
Extracting Structural Motifs from Pair Distribution Function Data of Nanostructures using Explainable Machine Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Characterization of material structure with X-ray or neutron scattering using e.g.
2 Pair Distribution Function (PDF) analysis most often rely on refining a structure
3 model against an experimental dataset. However, identifying a suitable model is
4 often a bottleneck. Recently, new automated approaches have made it possible to
5 test thousands of models for each dataset, but these methods are computationally
6 expensive, and analysing the output, i.e., extracting structural information from the
7 resulting fits in a meaningful way is challenging. Our Machine Learning based
8 Motif Extractor (ML-MotEx) trains an ML algorithm on thousands of fits, and uses
9 SHAP (SHapley Additive exPlanation) values to identify which model features are
10 important for the fit quality. We use the method for 4 different chemical systems
11 including disordered nanomaterials and clusters. ML-MotEx opens for a new type
12 of modelling where each feature in a model is assigned an importance value for the
13 fit quality based on explainable ML.

14 1 Introduction

15 The development of advanced, functional materials builds on an understanding of the intricate relation-
16 ship between material structure and properties, and over the past century, crystallographic methods
17 using scattering and diffraction have thus been essential for materials science. Crystallography allows
18 *ab initio* determination of crystal structures from diffraction data, and has provided us with the vast
19 knowledge of crystal chemistry that is now used in design of functional materials. However, in the
20 case of nanomaterials with limited long-range order, crystallographic methods are challenged, and *ab*
21 *initio* structure determination, or structure solution, is not currently possible. Over the past decades,
22 total scattering with Pair Distribution Function (PDF) analysis has become an essential tool for
23 characterisation of nanomaterial structure.[1, 2] The PDF is the Fourier transform of normalized and
24 corrected X-ray, neutron, or electron scattering intensities, and is a function in real space representing
25 a histogram of interatomic distances in the sample. Compared to crystallographic methods relying on
26 long-range order, PDF analysis can be applied for nanomaterials,[3-5] disordered[1, 6, 7] or amor-
27 phous materials.[3, 5, 8] However, structure solution from the PDF is not possible except in a very
28 few simple cases,[9] using either the Reverse Monte Carlo method[10] or the LIGA algorithm.[11,
29 12] In the absence of broadly applicable *ab initio* nanostructure determination methods, it is therefore
30 necessary to propose reasonable starting models and to then ‘refine’ the model parameters against
31 the data using local minimization methods. The step of finding a starting model can be a major
32 challenge and is thus a bottleneck in complex material characterization. In the case of PDF analysis
33 of nanomaterials, such models are often guessed at by considering related bulk materials, however
34 these are often not good starting models for very small clusters and nanoparticles, where significant
35 structural changes may take place.[3, 5, 13, 14] A way of building plausible starting models is thus

36 needed, where structure models can be built capturing local bonding topologies suggested by known
37 chemistries.

38 Recently, automated methods such as ‘structure mining’ and ‘cluster mining’ have appeared in the
39 literature to help overcome this challenge.[15-17] In a study of the structure of metallic nanoparticles,
40 Banerjee et al. automatically generated thousands of discrete metal nanocluster structures and fitted
41 PDFs from each of them to experimental data to identify the best model in an automated manner.[17]
42 In a recent study of molybdenum oxide nanomaterials, a new approach were introduced, where a large
43 number of MoO_x cluster structure models were automatically generated and compared their PDFs to
44 experimental data in order to identify dominating structural motifs in the sample, i.e. arrangements of
45 atoms that dominate the material structure on the local scale.[7] The authors hypothesised that the
46 structural motifs present in amorphous molybdenum oxides can also be found in crystalline structures,
47 and therefore used crystal structures of molybdenum oxides as starting models. From these models,
48 they cut out thousands of different cluster structure models of different sizes to build a ‘catalogue’ of
49 structure candidates. These models were all tested against the experimental PDFs to identify the best
50 fitting structural motif. In another study, a similar approach were used for identification of a bismuth
51 oxido cluster intermediate structure in a study of cluster growth.[18]

52 While these approaches can extend the structural space searched when identifying models for structure
53 refinement, new challenges arise. Firstly, the refinement processes can be computationally heavy,
54 which can limit the number of catalogue structures that are tested. For example, our brute force
55 approach for cluster identification above generates $2^N - 1$ structures for starting model sizes with N
56 atoms. Each structure must have its PDF computed and then refined against the target measured PDF,
57 so that its fit quality can be evaluated. This process is computationally costly and does not scale well
58 with number of structure candidates. Furthermore, for disordered, amorphous, and nanostructured
59 systems many hundred models may provide similar fit qualities, and if only reporting a few of them,
60 it is difficult to assess which structural features of these models are important. We therefore need
61 effective and unbiased methods to compare many fits to extract structural information. Here, we
62 introduce a completely new approach that uses an explainable Machine Learning (ML) model that,
63 after training, will predict the agreement factor for a test cluster with a given dataset. Furthermore,
64 the use of explainable ML informs which features in the model are important for the agreement
65 factor.[19-24] Our Machine Learning based Motif Extractor (ML-MotEx) model is illustrated in
66 Figure 1. Firstly, it builds a large catalogue of thousands of candidate structural motifs, which are
67 ‘cut outs’ from a chosen bulk structure[7, 18] (step 1). The PDF is then computed from each one, and
68 each model is fit to the target dataset (step 2). The structures and R_{wp} values from each fit are handed
69 to an ML algorithm applying gradient boosting decision trees (GBDTs),[25] which learns to predict
70 R_{wp} values for new fits based on an atomic structure model (step 3). The ML-MotEx algorithm then
71 outputs quantified values of how important each atom or feature in the starting structure is for the
72 fit to yield a low R_{wp} value with the given fitting-algorithm (step 4). This is done by using SHAP
73 (Shapley Additive exPlanation)[26, 27] values, which is a known method for explaining tree-based
74 ML models. The amplitude of the SHAP value reflects how important a structural feature is for the fit
75 quality, while the sign of the SHAP value reflects whether the feature affects the R_{wp} value of the fit
76 towards 1 (poor fit) or 0 (perfect fit), in other words why it is important.

