching the lattice
31 packing space of a single conformation we instead search with a pre-calculated pool of rigid conformations. Structures are
32 then generated by randomly selecting a conformation from the pool. For resorcinol, the pool was generated by fixing one of
33 the hydroxyl group torsions in an anti position while stepping the other through 360 degrees in 40 degree increments. The
34 space group of the experimental alpha and beta forms, $Pna2_1$, was then searched generating 10,000 valid structures. These
35 were lattice energy minimized initially by the same 3-stage protocol, however, after clustering the unique structures were
36 further optimised with D3 dispersion-corrected tight-binding DFT (DFTB+D3) as implemented in DFTB+(10) using the 3ob
37 parameter set to allow the conformations to relax within the the lattice.(11)

*Manuscript in preparation

38 2. Threshold Clustering

39 For benzene, acrylic acid, and resorcinol the lowest energy structures from the CSP landscapes were included in the threshold
40 simulations. The structures were expanded to *P1* cells to remove symmetry constraints on the sampling, with structures in
41 higher symmetry, centered space groups, such that the corresponding *P1* cells had more than 8 molecules in the cell, reduced
42 to smaller primitive cells for computational efficiency. The Monte Carlo trajectories in all simulations employed the same
43 core moveset consisting of molecular translations and rotations, cell length and angle changes, and cell volume changes. The
44 resorcinol simulations included additionally torsional moves. The maximum move amounts were calibrated for each system by
45 running short simulations and adjusting the cutoffs to achieve an absolute average energy change of around 1 kJ mol⁻¹. If
46 necessary, these were refined further to produce an acceptance ratio of approximately 50-60%. The probability of choosing a
47 given move was set equal to the proportion of the total degrees of freedom it represented. In the case of torsion moves the
48 probability was increased three times this base value in order to better reflect the importance of this degree of freedom.

49 The number and length of Monte Carlo trajectories was set for each system based on the expected complexity of the energy
50 surface following smaller preliminary simulations of the systems. For benzene 3 trajectories of 30,000 steps (15,000 with 2.5 kJ
51 mol⁻¹ lid then 15,000 with 5.0 kJ mol⁻¹ lid) were initiated from each CSP structures. For acrylic acid and resorcinol the same
52 number of trajectories were initiated per CSP structures but the length of the trajectories was 20,000 steps (single lid of 5.0
53 kJ mol⁻¹). In all simulations, the single point evaluations of the perturbed structures and the lattice energy minimizations
54 of accepted structures were performed with the same energy model as used to calculate the CSP landscapes. For benzene
55 and acrylic acid this was FIT+DMA and for resorcinol this was DFTB+D3. Following the simulations, the trajectories were
56 combined by comparing pairs of energy minimized structures across all trajectories in an iterative procedure consisting of initial
57 comparisons of simulated pXRD patterns generated by PLATON(8) followed by structural overlays using the CSD API(9)
58 where matches were identified by root-mean-square differences (RMSD) in atomic positions of less than 0.3 Å using a molecular
59 cluster of 15 molecules. In the case of resorcinol, due to the poor sensitivity of pXRD to hydrogen position, the equivalent
60 structures from pXRD comparisons were checked to ensure they were the same conformation. From the clustering, overlapping
61 trajectories are readily identifiable and the disconnectivity graph was then constructed assuming that all overlapping trajectories
62 represent a single basin.

63 It was observed that the initial *P1* minimizations at the start of the threshold simulations led to some minima coalescing.
64 Consequently, the disconnectivity graphs appear to have fewer 'initial minima' than expected. For benzene the 5.0 kJ mol⁻¹
65 disconnectivity graph indicates 90 distinct initial structures from the 100 CSP structures, for acrylic acid it is 97 from 100, and
66 for resorcinol it is 45 from 50. In the case of benzene and resorcinol, this can be partly attributed to combining CSP landscapes
67 generated by searching selected space groups with simulations in *P1* because some minima in the space group symmetry may
68 not be minima when the symmetry is removed (due to the additional degrees of freedom). However, this cannot be entirely
69 the cause considering the same behaviour is observed for the resorcinol structures, which were optimised in *P1* as part of the
70 DFTB stage before the simulations. From this, it seems that some of the coalescing minima may then also be due to instability
71 in the minimizer algorithms. Removing the duplicate trajectories of coalesced structures was not found to alter the results,
72 which is expected since in order to coalesce the two structures must be within the same basin. Nevertheless, this does represent
73 inefficient sampling, and thus ideally threshold clustering should be combined with CSP landscapes wherein the structures are
74 true minima on the *P1* energy surface.

