

---

# DeepStruc: Towards structure solution from pair distribution function data using deep generative models

---

**Emil T. S. Kjær\***

Department of Chemistry and Nano-Science Center  
University of Copenhagen  
2100 Copenhagen Ø, Denmark  
etsk@chem.ku.dk

**Andy S. Anker\***

Department of Chemistry and Nano-Science Center  
University of Copenhagen  
2100 Copenhagen Ø, Denmark  
andy@chem.ku.dk

**Marcus N. Weng**

Department of Chemistry and Nano-Science Center  
University of Copenhagen  
2100 Copenhagen Ø, Denmark  
qkv769@alumni.ku.dk

**Simon J. L. Billinge**

Department of Applied Physics and Applied Mathematics  
Columbia University, New York  
NY 10027, USA  
Condensed Matter Physics and Materials Science Department  
Brookhaven National Laboratory, Upton  
NY 11973, USA  
sb2896@columbia.edu

**Raghavendra Selvan**

Department of Computer Science  
University of Copenhagen  
2100 Copenhagen Ø, Denmark  
Department of Neuroscience  
University of Copenhagen  
2200 Copenhagen N, Denmark  
raghav@di.ku.dk

**Kirsten M. Ø. Jensen**

Department of Chemistry and Nano-Science Center  
University of Copenhagen  
2100 Copenhagen Ø, Denmark  
kirsten@chem.ku.dk

## Abstract

1 Structure solution of nanostructured materials that have limited long-range order  
2 remains a bottleneck in materials development. We present a deep learning algo-  
3 rithm, DeepStruc, that can solve a simple nanoparticle structure directly from a

---

\*Both authors contributed equally to this work

4 Pair Distribution Function (PDF) obtained from total scattering data by using a  
5 conditional variational autoencoder. We first apply DeepStruc to PDFs from seven  
6 different structure types of monometallic nanoparticles, and show that structures  
7 can be solved from both simulated and experimental PDFs, including PDFs from  
8 nanoparticles that are not present in the training distribution. We also apply Deep-  
9 Struc to a system of *hcp*, *fcc* and stacking faulted nanoparticles, where DeepStruc  
10 recognizes stacking faulted nanoparticles as an interpolation between *hcp* and *fcc*  
11 nanoparticles and is able to solve stacking faulted structures from PDFs. Our  
12 findings suggests that DeepStruc is a step towards a general approach for structure  
13 solution of nanomaterials.

## 14 1 Introduction

15 Crystallographic methods, such as single crystal and powder diffraction, have been foundational in  
16 the development of functional materials over the past century. They yield atomic-scale structural  
17 models for crystalline materials and allow establishing the links between material structure and  
18 properties that are at the heart of materials development.[1,2] However, other approaches for structure  
19 determination are needed for nanostructured materials that have limited long-range order, and total  
20 scattering methods such as atomic pair distribution function (PDF) analysis have become increasingly  
21 important tools.[3-7] Currently, PDF analysis is mainly done by fitting a known starting model to  
22 an experimental PDF, a process known as structure refinement. Recent developments in automated  
23 modelling[8-10] have made it possible to extend the searched structural space, but identifying a  
24 model or solving a structure *de novo* from a PDF is still an enormous challenge. So far, only highly  
25 symmetrical nanostructures such as the  $C_{60}$  buckyball have been solved *ab initio* from a PDF.[11-15]  
26 Determining the structure of less symmetrical nanostructures is limited by the lost information caused  
27 by PDF peak overlap, which challenges the use of PDF for structure solution of more complicated  
28 nanomaterials.

29 An approach to handle the challenges due to the information barrier in PDFs is to employ supervised  
30 machine learning (ML) methods that can learn from well-known PDF-structure pairs. While deter-  
31 mining a unique structure from a PDF is not always a solvable problem, as several different structures  
32 may give rise to identical PDFs, ML methods can still learn to capture the relationship between PDF  
33 and structure and thereby push the boundaries of nanostructure solution from PDF. When there is not  
34 enough information in the PDF to provide a unique structure solution, ML methods may provide a  
35 distribution of starting models which can aid in further structure analysis. In this work, we use deep  
36 generative models (DGMs). DGMs are a class of ML models that can estimate the underlying data  
37 distribution from a reasonably small set of training examples.[16] A well-known use case of DGMs  
38 is in the generation of synthetic ‘deep-fake’ images[17,18] based on large datasets of real images. We  
39 here train our DGM to identify new structure models by training on known chemical structures. The  
40 DGM learns the relation between PDF and atomic structure, which enables it to solve a structure,  
41 based on a PDF it has not seen before and its learned chemical knowledge.