77 Compared to the automated, brute-force methods previously introduced for PDF analysis,[7, 15-17]
78 we can much faster screen a larger number of structures. Our method only needs to screen a sub-
79 sample ($\approx 10,000$) of the much larger number of motifs that can be generated from a bulk material to
80 learn how to predict which structures provide a good agreement with the data. The analysis done for
81 the examples presented below would take ≈ 24 days for starting models with 24 atoms, $\approx 3 \cdot 10^6$
82 years for starting models with 48 atoms and $\approx 6 \cdot 10^{13}$ years for starting models with 72 atoms
83 using a brute-force approach (section A in the SI), while ML-MotEx analysis is done in minutes or
84 hours. Furthermore, the use of explainable ML provides a way to better analyse the output of the
85 screening: instead of just identifying the model that provides the lowest R_{wp} value, we are able to
86 output a measure of how important each atom or feature (e.g. size or shape) in the starting model is
87 for the fit to yield a low R_{wp} value (step 4). This procedure is automated, can be done in quasi-real
88 experimental time and without human bias.

89 We illustrate the use of ML-MotEx using 4 different examples. We first show the principles of the
90 method using a simple model system based on simulated X-ray PDF data from a C_{60} buckyball.
91 We further demonstrate the use of ML-MotEx on experimental X-ray PDF data from amorphous,
92 disordered molybdenum oxides[7] and tungstate α -Keggin clusters in solution,[28] where it allows

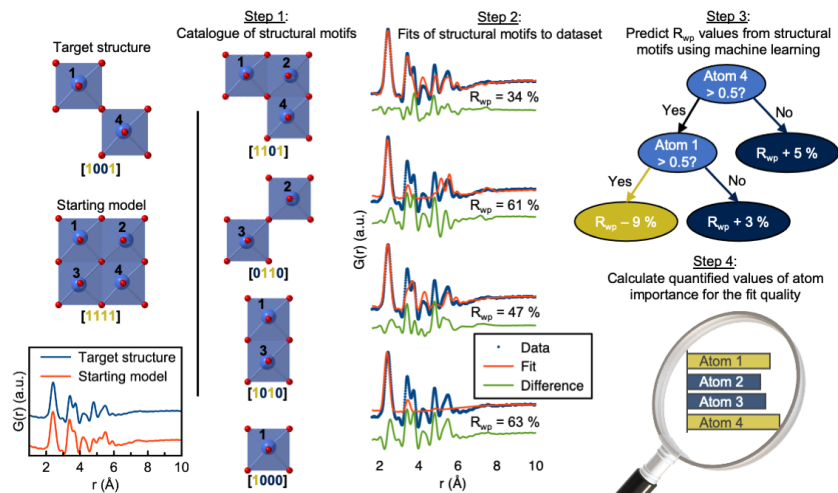


Figure 1: Illustration of the ML-MotEx process. Firstly, a starting model is provided. Using this starting model, a structure catalogue is generated, and the structures in the catalogue are fitted to the experimental data in question. An ML algorithm is then trained to predict R_{wp} values and finally calculating quantified values of feature importance for the fit quality.

93 identifying the main structural motifs present in the samples using different starting models. Lastly,
 94 we extend the method to use a ‘cookie-cutter’ strategy to generate structures for the catalogue of
 95 candidate motifs. Here, the algorithm is used to identify a bismuth oxido cluster by using a cut-out
 96 of the β - Bi_2O_3 structure as starting model. The examples illustrate that it is possible to obtain
 97 knowledge of dominating structural motifs from PDF in an automated manner using ML.

98 2 Results

99 2.1 ML-MotEx algorithm

100 ML-MotEx consists of four steps. These four steps are shown in Figure 1. In the first step, a starting
 101 structure model is used to generate a catalogue of candidate structure motifs. As detailed in the
 102 Methods section, the structures are generated by removing different numbers of atoms from the
 103 original starting structure which results in thousands of smaller, candidate structure motifs. In the
 104 second step, a fitting script is used to fit the generated candidate structures to the dataset. In the third
 105 step, the fitting results are handed to the explainable ML algorithm which is optimised and trained.
 106 By using this information, SHAP values of the atoms or structural features in the starting model are
 107 calculated in the fourth step. The output of the algorithm is thus the starting model along with SHAP
 108 values, indicating the importance of each individual atom in the structure for the fit quality, or in other
 109 words; how much each individual atom or feature affects the R_{wp} value either positively or negatively.
 110 We refer to this value as the “atom contribution value”. We furthermore define the ratio between the
 111 atom contribution value and its uncertainty as the “confidence factor”. Further definitions and
 112 descriptions of the individual steps of the algorithm are given in the Methods section.

113 2.2 Example 1: Proof-of-concept: Identification of the C_{60} buckyball

114 We first show the use of ML-MotEx with a simple, proof-of-concept example, using a calculated
 115 PDF from an ideal C_{60} buckyball (Figure 2A). The aim is to identify the structural motif, the C_{60}
 116 buckyball, from the data. We first need a starting structure that contains the motifs we are looking for.
 117 In this simplified example, we use a single unit cell of the crystal structure of C_{60} .^[29] However, we
 118 discarded all symmetry and generated a discrete structure model corresponding to the 132 atoms in
 119 one unit cell. This model is shown in Figure 2B, where one whole C_{60} structure (Figure 2A) is seen
 120 along with fragments of the neighbouring C_{60} buckyballs. The simulated PDF of the C_{60} buckyball
 121 and the starting model are shown in Figure 2C. We can now use this starting model to generate
 122 a catalogue of structures, which are all fitted to the data. The structures are created by removing

123 different numbers of atoms from the original starting structure, which results in thousands of smaller,
124 candidate structure motifs. This model generation and fitting steps are identical to our previously
125 reported brute-force approach, where we simply compare the R_{wp} values of all the fits to identify the
126 best structure motif. We first consider this simple approach. One of the limitations of the brute-force
127 method is that the possible candidate structures is exponential in N , the number of atoms in the model.
128 Since each atom in the starting model can be present or absent, the number of possible sub-clusters
129 is equal to $2^N - 1$. For large models such as the C_{60} starting model containing 132 atoms, this is
130 $\approx 10^{40}$, a gigantic number, making it impossible to investigate all candidate structures. For this
131 example, we used 384,260 structures to train ML-MotEx, which is only a very small fraction of the
132 $2^{132} - 1$ possible candidate structures. Note that the model with a single C_{60} buckyball was not in
133 the generated structure catalogue. All these 384,260 structures were fitted to the PDF calculated from
134 the C_{60} cluster. Only a scale factor, an isotropic expansion/contraction factor, and isotropic Atomic
135 Displacement Parameters (ADPs) were refined, as detailed in the Methods section. We note that
136 refinement of the atom positions can be added to the fitting procedure to expand the chemical space
137 that is investigated. However, this would be computationally expensive and it would allow deviations
138 from the chemical topologies set up in the starting model.