75 **3. Comparison of Threshold Clustering and MD-based Methodologies**

76 To further expand on the comparison of methods for reducing CSP overprediction we have collected here the advantages and
77 disadvantages of the proposed threshold clustering method and the established MD with enhanced sampling approach.

78 Firstly, as stated in the main text in terms of complexity of the workflow and software required, threshold clustering is
79 much simpler than the MD approaches so far proposed. However, MD software is more well developed and accessible, which
80 has practical benefits.

81 The enhanced sampling dynamics approaches, such as metadynamics would be expected to more efficiently search high
82 dimensional configurational spaces, such as that of highly flexible molecules, than random walkers. However, these typically
83 rely on enhanced sampling along a collective variable (CV), which introduces sensitivity of the results to the choice of CV. Of
84 course, this is not a like-for-like comparison and there are more efficient MC methods for exploring complex spaces that could
85 be applied to improve the performance of threshold clustering.

86 Finally, threshold clustering is readily compatible with the static energy methods that are ubiquitous in the lattice energy
87 approach of conventional CSP. While MD methods could in principle use these energy models, practically it is not an insignificant
88 challenge, which is evidenced by the lack of examples in the literature. Overall, the performance of these methods will likely
89 vary depending on the system and a perfect comparison will be challenging. Moreover, we do not expect such comparisons to
90 be helpful in regard to advancing methods for reducing overprediction and thus such work is not in our future plans.

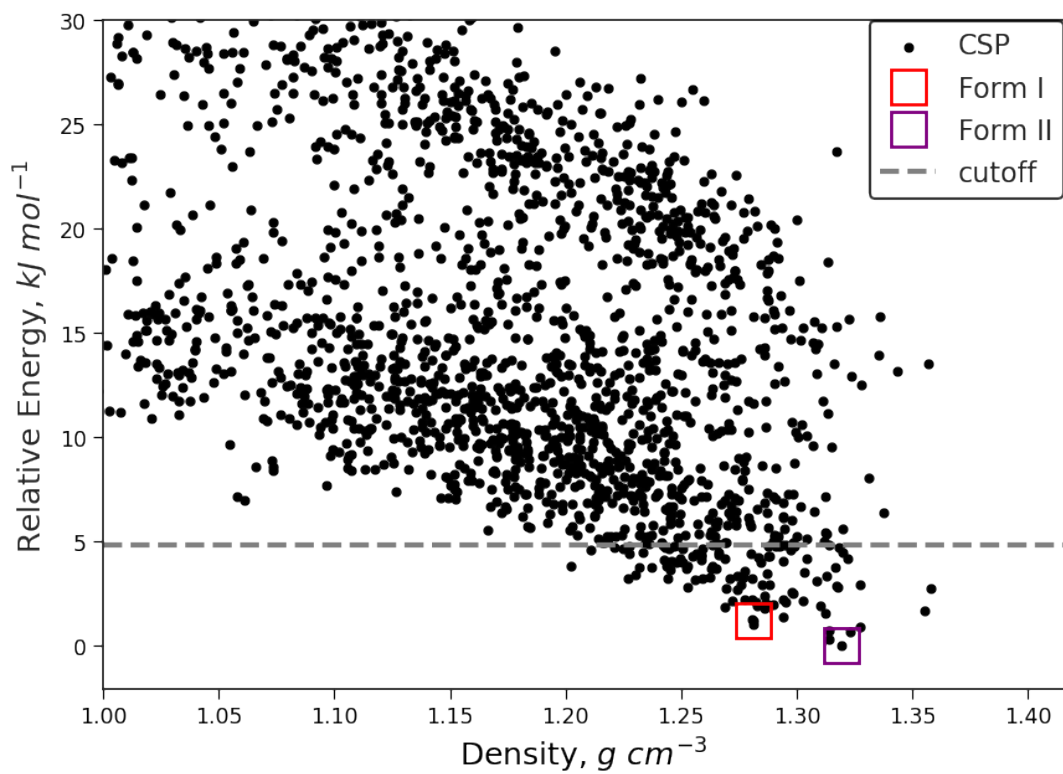


Fig. S1. Initial predicted landscape for acrylic acid. The matches to the experimental polymorphs and the energy cutoff for the structures included in threshold clustering simulations are shown.