42 We apply our DGM, which we refer to as ‘DeepStruc’, for structural analysis of a model system of  
43 monometallic nanoparticles (MMNPs) with seven different structure types (Fig. 1a) and demonstrate  
44 the method for both simulated and experimental PDFs. DeepStruc is generative, which means that it  
45 can be used to construct structures that are not in the training set, i.e., solve a structure from a PDF.  
46 We demonstrate this capability on a dataset of face-centered cubic (*fcc*), hexagonal closed packed  
47 (*hcp*) and stacking faulted structures, where DeepStruc can recognize the stacking faulted structures  
48 as an interpolation between *fcc* and *hcp* and construct new structural models based on a PDF.

## 49 2 Results

### 50 2.1 Training DeepStruc to determine the structure of MMNPs from PDF data

51 DeepStruc, illustrated in Fig. 1a and discussed below, is a graph-based conditional variational autoen-  
52 coder (graph CVAE). Autoencoders are a class of deep learning (DL) methods where high-dimensional  
53 inputs, such as chemical structures,[19,20] are reduced in dimensionality. The transformation into 2  
54 or 3 dimensional vectors is achieved using an information bottleneck by an encoder neural network  
55 (NN),[19,21,22] and the resulting lower-dimensional, compressed feature space is known as the latent

56 space. A decoder NN can reconstruct the input from these low-dimensional representations. When  
57 the latent space is regularized (smoothed) using normal distributions instead of discrete points we  
58 obtain a variational autoencoder (VAE). It has previously been demonstrated that VAEs do a better  
59 job interpolating in the latent space compared to deterministic AEs.[19] The VAE can be made to be  
60 dependent (conditioned) on additional information by the prior NN resulting in a CVAE.[22]

61 We here use MMNP structures (Fig. 1b) as input, and condition them on their simulated PDFs (Fig.  
62 1c). The MMNP structures span seven different structure types computed using a variety of metals  
63 to emulate the variability in bond lengths in real metallic nanoparticle samples. The structure types  
64 are simple cubic (*sc*), body-centered cubic (*bcc*), face-centered cubic (*fcc*), hexagonal closed packed  
65 (*hcp*), decahedral, icosahedral, and octahedral, and all structure types have been constructed in sizes  
66 from 5 to 200 atoms. We used 3743 MMNP structures, which were randomly split into training-  
67 (60 %), validation- (20 %) and testing-sets (20 %). Note that the validation and test sets are derived  
68 from the same underlying data distribution as the training set, and serve as intermediaries to the  
69 actual test set which is based on the experimental PDF data. A histogram of the distribution of  
70 the seven structure types are provided in section A in the Supplementary Information. During the  
71 training process (blue + green region Fig. 1a), DeepStruc learns to map the conditioning PDFs to  
72 their structures in the latent space. After the training process is complete, DeepStruc can be used on  
73 data that have not been part of the training set, which is referred to as ‘inference’. Further details  
74 about the DeepStruc network can be found in section B in the Supplementary Information.

## 75 2.2 Mapping of structures in a latent space

76 We first evaluate DeepStruc’s ability to map the MMNP structures in a low-dimensional latent space  
77 by investigating structural trends and clustering. Fig. 2 shows a visualization of the two-dimensional  
78 latent space with selected MMNP reconstructions indicated. The colour of the points indicates the  
79 structure type, and the relative point size indicates the size of the MMNP cluster. We observe that  
80 DeepStruc learns to map the chemical structures in the latent space by size and symmetry. It maps the  
81 cubic structure types (*sc*, *bcc*, and *fcc*) together, and it learns that the octahedral MMNPs are closely  
82 related to the *fcc* structure type. Interestingly, DeepStruc also allocates the decahedral structures to  
83 be in between the *fcc* and *hcp* structures. This can be rationalized by considering that decahedral  
84 structures are constructed from five tetrahedrally shaped *fcc* crystals which are separated by {111}  
85 twin boundaries that resemble stacking faults.[9,23,24] The twin boundaries will resemble stacking  
86 faulted regions of *fcc* justifying that they exist in the latent space between *fcc* and *hcp*.