139 To get an overview of the results from these fits, we plot the R_{wp} value versus the number of atoms in
140 the structure, Figure 2D. To further investigate the results, one must visually inspect the fits of the
141 catalogue of candidate structure motifs and their R_{wp} value. Some of the candidate structure motifs
142 are shown as inserts in Figure 2E, where transparent grey atoms represent atoms deleted from the
143 models. The fits of these structures to the dataset are presented in Figure 2E, along with the R_{wp}
144 values. The R_{wp} value appears to drop when the ‘outer’ atoms are removed, while it increases when
145 the atoms that are part of the center C_{60} buckyball are removed. From investigating these few, but
146 manually selected, structures and their corresponding fitted R_{wp} value, one can hypothesize that
147 the structure giving the best fit should be the C_{60} buckyball. However, this method can be biased
148 by human interaction, and it is time-consuming and difficult to go through the many fits to extract
149 structural information. We therefore move on to the ML-MotEx method. Using the catalogue of
150 candidate structure motifs and the corresponding R_{wp} values obtained above, we train a GBDT model
151 on the training set to predict the R_{wp} value of the candidate structure motifs. Figure 2F shows the
152 predicted R_{wp} values of the ML algorithm versus the R_{wp} value of the structures when they are fitted
153 to the simulated C_{60} dataset in DiffPy-CMI.[30] For the structures used in the test set, the GBDT
154 model predicts the R_{wp} value with a mean absolute error of 2.0 %. We now use explainable ML to
155 explain R_{wp} values with the use of the feature importance tool SHAP values.[27] As described in
156 detail in the Methods section, a SHAP value is calculated for each structural feature (here each atom
157 and the cluster size) for each candidate structure motif that is fitted to the PDF during the training
158 process. The amplitude of the SHAP value reflects how important a structural feature is for the fit
159 quality, while the sign of the SHAP value reflects whether the feature affects the R_{wp} value of the fit
160 towards 1 (poor fit) or 0 (perfect fit), in other words why it is important.

161 Figure 3A shows the most important results from the SHAP value analysis. The first feature we
162 consider is the number of atoms, with SHAP values shown in the top part of Figure 3A. The plot
163 represents SHAP values for the cluster size feature with the size shown on a colour scale, going from
164 small (blue) to large clusters (red). From the large amplitude of some of the SHAP values observed
165 from this feature, we see that the number of atoms in the structure motif is the most important feature
166 for the R_{wp} value. All small clusters (0–34 atoms, plotted in blue colours) show a large positive
167 SHAP value, which implies that the R_{wp} value of the fit to the PDF data is high, i.e. the fit quality is
168 low. All small clusters can thereby be discarded as structural models for satisfyingly describing the
169 data. Next, we can investigate the SHAP values obtained for the individual atoms in the structure. We
170 first consider atom 13, as labelled in the structure drawing in Figure 3B. The SHAP values obtained
171 from this atom for each of the fits in the training set are all plotted on the SHAP axis. For the models
172 where the atom is not present in the model, the SHAP value is shown in blue, while it is shown in
173 red for the atoms where it is present in the model. If first considering the cases where the atom is
174 kept in the model, the atom 13 SHAP values are generally negative, which means that the presence
175 of this atom pushes the R_{wp} value towards 0. We interpret this as ML-MotEx wants to keep the
176 atom in the model. The SHAP values obtained for the fits without the atom present are positive,
177 which confirms that if removing the atom, the fit quality becomes worse. Based on the SHAP values
178 obtained for the atom in each fit, we calculate an atom contribution value. The atom contribution
179 value is defined in the Methods section, and is calculated as the difference between the average SHAP
180 values obtained for the atom when kept in the model, and when removed from the model. A negative

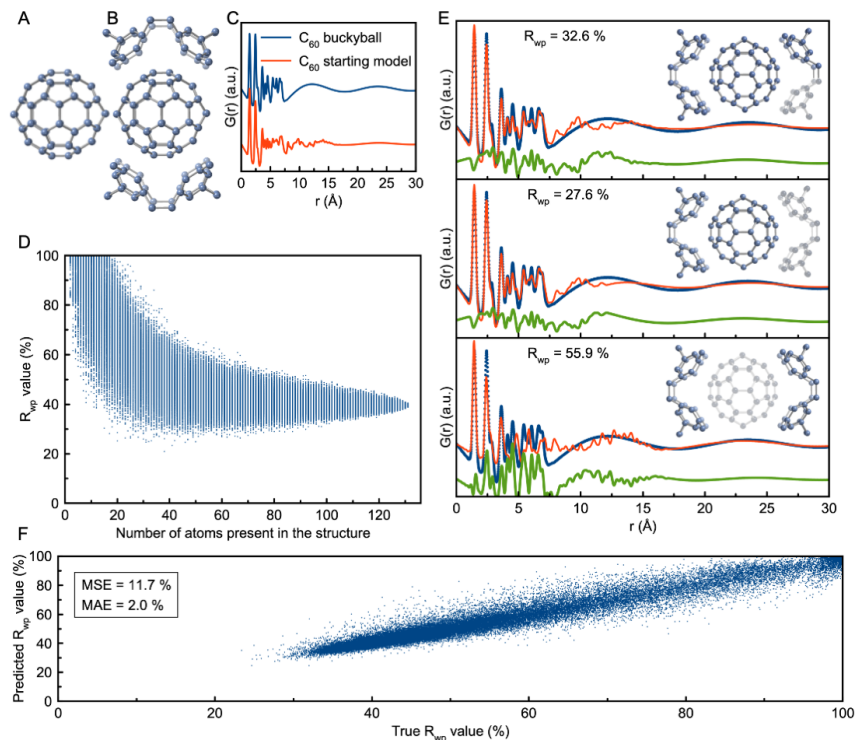


Figure 2: A) C_{60} buckyball, B) single C_{60} unit cell,[29] treated as a discrete structure with 132 atoms and C) their simulated PDFs. The simulation parameters (presented in section B in the SI) mimic typical values of a PDF dataset. D) R_{wp} values obtained in the fits using the C_{60} structure catalogue, plotted as a function of number of atoms in the structure motifs. Note that the model with a single C_{60} buckyball is not included in the set of 384,260 structures tested. This would result in a perfect fit with $R_{wp} = 0\%$. E) Examples of candidate structure motifs with their corresponding fits to the simulated C_{60} buckyball data. Grey, semitransparent atoms are removed from the starting model. F) Predicted R_{wp} values versus true R_{wp} values. R_{wp} values from the fits of the catalogue structures to the simulated C_{60} dataset, plotted versus the predicted R_{wp} values from the GBDT model from the same structures. The mean squared error (MSE) and the mean absolute error (MAE) are based on all 76,852 predictions in the test set, which are structures the model has not been trained on.