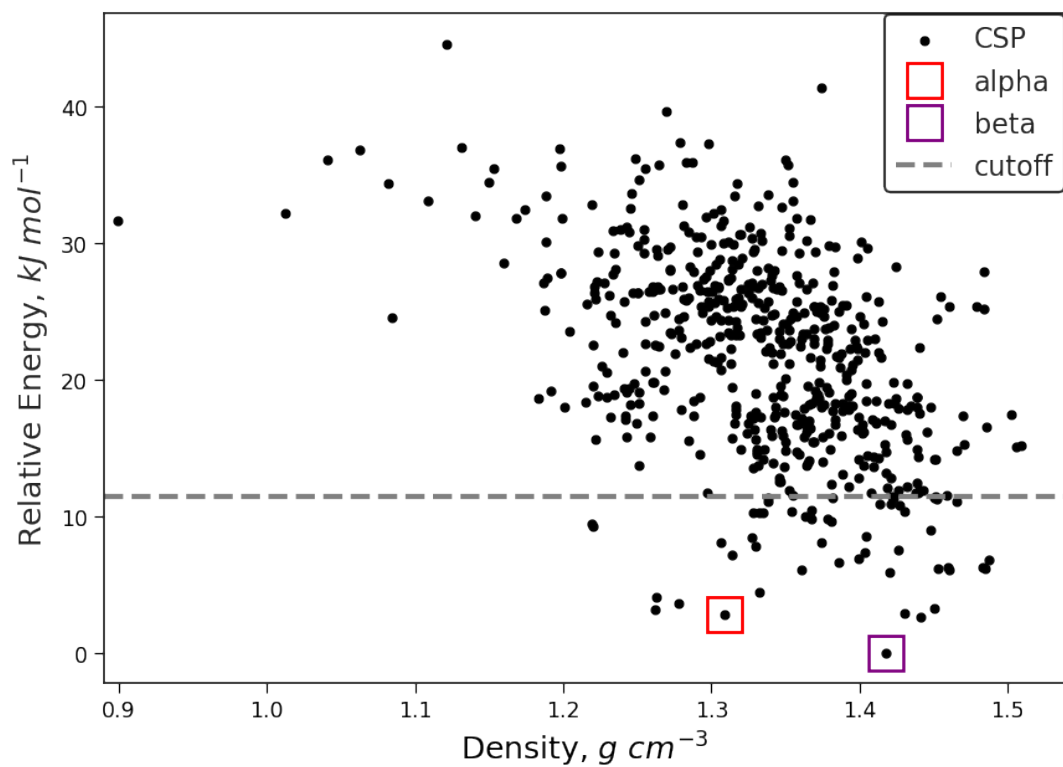


Fig. S2. Initial predicted landscape for resorcinol. The matches to the experimental polymorphs and the energy cutoff for the structures included in threshold clustering simulations are shown.