## 87 2.3 DeepStruc for structure determination from PDF

88 We now move on to identify structures directly from a PDF. The results of using DeepStruc on seven  
89 simulated PDFs of MMNPs not used in the training process are illustrated in Fig 3. Here, we show  
90 the structure that the input PDF was calculated from (left), the reconstructed structure (right), and  
91 its agreement with the input PDF after structure refinement (middle, discussed below). In all seven  
92 cases, the structures are correctly reconstructed from the PDF input. Before structure refinement, the  
93 mean absolute error (MAE) of the atom positions is  $0.128 \pm 0.073 \text{ \AA}$  as described in section C in the  
94 Supplementary Information. However, the MAE is artificially high due to a common aberration by  
95 DeepStruc, where it predicts the right geometric atomic arrangement, but isotropically contracted  
96 or expanded compared to the original structure. We do not yet understand why DeepStruc has this  
97 aberration, but it is easily solvable by refining an expansion/contraction variable as a post processing  
98 step to DeepStruc. After refining the structure to the PDF[25] by fitting a contraction/expansion  
99 factor, a scale factor and an isotropic atomic displacement parameter (ADP), as described in section  
100 C in the Supplementary Information, the MAE of the atom positions is reduced to  $0.093 \pm 0.058$   
101  $\text{\AA}$ . The inference is thus robust against moderate changes in lattice parameter between a provided  
102 PDF and the structures that DeepStruc were trained on. The reconstructed structures exhibit some  
103 artificial positional atomic disorder that broadens the PDF peaks. The fitted ADP values (section C in  
104 the Supplementary Information) are thus lower than the ADP values of the conditioning PDFs.

105 Having established that DeepStruc works for structures highly resembling those in the training set, we  
106 now consider more challenging cases and explore the capabilities of DeepStruc on an actual test set  
107 which is far from the training distribution. As described above, the largest structures in the training  
108 set contained only 200 atoms.