181 atom contribution value means that the atom pushes the R_{wp} value down if kept in the structure. The
 182 atom contribution value obtained for atom 13 is negative, and we therefore colour it yellow in the
 183 structural representation in Figure 3B to indicate that it should be kept in the model. We use this
 184 strategy to automatically go through all the atoms in the starting model and colour them yellow/black
 185 based on their impact on the R_{wp} value. The result can be seen in Figure 3B where the 60 atoms with
 186 the lowest atom contribution values are coloured yellow. The results are also shown in section C in
 187 the SI, where the atom contribution values are plotted using a continuous colour bar. The results
 188 show that ML-MotEx mainly favours the atoms comprising the central buckyball. While the average
 189 confidence factor (as defined in the Methods section) is 1.26 for all of the atoms in the starting
 190 model, we observe that the average confidence factor of the mislabelled atoms is 0.37, meaning that
 191 ML-MotEx is less confident about the atom contribution values of those.

192 The ML-MotEx algorithm thus provides an unbiased method to extract important motifs from PDF
 193 data, without any inputs other than a starting model and a fitting script. We emphasize that the
 194 structural motifs extracted with ML-MotEx are based on the R_{wp} value of the fits and are thereby not
 195 necessarily physically reasonable. It is therefore important to still critically consider the extracted
 196 motif with chemical knowledge, in the same manner as for conventional PDF refinements. In this
 197 process, one could refine additional parameters such as atom positions. Consequently, in Figure 3B,
 198 the user should identify the full C_{60} buckyball as the structural motif rather than just choosing the

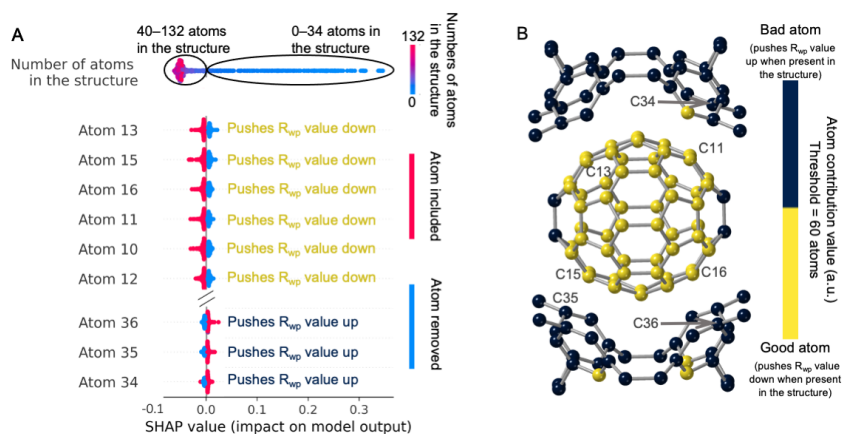


Figure 3: A) Plot of the SHAP values obtained in the C_{60} analysis, showing if atoms in the starting model are favourable for the fit quality. For the models where the atom is not present in the model, the SHAP value is shown in blue, while it is shown in red for the atoms where it is present in the model. The SHAP values are plotted as a violin plot. B) Structural visualisation of kept and removed atoms. The atoms with the 60 lowest atoms contribution values have been coloured yellow, while the rest are coloured black. Section C in the Supporting Information (SI) shows a similar representation using a continuous colorbar for the atom contribution values.

199 motif of the yellow atoms. Another approach to avoid unphysically results from ML-MotEx would
 200 may be to include e.g., density function theory (DFT) calculations in the goodness-of-fit value.

201 2.3 Example 2: Identification of the ionic cluster structure from PDFs

202 To investigate the reproducibility of the ML-MotEx method, we investigate if similar re-
 203 sults are achieved with different starting models, all containing the correct structure motif.
 204 We here model a PDF obtained from a solution of 0.05 M ammonium metatungstate hy-
 205 drate, $(NH_4)_6[H_2W_{12}O_{40}] \cdot H_2O$ in water, which dissolves to form monodisperse α -Keggin
 206 clusters.[28] Experimental details are provided in section D in the SI. To test the ML-
 207 MotEx method we use four different starting models of tungstate oxide crystals, all includ-
 208 ing the α -Keggin cluster motif with varying complexity. Unit cells from the 4 following
 209 crystal structures were used as starting models: $[Hpy]_4H_2[H_2W_{12}O_{40}]$ (py=pyridine) (1),[31]
 210 $(CH_3)_4N)_4SiW_{12}O_{40}$ (2),[32] $((CH_3)_2NH_2)_6(Cu(HCON(CH_3)_2)_4)(GeW_{12}O_{40})_2)(HCON(CH_3)_2)_2$
 211 [33] (3), and $(CH_3)_2NH_2)_3(PW_{12}O_{40})$ (4).[34] Again, we discarded all symmetry and generated
 212 discrete structure model corresponding to the atoms in one single unit cell. All other atoms than
 213 tungsten and oxygen were furthermore removed from the structures before catalogue structures were
 214 created. Figure 4A shows the experimental dataset with simulated PDFs from the 4 different starting
 215 models. Figure 4B illustrates a $W_{12}O_{40}$ α -Keggin structure.

216 Again, we first build structure catalogues based on the starting models (step 1) and fit them to the
 217 experimental PDF (step 2). In this case, we extract 104 structures from each starting model, which is
 218 just a small fraction of all possible structures that can be made from the starting models that have
 219 24 (2), 48 (1) and (3), and 72 (4) atoms that are permuted. Again, a GBDT model was trained to
 220 predict the R_{wp} values of the structures (step 3), and SHAP values were obtained to calculate atom
 221 contribution values (step 4). The resulting SHAP value plots can be seen in section D in the SI.
 222 While ML-MotEx takes about 100 seconds on an AMD Ryzen Threadripper 3990X with 64-core
 223 2.9/4.3GHz using 104 fits on a structure with 48 atoms, it would take about $\approx 3 \cdot 10^6$ years (section
 224 A in the SI) to make fits of all the 2^{48} – 1 possible structures using the brute-force approach. Table
 225 S1 in the SI shows the exact computer time of the fits on a MacBook Pro and a Threadripper, which
 226 clearly demonstrates the scalability of ML-MotEx.