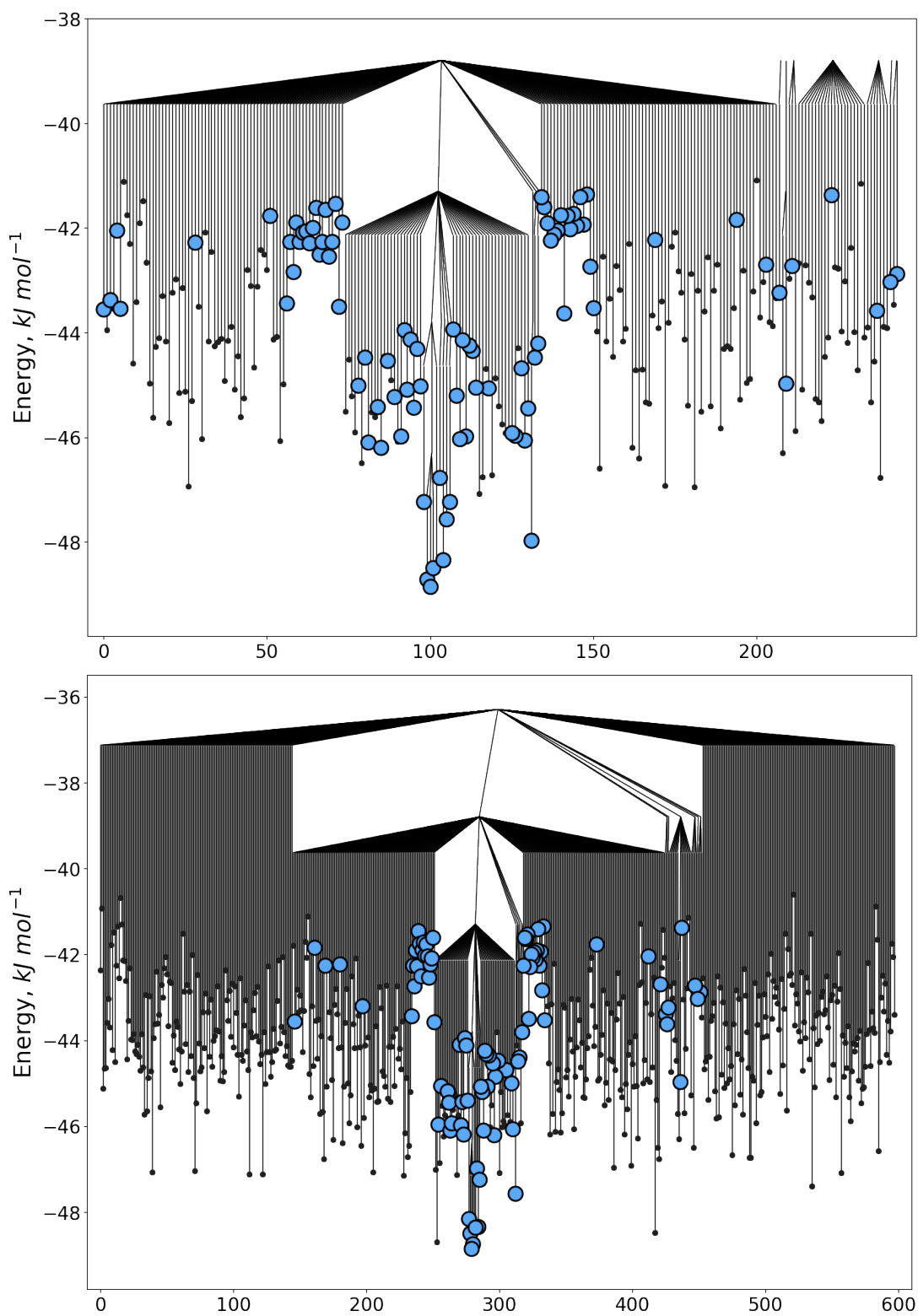


Fig. S3. Threshold clustering results of the lowest 100 structures predicted Benzene structures showing the disconnection graph under energy thresholds of 2.5 kJ mol^{-1} (above) and 5.0 kJ mol^{-1} (below). Initial Structures are coloured blue.