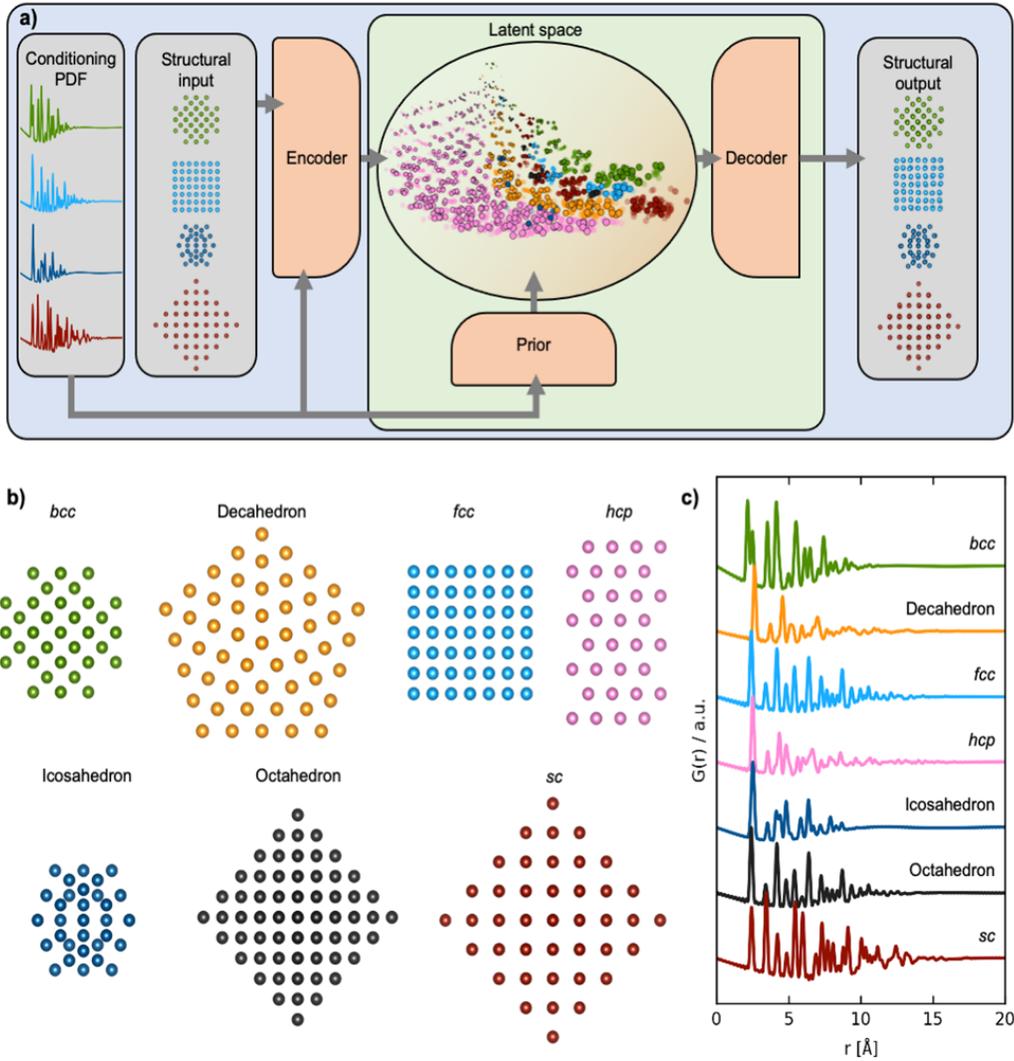


Figure 1: Training DeepStruc to determine the structure of MMNPs from PDFs. a) DeepStruc predicts the xyz-coordinates of the MMNP structure with conditional input provided in the form of a PDF. The encoder uses the structure and its PDF as input while the prior only takes the PDF as input. To obtain the structural output a latent space embedding is given as input to the decoder which produces the corresponding MMNP xyz-coordinates. During training of DeepStruc both the blue and green regions are used, while only the green region is used for structure prediction during the inference process. b) Examples of the seven different structure types which are used as input to DeepStruc together with their c) simulated PDFs used as conditioning in DeepStruc. Each structure type has been included in the training set with varying sizes of 5 to 200 atoms and with varying lattice constants. The 3743 structures were split into training- (60 %), validation- (20 %), and testing sets (20 %).

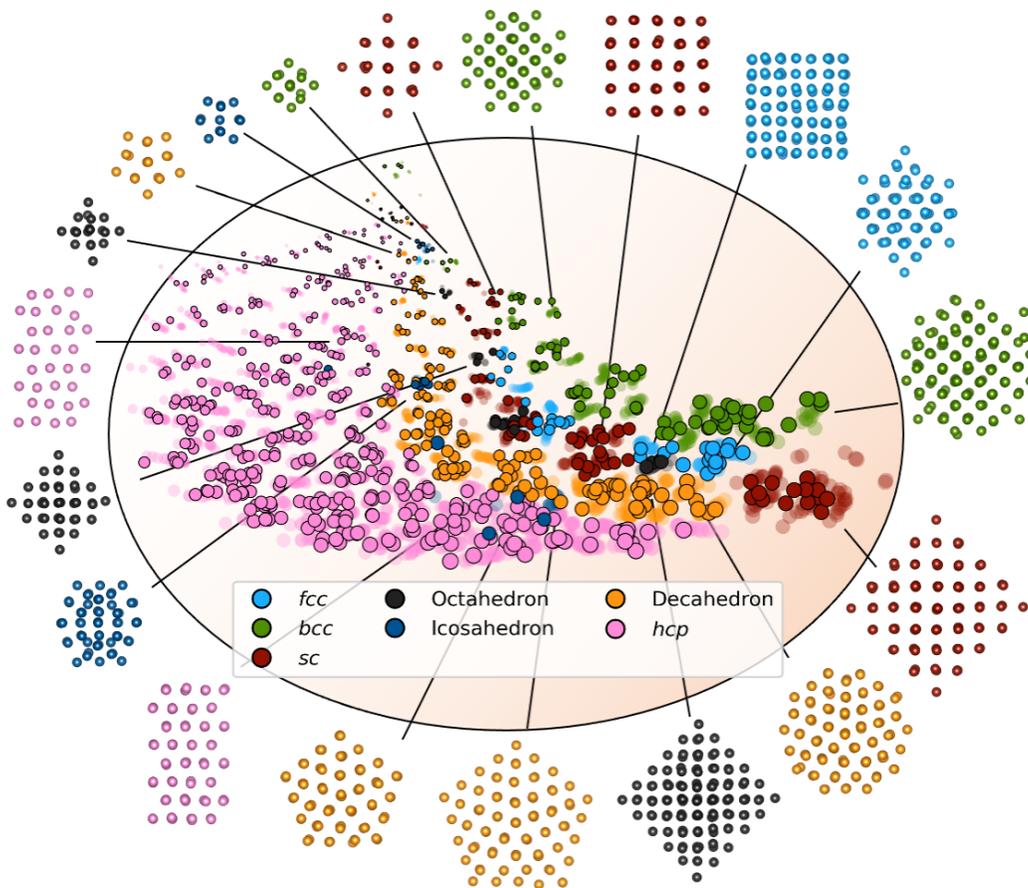


Figure 2: The two-dimensional latent space with structure reconstructions. The points in the latent space correspond to a structure and its simulated PDF. Data points from the test set are shown in solid colour and outlined. The points from the training and validation sets are shown as semi-transparent. The size of the points relates to the size of the embedded MMNP, and the orange background indicates the general size increase throughout the latent space. The colour of each point resembles its structure type, *fcc* (light blue), octahedral (dark grey), decahedral (orange), *bcc* (green), icosahedral (dark blue), *hcp* (pink), and *sc* (red). Note that the structures shown here are predicted by DeepStruc during inference on PDFs from the test set.