227 Figure 4C-F shows the results of applying ML-MotEx to the 4 different starting models. For structures
 228 (1), (3), and (4), the 24 atoms most preferred by ML-MotEx were coloured yellow, while the rest
 229 were coloured black. For structure (2), 12 atoms were coloured yellow. In all 4 examples, the yellow
 230 atoms have a motif of a α -Keggin cluster, however, in Figure 4E–F, we see a few mislabelled atoms

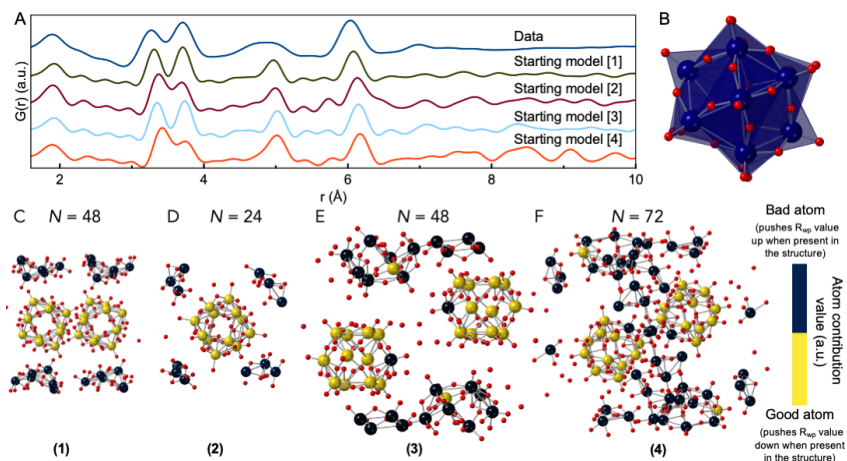


Figure 4: A) Comparison of experimental data from a 0.05 M ammonium metatungstate hydrate solution, and simulated PDFs from the four different starting models (1)–(4). The simulation parameters mimic typical values of a PDF dataset and can be seen in section B in the SI. B) The $W_{12}O_{40}$ α -Keggin structure. C-F) Results from ML-MotEx on a PDF from a solution of ammonium metatungstate hydrate, using four different starting models: C) $[Hpy]_4H_2[H_2W_{12}O_{40}]$ (py=pyridine),[31] D) $(CH_3)_4N_4SIW_{12}O_{40}$,[32] E) $((CH_3)_2NH_2)_6(Cu(HCON(CH_3)_2)_4)(GeW_{12}O_{40})_2(HCON(CH_3)_2)_2$,[33] F) $(CH_3)_2NH_2)_3(PW_{12}O_{40})$.[34] Atoms kept by ML-MotEx are shown in yellow while removed atoms are shown in black. The kept atoms were chosen as the 24 atoms (1), 12 atoms (2), 24 atoms (3), and 24 atoms (4) with the lowest atom contribution values. In section D in the SI a similar representation is shown using a continuous colorbar for the atom contribution values.

231 (2 of 24 atoms in the worst case). The mislabelled atoms are found in the starting models containing
 232 most atoms, i.e. with the highest permutation value N . To achieve a better prediction, we could have
 233 built larger catalogues of candidate structure motifs and thus performed more fits. We therefore
 234 conclude that the ML-MotEx method is not completely insensitive to the starting model, but that it
 235 yields very similar results for all the tested starting models if it contains similar motifs. Furthermore,
 236 the example shows that ML-MotEx can be used to investigate PDF data from clusters in solution,
 237 whose structure also is part of known crystal structures. As described in section E in the SI, we
 238 performed an identical analysis of a different dataset also obtained from a second solution of 0.05 M
 239 ammonium metatungstate hydrate. This analysis provided highly comparable results. This illustrates
 240 the reproducibility of the method. In section F in the SI, we discuss what happens if a poor starting
 241 model is used, and how one can identify if the starting model does contain the right motif using the
 242 confidence factor. In the SI, we describe two other examples where we have used ML-MotEx. Firstly,
 243 we have used the ML-MotEx method to identify the main structural motifs present in an amorphous,
 244 disordered molybdenum oxide[7] from its experimental X-ray PDF. This example is described in
 245 Section G of the SI. Secondly, we have identified a larger ionic cluster, namely $[Bi_{38}O_{45}]$, from an
 246 experimental PDF. Here, we use the β - Bi_2O_3 structure as starting model, and used a ‘cookie-cutter’
 247 strategy to generate structures for the motif catalogue. This example, and the ‘cookie-cutter’ approach,
 248 are described further in Section I of the SI.

249 3 Discussion

250 In the 4 examples presented above, we have shown how explainable ML can aid in identifying struc-
 251 tural motifs in nanostructured materials and presented a new approach to structure characterization.
 252 Traditional PDF analysis investigates how an entire structure model agrees with an experimental PDF,
 253 rather than identifying how different features in the model affect the fit quality. Instead, ML-MotEx
 254 provides a quantitative measure of how each atom or feature contributes to the fit. The use of ML
 255 furthermore allows screening of a large number of models in an automated and fast manner. In the
 256 examples described here, ML-MotEx has been used with various starting models with up to 256 metal
 257 atoms, however, the algorithm can handle larger systems, as it is highly scalable. In comparison, a

258 full brute-force approach is computationally restricted to systems with up to 15–30 atoms. For the
259 type of systems described here, it is possible to use the method in quasi-experimental time which
260 could, for example, be useful for analysis of time-resolved scattering data, where the structural motifs
261 present might change with time, which would be revealed by changing SHAP values.