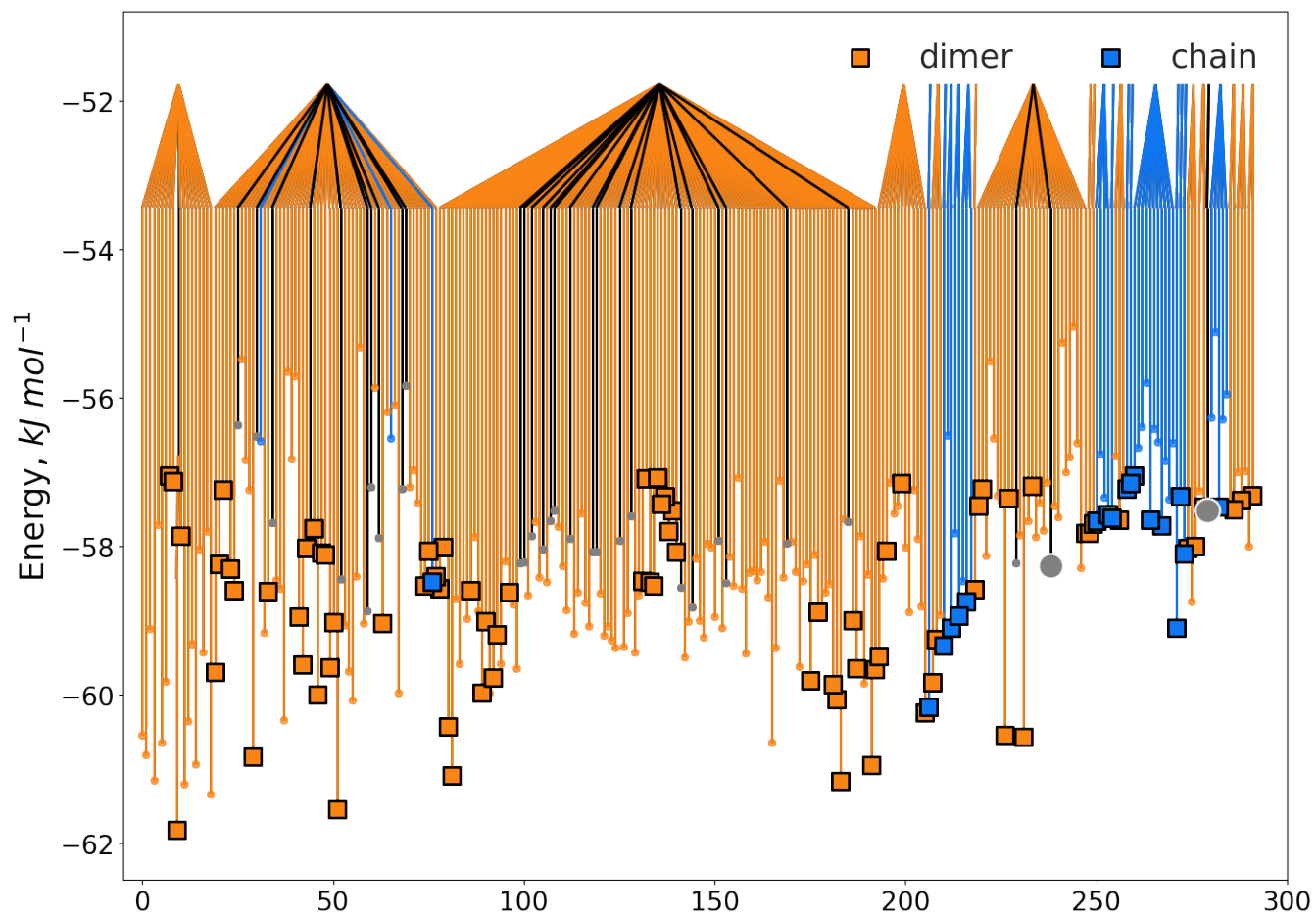


Fig. S4. Disconnection graph of the from the threshold simulations of the 100 lowest energy predicted acrylic acid structures. The structures are coloured by their corresponding hydrogen bonding motif. Structures that do not correspond to purely dimer or chain motifs are coloured grey.

