**Ultra-fast detection of (near-)duplicate structures across major crystal databases Daniel Widdowson, Vitaliy Kurlin**. Computer Science, *Materials Innovation Factory*, Liverpool

**Problem**: the incomplete definition of `*isostructural crystals*' allowed anyone to claim a new material by **disguising a known crystal** via almost any perturbation discontinuously changing a reduced cell [1]. Crystal structures are determined in a *rigid form* and hence are indistinguishable under *rigid motion*. *Rigid motion* = translations + rotations, *isometry* = any rigid motion + reflection. *New definition* [1] : a *periodic structure* is a *class* of periodic sets of atomic centres under *rigid motion* or (weaker) *isometry* 

Solution : the invariant descriptor Pointwise Distance Distribution

PDD(S;k) = matrix of weighted rows of distances from an atom of S to k nearest neighbours in the full S, invertible to any generic crystal S.

*Earth Mover's Distance* (**EMD**) on PDDs satisfies all metric axioms.



**Scale of duplication** : Google's GNoME claimed "2.2M new crystals – equivalent to nearly 800 years' worth of knowledge" [2], made 384K+ CIFs public, of which 43 were synthesised by Berkeley's A-lab [3]

Review [4] of [3] : "none of the 43 materials [of 58 attempted] produced by A-lab were new: the large majority were misclassified, and a smaller number were correctly identified but already known". Review [5] of [3] : all 43 materials have near-duplicates in ICSD or Materials Project used for training. Review [6] of [2] : "scant evidence for compounds that fulfill the trifecta of novelty, credibility, utility". Review [1] of [2] : 4 identical CIFs, 43 triples, 1089 pairs, many more thousands of digital duplicates.



[1] O.Anosova, V.Kurlin, M.Senechal. The importance of definitions in crystallography. *IUCrJ*, v.11 (4), p.453-463, 2024.

[2] A.Merchant et al. Scaling deep learning for materials discovery. Nature 624 (7990), 80-85, November 2023.

[3] N.Szymanski et al. An autonomous laboratory for the accelerated synthesis of novel materials. Nature 624 (7990), 86–91.

[4] J.Leeman et al. Challenges in High-Throughput Inorganic Materials Prediction. PRX Energy 3, 011002, March 2024.

[5] **D.Widdowson, V.Kurlin**. Navigation maps of the materials space for automated self-driving labs, *arxiv:2410.13796*.

[6] A.Cheetham, R.Seshadri. Artificial Intelligence Driving Materials Discovery? Chemistry of Materials 36, 3490–3495, 2024.

[7] **D.Widdowson et al**. Average Minimum Distances of periodic point sets. *MATCH*, v.87 (3), 529-559, 2022.

[8] D.Widdowson, V.Kurlin. Resolving the data ambiguity for periodic crystals. *NeurIPS*, v.35, 24625-24638, 2022.

[9] D.Widdowson, V.Kurlin. Continuous invariant-based maps of the CSD. Crystal Growth & Design, v.24, 5627–5636, 2024.

## Geometric Data Science (GDS) develops continuous maps of data objects

**The vision** is to map (continuously parametrize) the space of any data objects considered up to practical equivalences. While Geometric Deep Learning experimentally outputs equivariant descriptors of clouds or graphs, GDS developed analytic, complete and continuous invariants for any finite and generic periodic sets of unordered points in  $\mathbb{R}^n$ , see the papers in NeurIPS 2022 and CVPR 2023 at http://kurlin.org/research-papers.php#Geometric-Data-Science.

**The key obstacle** for periodic crystals was the ambiguity of conventional data based on minimal or reduced cells that are discontinuous under atomic displacements. Without continuously quantifying the crystal similarity, the brute-force Crystal Structure Prediction produces millions of nearly identical approximations to numerous local energy minima, see red peaks in Fig. 1.



Figure 1: Left: energy landscapes show crystals as isolated peaks of height= -energy. To see beyond the 'fog', we need a map parametrized by invariant coordinates with a continuous metric. **Right**: R. Feynman's first lecture showed that 7 cubic crystals differ by side lengths, while our invariants distinguished all 850K+ periodic crystals in the CSD. These crystals have unique positions in a common *Crystal Isometry Space* whose one 2D projection is in Fig. 2.



Figure 2: Carbon allotropes on a continuous map of periodic crystals from the Cambridge Structural Database (CSD), Crystallography Open Database (COD), Inorganic Crystal Structure Database (ICSD), and Materials Project (MP). The original invariant  $AMD_k$  is the average distance to the *k*-th atomic neighbor whose asymptotic explains the density. The color indicates the number of crystals whose invariants  $ADA_k$  (average deviation from asymptotic) are discretized to each pixel. Any hot spot can be visualized in many other explicit coordinates.