262 ML-MotEx shares some similarities with the cluster build-up algorithm LIGA,[11, 12] which au-
263 tomatically builds clusters of different sizes based on information that is contained in inter-atomic
264 distance lists extracted from the PDF. LIGA has shown to be successful at automatically reconstruct-
265 ing clusters (up to 150 atoms) with no user input except the interatomic distance list, extracted from
266 an experimental PDF, and at low computational cost. However, its use has not caught on because
267 extracting the distance list from the data presents significant practical difficulties, and is not unique.
268 As with ML-MotEx it uses the error each atom in a cluster contributes to the fit to weight the decision
269 about which atom to include in the model. Presumably, part of the success of LIGA and ML-MotEx is
270 its use of this atom contribution for rapidly finding good candidate motifs. Unlike LIGA, ML-MotEx
271 requires a starting model that contains the target structural motif, and it leverages ML to rapidly
272 compute the atom contributions. It can therefore be positioned between traditional refinement (where
273 the complete starting model is needed) and LIGA (which is *ab initio*) as it finds structural motifs from
274 within a larger model as a starting model for a subsequent refinement. However, it has the significant
275 advantage over LIGA that it works directly on the measured PDF and does not require the inter-atomic
276 distance list to be extracted from the PDF data and we expect it to be of great practical value. It may
277 be considered as a significant drawback that ML-MotEx requires as an input a structure fragment that
278 contains the target motif within it in order to work. We provide a confidence factor for the starting
279 model but ML-MotEx still requires significant chemical/structural knowledge and intuition to be of
280 use. We first note that such intuition is widespread in the chemistry community and is unlikely to be
281 a significant drawback in practice. For example, ML-MotEx has recently been used to identify the
282 structure of intermediates in the formation of transition metal tungstates from polyoxometalate ions
283 using in situ PDF data, and for identifying stacking fault domain sizes in manganese oxides from
284 PDF and PXRD.[36, 37] We also note that the method is sufficiently fast that it would be possible to
285 combine it with structural screening applications such as structureMining@PDFitc.[15, 37] Given
286 chemical information about elements that are present, structureMining searches structural databases
287 for candidate structures. These are then refined to a target dataset and a rank ordered list returned
288 to the user. If the PDF represents a signal from a short-range ordered structural motif, we could
289 insert ML-MotEx between the database mining and refinement steps to search over sets of plausible
290 structures to look for structural sub-motifs. It may be possible to first use structure mining to identify
291 starting models, which could then be used for ML-MotEx analysis. The models could then be further
292 evaluated using both the resulting R_{wp} values and confidence factor. The ML-MotEx method is
293 currently limited to PDF analysis in the fitting procedure of the algorithm (step 2), however, the rest
294 of ML-MotEx (step 1+3+4) is ready to use with data from other techniques. We are confident that a
295 similar approach, taking advantage of explainable ML and SHAP values can be broadly useful for
296 enhancing and developing how models for data analysis are identified and constructed.

297 4 Methods

298 4.1 Step 1: Creation of a catalogue of candidate structure motifs

299 The first step in ML-MotEx is to use a starting structure model to generate a catalogue of candidate
300 structure motifs, which are all fitted to the data. The structures are generated by removing different
301 numbers of atoms from the original starting structure resulting in thousands of smaller, candidate
302 structure motifs. This process, which we refer to as ‘structure permutation’, is illustrated in Figure 1,
303 step 1. Here, the starting model contains 4 metal atoms, which are each bonded to 6 oxygen atoms.
304 Before candidate structure motifs are generated, we select which atom type should be included in the
305 permutation process. For the project discussed here, this selection is based on the X-ray scattering
306 power of the atoms (i.e., heavier atoms scatter X-rays strongly, while lighter ones do not), and we
307 therefore choose to permute over the 4 metal atoms in the structure rather than oxygen atoms. The
308 total number of atoms that are selected for permutation (here 4) is referred to as the permutation
309 number, N . Note that we do not take symmetry into account in this process. The selected atoms are
310 removed or kept in the model by randomly associating them with zeros and ones, where 0 means
311 that we remove the atom and 1 means we keep it. This is repeated multiple times to generate a large
312 catalogue of candidate structure motifs. The total number of possible motifs from the permutations

313 is equal to $2^N - 1$, but only a small fraction of these needs to be produced for ML-MotEx to provide
314 satisfactory results. In section J in the SI, we discuss how large a catalogue of candidate structure
315 motifs ML-MotEx needs as training data to output reasonable results. This is likely to be highly
316 system dependent and especially dependent on N and structure symmetry. For the examples presented
317 in the paper, we use ≈ 140 –3000 structure motifs per N. The atoms which were not chosen for
318 permutation, in this case oxygen, are removed if they are not within a distance threshold from any
319 other atom. The threshold is user-defined and can be set according to PDF peaks and/or chemically
320 valid distances (i.e., bond lengths) for the expected compounds.

321 4.2 Step 2: Fitting the catalogue of candidate structure motifs to the data

322 We fit each of the candidate structures in the catalogue to the experimental PDF using the Python-
323 based program DiffPy-CMI[30, 38-40]. We apply the Debye equation for calculation of scattering
324 intensities and PDFs from the structures. The fitting strategies and parameters for all 4 examples are
325 listed in section K in the SI including a description of the fit quality measure, R_{wp} .

326 4.3 Step 3: Predicting R_{wp} values using Gradient Boosting Decision Trees

327 GBDTs[25] are a tool that can do classification or regression using decision trees. In this work, we
328 are using XGBoost[25] as the GBDT algorithm to do the regression task of predicting the fit quality
329 (step 2) based on the structural input given as zeros or ones (step 1) and the number of atoms in the
330 structure. Section L in the SI demonstrates how the structure can be given as an input to the GBDT
331 model. The optimisation is done by making trees of ‘yes’ and ‘no’ questions on whether to keep an
332 atom in the structure or not, based on the resulting R_{wp} value. A hypothetical example of a simple
333 tree can be seen in Figure 1, step 3. When atom 4 is present in the structure, the GBDT model will
334 predict a R_{wp} value which is 5 % lower than if atom 4 is not present in the structure. In the same
335 way, it will predict an R_{wp} value which is 12 % lower if atom 1 is present in the structure. In the
336 decision tree, the algorithm will therefore say ‘yes’ to keep both atoms 1 and 4 in the structure. In
337 this project, the GBDT model predicts the R_{wp} value using a weighted average of 100 trees. The
338 GBDT model performance is improved with a large amount of training data, which in this tool is
339 provided by creating a larger catalogue of candidate structure motifs and fitting them to the data. The
340 GBDT model is trained on 80 % of the data, which is referred to as the training set. XGBoost[25]
341 were used with default parameters except for learning rate and max depth, which were optimised
342 with the use of Bayesian optimization using 50 iterations and cross-validation split on 3.[41, 42]
343 While this procedure automates the hyperparameter tuning, we demonstrate in section M in the SI
344 that similar results are achieved across various hyperparameters. The last 20 % of the data is used to
345 evaluate the performance of the algorithm and is referred to as test set.

346 4.4 Step 4: Quantifying the contribution of each atom using SHAP values

347 SHAP values are used to analyse the R_{wp} values resulting from the process described above. For
348 each fit (step 2), each atom in the starting model is assigned a SHAP value. The amplitude of
349 the SHAP value reflects how important a structural feature is for the fit quality, while the sign of
350 the SHAP value reflects whether the feature affects the R_{wp} value of the fit towards 1 (poor fit)
351 or 0 (perfect fit), in other words why it is important. Each atom in the starting model will thus
352 get F number of SHAP values, where F corresponds to the number of fits made in step 2 of the
353 algorithm. We divide the F number of SHAP values into two categories; firstly the ones where the
354 atom was kept in the structure motif (kept atom SHAP value list) and secondly the ones where the
355 atom was removed to create the structure motif (removed atom SHAP value list). From each of the
356 two lists, an average SHAP value for the atoms can be calculated, defined as $\text{SHAP}_{\text{average-kept}}$
357 and $\text{SHAP}_{\text{average-removed}}$. We then define an atom contribution value, which is calculated as the
358 difference between two average SHAP values, i.e. atom contribution value = $\text{SHAP}_{\text{average-kept}}$
359 - $\text{SHAP}_{\text{average-removed}}$. We also define the uncertainty on this value as: atom contribution value
360 $\text{RMS} = (\text{SHAP}_{\text{average-kept}}^2 - \text{SHAP}_{\text{average-removed}}^2)^{0.5}$. We define a confidence factor for each atom
361 that describes how confident we can be about including/excluding that atom in a structural motif;
362 Confidence factor = atom contribution value / atom contribution value RMS. ML-MotEx outputs a
363 VESTA[43] and CrystalMaker[44] file where all the atoms are coloured with regard to their atom
364 contribution value.