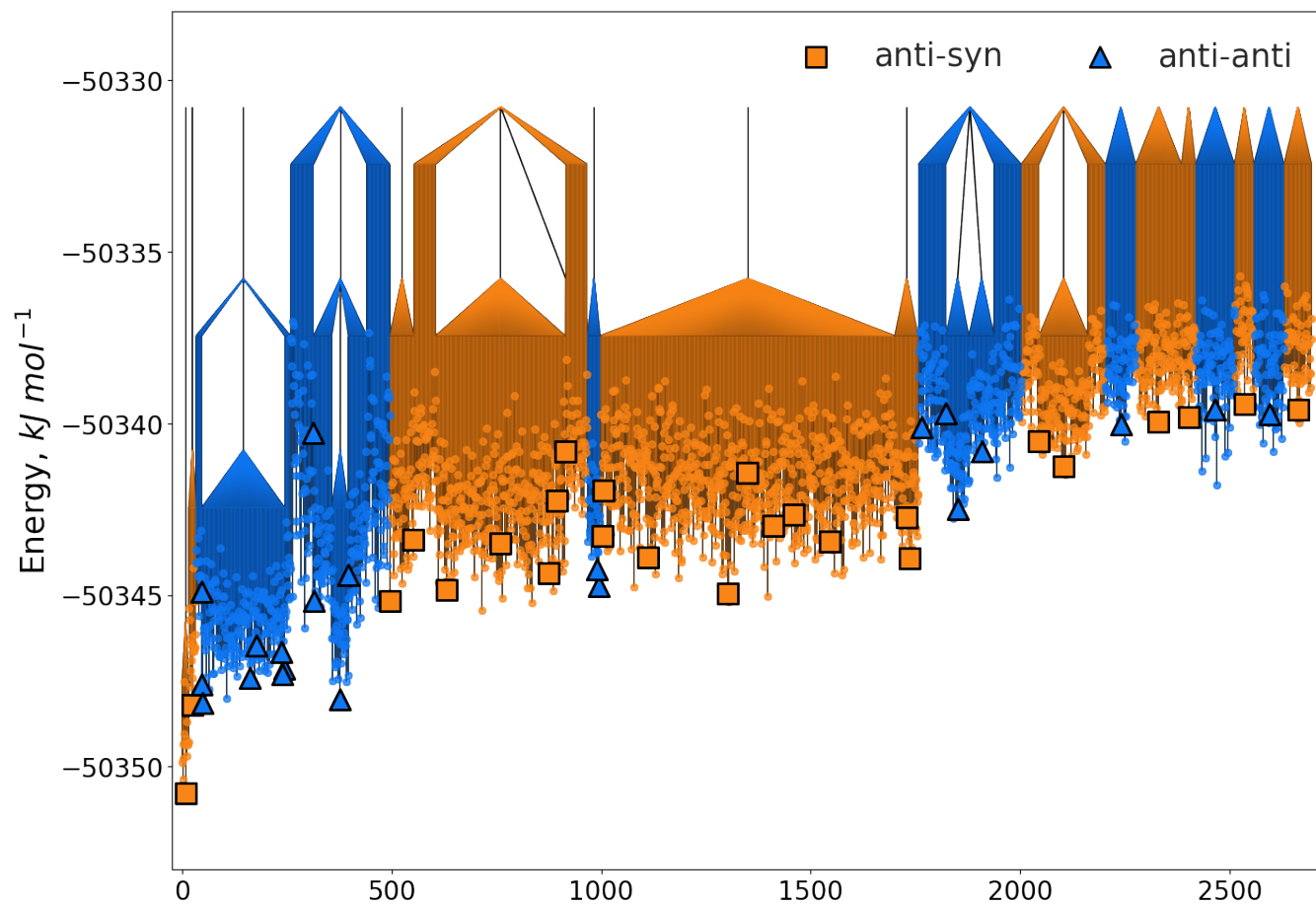


Fig. S5. Disconnection graph of the from the threshold simulations of the 100 lowest energy predicted acrylic acid structures. The structures are coloured depending whether closer to the anti-anti conformation or the anti-syn conformation.