365 **5 References**

- 366 1. Billinge, S.J.L. and M.G. Kanatzidis, Beyond crystallography: the study of disorder, nanocrystallinity and crystallographically challenged materials with pair distribution functions. *Chemical Communications*, 2004(7): p. 749-760.
- 367
368
- 369 2. Keen, D.A. and A.L. Goodwin, The crystallography of correlated disorder. *Nature*, 2015. 521(7552): p. 303-309.
- 370
- 371 3. Christiansen, T.L., S.R. Cooper, and K.M.Ø. Jensen, There's no place like real-space: elucidating size-dependent atomic structure of nanomaterials using pair distribution function analysis. *Nanoscale Advances*, 2020. 2(6): p. 2234-2254.
- 372
373
- 374 4. Billinge, S.J.L. and I. Levin, The Problem with Determining Atomic Structure at the Nanoscale. *Science*, 2007. 316(5824): p. 561-565.
- 375
- 376 5. Juelsholt, M., et al., Size-induced amorphous structure in tungsten oxide nanoparticles. *Nanoscale*, 2021. 13(47): p. 20144-20156.
- 377
- 378 6. Yang, X., et al., Confirmation of disordered structure of ultrasmall CdSe nanoparticles from X-ray atomic pair distribution function analysis. *Physical Chemistry Chemical Physics*, 2013. 15(22): p. 8480-8486.
- 379
380
- 381 7. Christiansen, T.L., et al., Structure analysis of supported disordered molybdenum oxides using pair distribution function analysis and automated cluster modelling. *Journal of Applied Crystallography*, 2020. 53(1): p. 148-158.
- 382
383
- 384 8. Bennett, T.D. and A.K. Cheetham, Amorphous Metal–Organic Frameworks. *Accounts of Chemical Research*, 2014. 47(5): p. 1555-1562.
- 385
- 386 9. Kjær, E.S.T., et al., DeepStruc: Towards structure solution from pair distribution function data using deep generative models. *ChemRxiv*, 2022. doi:10.26434/chemrxiv-2022-0zrdl.
- 387
- 388 10. Cliffe, M.J., et al., Structure determination of disordered materials from diffraction data. *Physical review letters*, 2010. 104(12): p. 125501.
- 389
- 390 11. Juhás, P., et al., *ab initio* determination of solid-state nanostructure. *Nature*, 2006. 440(7084): p. 655-658.
- 391
- 392 12. Juhás, P., et al., The Liga algorithm for *ab initio* determination of nanostructure. *Acta Crystallogr A*, 2008. 64(Pt 6): p. 631-40.
- 393
- 394 13. Christiansen, T.L., et al., Size Induced Structural Changes in Molybdenum Oxide Nanoparticles. *ACS Nano*, 2019. 13(8): p. 8725-8735.
- 395
- 396 14. Aalling-Frederiksen, O., et al., Formation and growth mechanism for niobium oxide nanoparticles: atomistic insight from in situ X-ray total scattering. *Nanoscale*, 2021. 13(17): p. 8087-8097.
- 397
- 398 15. Yang, L., et al., Structure-mining: screening structure models by automated fitting to the atomic pair distribution function over large numbers of models. *Acta Crystallographica Section A*, 2020. 76(3): p. 395-409.
- 399
400
- 401 16. Anker, A.S., et al., Characterising the Atomic Structure of Mono-Metallic Nanoparticles from X-Ray Scattering Data Using Conditional Generative Models. 2020: Proceedings of the 16th International Workshop on Mining and Learning with Graphs (MLG).
- 402
403
- 404 17. Banerjee, S., et al., Cluster-mining: an approach for determining core structures of metallic nanoparticles from atomic pair distribution function data. *Acta Crystallographica Section A*, 2020. 76(1): p. 24-31.
- 405
406
- 407 18. Anker, A.S., et al., Structural Changes during the Growth of Atomically Precise Metal Oxide Nanoclusters from Combined Pair Distribution Function and Small-Angle X-ray Scattering Analysis. *Angewandte Chemie International Edition*, 2021. 60: p. 2-12.
- 408
409
- 410 19. Butler, K.T., et al., Interpretable, calibrated neural networks for analysis and understanding of inelastic neutron scattering data. *Journal of Physics: Condensed Matter*, 2021. 33(19): p. 194006.
- 411

- 412 20. Suzuki, Y., et al., Symmetry prediction and knowledge discovery from X-ray diffraction patterns
413 using an interpretable machine learning approach. *Scientific Reports*, 2020. 10(1): p. 21790.
- 414 21. Torrisi, S.B., et al., Random forest machine learning models for interpretable X-ray absorption
415 near-edge structure spectrum-property relationships. *npj Computational Materials*, 2020. 6(1): p.
416 109.
- 417 22. Oviedo, F., et al., Interpretable and Explainable Machine Learning for Materials Science and
418 Chemistry. arXiv preprint arXiv:2111.01037, 2021.
- 419 23. Schmidt, J., et al., Recent advances and applications of machine learning in solid-state materials
420 science. *npj Computational Materials*, 2019. 5(1): p. 83.
- 421 24. Lee, K., et al., Phase classification of multi-principal element alloys via interpretable machine
422 learning. *npj Computational Materials*, 2022. 8(1): p. 25.
- 423 25. Chen, T. and C. Guestrin, XGBoost: A Scalable Tree Boosting System, in *Proceedings of the*
424 *22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2016,
425 Association for Computing Machinery: San Francisco, California, USA. p. 785–794.
- 426 26. Lundberg, S.M., et al., From local explanations to global understanding with explainable AI for
427 trees. *Nature Machine Intelligence*, 2020. 2(1): p. 56-67.
- 428 27. Lundberg, S.M. and S.-I. Lee, A Unified Approach to Interpreting Model Predictions. *Proceedings*
429 *of the 31st International Conference on Neural Information Processing Systems*, 2017: p. 4765-4774.
- 430 28. Juelsholt, M., T. Lindahl Christiansen, and K.M.Ø. Jensen, Mechanisms for Tungsten Oxide
431 Nanoparticle Formation in Solvothermal Synthesis: From Polyoxometalates to Crystalline Materials.
432 *The Journal of Physical Chemistry C*, 2019. 123(8): p. 5110-5119.
- 433 29. Chen, X. and S. Yamanaka, Single-crystal X-ray structural refinement of the ‘tetragonal’ C₆₀
434 polymer. *Chemical Physics Letters*, 2002. 360(5): p. 501-508.
- 435 30. Juhás, P., et al., Complex modeling: a strategy and software program for combining multiple
436 information sources to solve ill posed structure and nanostructure inverse problems. *Acta Crystallogr*
437 *A Found Adv*, 2015. 71(Pt 6): p. 562-568.
- 438 31. Niu, J., et al., Syntheses, spectroscopic characterization, thermal behavior, electrochemistry and
439 crystal structures of two novel pyridine metatungstates. *Journal of Coordination Chemistry*, 2004.
440 57(11): p. 935-946.
- 441 32. Joachim, F., T. Axel, and P. Rosemarie, Strukturen und Schwingungsspek-
442 tren des Tetramethylammonium- α -dodekawolframatosilikats und des Tetrabutylammonium- β -
443 dodekawolframatosilikats: Structures and Vibrational Spectra of Tetramethylammonium α -
444 Dodecatungstosilicate and Tetrabutylammonium β -Dodecatungstosilicate. *Zeitschrift für Natur-*
445 *forschung B*, 1981. 36(2): p. 161-171.
- 446 33. Niu, J.-Y., Q.-X. Han, and J.-P. Wang, A Novel Keggin Units-Supported Complex: Synthesis,
447 Characterization and Crystal Structure of $[(\text{CH}_3)_2\text{NH}_2]_6[\text{Cu}(\text{DMF})_4(\text{GeW}_{12}\text{O}_{40})_2] \cdot 2\text{DMF}$. *Journal*
448 *of Coordination Chemistry*, 2003. 56(6): p. 523-530.
- 449 34. Busbongthong, S. and T. Ozeki, Structural Relationships among Methyl-, Dimethyl-, and
450 Trimethylammonium Phosphododecatungstates. *Bulletin of the Chemical Society of Japan*, 2009.
451 82(11): p. 1393-1397.
- 452 35. Skjærø, S.L., et al., Atomic structural changes in the formation of transition metal tungstates:
453 The role of polyoxometalate structures in material crystallization. 2022.
- 454 36. Magnard, N., et al., Characterisation of intergrowth in metal oxide materials using structure-
455 mining: the case of γ -MnO₂. 2022. 37. Yang, L., et al., A cloud platform for atomic pair distribution
456 function analysis: PDFitc. *Acta Crystallographica Section A*, 2021. 77(1): p. 2-6.
- 457 38. Proffen, T. and R.B. Neder, DISCUS, a program for diffuse scattering and defect structure
458 simulations – update. *Journal of Applied Crystallography*, 1999. 32(4): p. 838-839.
- 459 39. Proffen, T. and R.B. Neder, DISCUS: a program for diffuse scattering and defect-structure
460 simulation. *Journal of Applied Crystallography*, 1997. 30(2): p. 171-175.

- 461 40. Coelho, A.A., TOPAS and TOPAS-Academic: an optimization program integrating computer
 462 algebra and crystallographic objects written in C++. Journal of Applied Crystallography, 2018. 51(1):
 463 p. 210-218.
- 464 41. Nogueira, F., Bayesian Optimization: Open source constrained global optimization tool for
 465 Python. 2014.
- 466 42. Putatunda, S. and K. Rama, A Comparative Analysis of Hyperopt as Against Other Approaches
 467 for Hyper-Parameter Optimization of XGBoost. Proceedings of the 2018 International Conference
 468 on Signal Processing and Machine Learning, 2018: p. 6-10.
- 469 43. Momma, K. and F. Izumi, VESTA 3 for three-dimensional visualization of crystal, volumetric
 470 and morphology data. Journal of Applied Crystallography, 2011. 44(6): p. 1272-1276.
- 471 44. Palmer, D.C., Visualization and analysis of crystal structures using CrystalMaker software.
 472 Zeitschrift für Kristallographie - Crystalline Materials, 2015. 230(9-10): p. 559-572.

473 Checklist

- 474 1. For all authors...
- 475 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
 476 contributions and scope? [Yes]
- 477 (b) Did you describe the limitations of your work? [Yes] Line 287 - 294.
- 478 (c) Did you discuss any potential negative societal impacts of your work? [No] We have
 479 not identified any potential negative societal impacts of this work.
- 480 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 481 them? [Yes]
- 482 2. If you are including theoretical results...
- 483 (a) Did you state the full set of assumptions of all theoretical results? [N/A] We do not
 484 include theoretical results
- 485 (b) Did you include complete proofs of all theoretical results? We do not include theoretical
 486 results
- 487 3. If you ran experiments...
- 488 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 489 mental results (either in the supplemental material or as a URL)? [Yes] The code, data
 490 and instructions needed to reproduce the results are uploaded to Github but not shared
 491 here due to double-blind reviewing.
- 492 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 493 were chosen)? [Yes]
- 494 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 495 ments multiple times)? [Yes] Not in a traditional way but we have reported the results
 496 using many different seeds, starting models and hyperparameters and shown that it
 497 nearly provide identical results.
- 498 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 499 of GPUs, internal cluster, or cloud provider)? [Yes]
- 500 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 501 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 502 (b) Did you mention the license of the assets? [No] The assets we have used are standard
 503 open-source software as XGBoost, scikit learn ect. Our software is also open-source.
- 504 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 505 (d) Did you discuss whether and how consent was obtained from people whose data you're
 506 using/curating? [No] Consent was obtained from all people whose data is used. They
 507 are all part of the author list.
- 508 (e) Did you discuss whether the data you are using/curating contains personally identifi-
 509 able information or offensive content? [No] We have not identified any offensive or
 510 identifiable information in our X-ray scattering data.

- 511 5. If you used crowdsourcing or conducted research with human subjects...
- 512 (a) Did you include the full text of instructions given to participants and screenshots, if
- 513 applicable?
- 514 (b) Did you describe any potential participant risks, with links to Institutional Review
- 515 Board (IRB) approvals, if applicable?
- 516 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 517 spent on participant compensation?