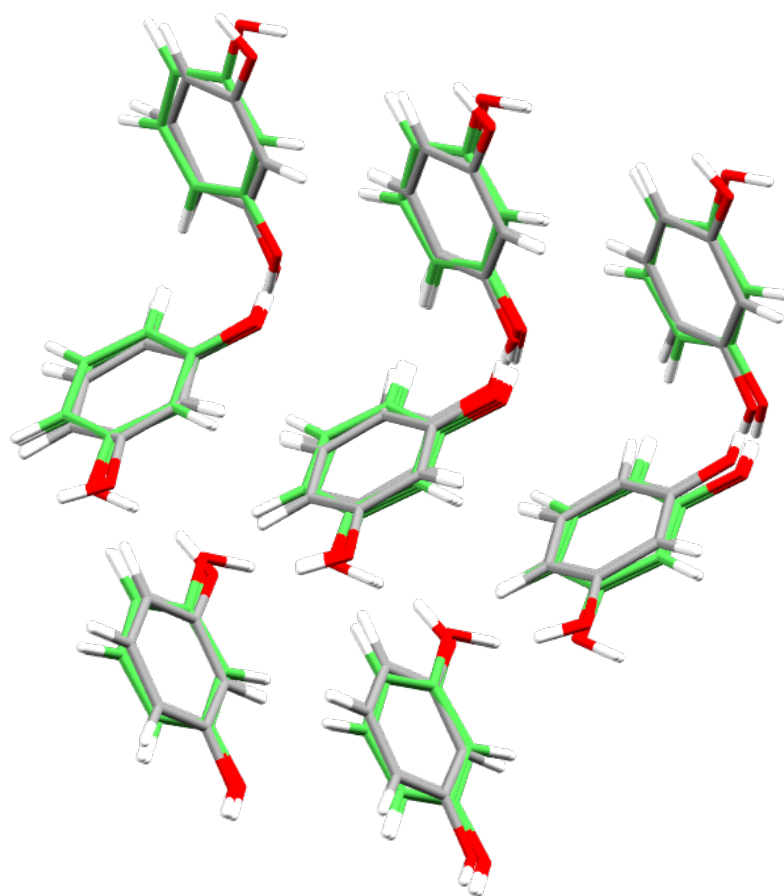


Fig. S6. Crystal structure overlay of the second and fourth lowest energy predicted resorcinol structures which have the same packing but different conformations. The conformation determines the direction of a hydrogen bonding chain.

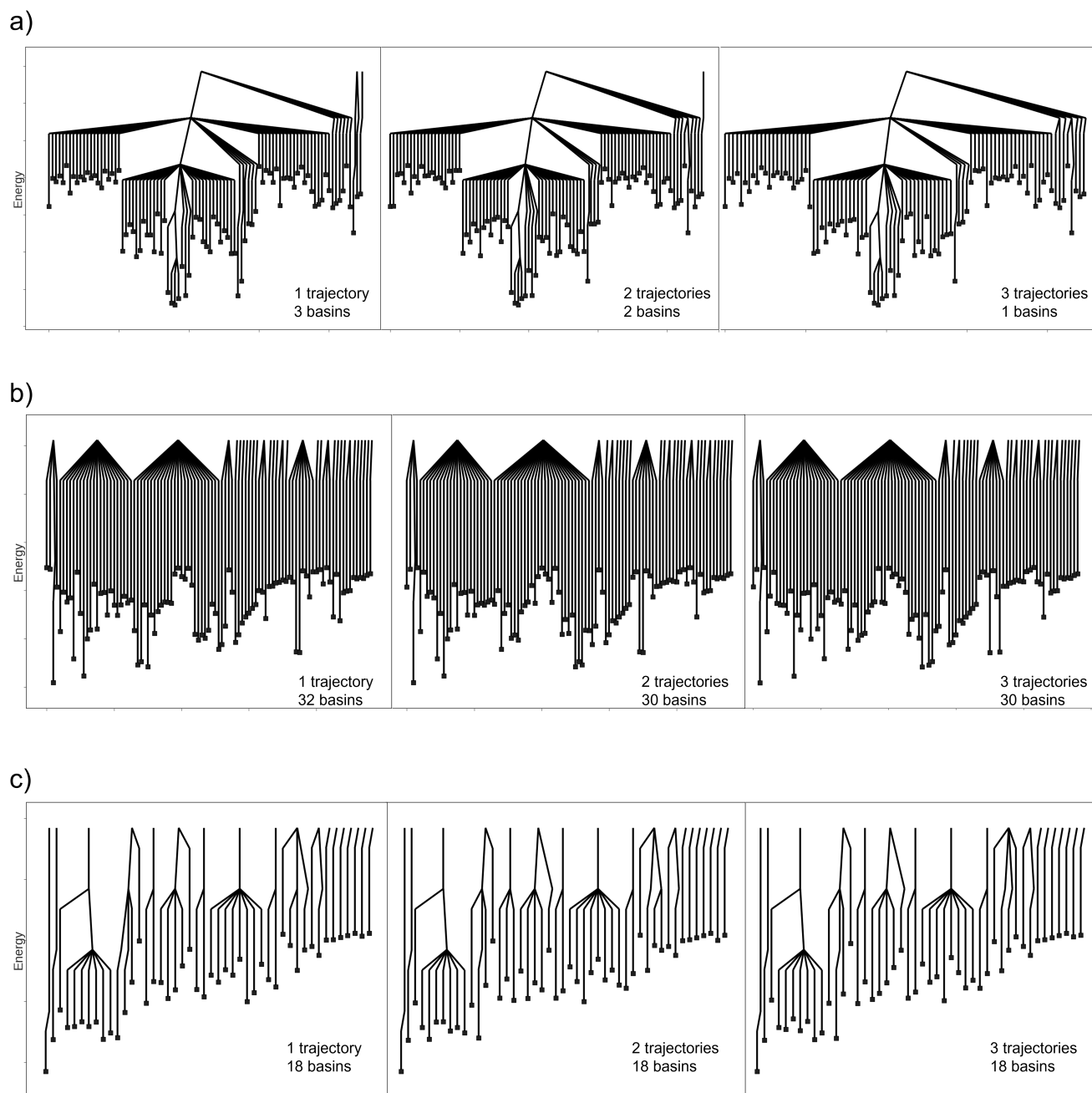


Fig. S7. Convergence of the connections between the initial structures for benzene (a), acrylic acid (b), and resorcinol (c) with increasing number of trajectories initiated per structure: one (left), two (middle), and three (right). MC Structures that were not part of the initial CSP have been omitted from the disconnectivity graphs for clarity.

Table S1. Breakdown of the cost in CPU hours to apply threshold clustering to each of the CSP landscapes. Considering the convergence presented above, we have provided an estimated total cost for running a single trajectory per initial structure.

	Benzene	Acrylic Acid	Resorcinol
Monte Carlo Simulations	2877	1905	2052
Lattice Energy Minimisations	14049	9111	47100
Clustering	272	309	893
Total	17198	11325	50045
Estimate for single trajectory	5732	3775	16681

91 **References**

- 92 1. DH Case, JE Campbell, PJ Bygrave, GM Day, Convergence Properties of Crystal Structure Prediction by Quasi-Random
93 Sampling. *J. Chem. Theory Comput.* **12**, 910–924 (2016).
- 94 2. DE Williams, SR Cox, Nonbonded potentials for azahydrocarbons: The importance of the Coulombic interaction. *Acta*
95 *Cryst B* **40**, 404–417 (1984).
- 96 3. MJ Frisch, et al., Gaussian09 Revision D.01 (Gaussian Inc.) (2013).
- 97 4. AJ Stone, Distributed Multipole Analysis: Stability for Large Basis Sets. *J. Chem. Theory Comput.* **1**, 1128–1132 (2005).
- 98 5. GG Ferenczy, Charges derived from distributed multipole series. *J. Comput. Chem.* **12**, 913–917 (1991).
- 99 6. JR Holden, Z Du, HL Ammon, Prediction of possible crystal structures for C-, H-, N-, O-, and F-containing organic
100 compounds. *J. Comput. Chem.* **14**, 422–437 (1993).
- 101 7. SL Price, et al., Modelling organic crystal structures using distributed multipole and polarizability-based model inter-
102 molecular potentials. *Phys. Chem. Chem. Phys.* **12**, 8478–8490 (2010).
- 103 8. AL Spek, Single-crystal structure validation with the program PLATON. *J Appl Cryst* **36**, 7–13 (2003).
- 104 9. CR Groom, IJ Bruno, MP Lightfoot, SC Ward, The Cambridge Structural Database. *Acta Cryst B* **72**, 171–179 (2016).
- 105 10. B Hourahine, et al., DFTB+, a software package for efficient approximate density functional theory based atomistic
106 simulations. *J. Chem. Phys.* **152**, 124101 (2020).
- 107 11. JG Brandenburg, S Grimme, Accurate Modeling of Organic Molecular Crystals by Dispersion-Corrected Density Functional
108 Tight Binding (DFTB). *J. Phys. Chem. Lett.* **5**, 1785–1789 (2014).