109 We now evaluate it on a test set of simulated MMNPs with 5 to 1000 atoms, i.e., containing much  
 110 larger particles. The latent space obtained from this new test set is plotted using diamond markers in  
 111 Fig. 4, where the latent space from the training process is shown with semi-transparent markers. We  
 112 observe that the trends in the training area are comparable for the training set and the test set of larger  
 113 MMNPs. Notably, the trends of both the size and the structure types continue beyond the training area  
 114 to structures containing about 400 atoms. Beyond 400 atoms, all structure types collapse onto a line,  
 115 however, DeepStruc still estimates the size of the structure. Of course, DeepStruc could be retrained  
 116 on a larger training set if reconstructions are desired on clusters larger than 200 atoms. However, this  
 117 experiment shows that DeepStruc can extrapolate significantly in the latent space. It can thereby give  
 118 useful information about PDFs from structures not represented in the training set and is generative in  
 119 a meaningful way. This can be compared to, for example, a tree-based ML-classifier, which is limited  
 120 to a predefined structural database and cannot extrapolate. The capability of DeepStruc to extrapolate  
 121 arises from each structure in the latent space being predicted as a normal distribution instead of a  
 122 discrete point.

123 In practice, DeepStruc must be able to yield valid reconstructed structures from experimental data  
 124 that contain noise and other aberrations. We therefore use DeepStruc to infer structures from

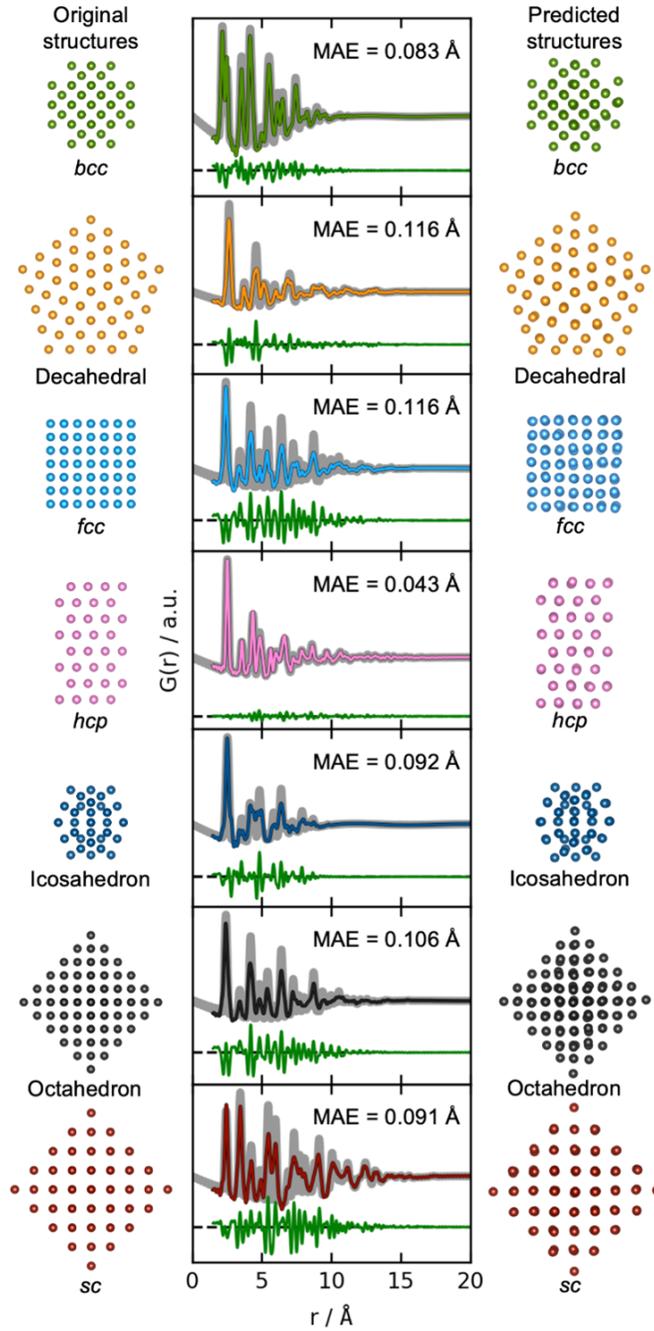


Figure 3: Structure determination from PDFs. Simulated PDFs (grey) from the original structures of the seven different structure types (left) are used during inference for structure prediction (right). The middle column shows the fitted PDFs of the predicted structures to the simulated PDFs of the original structures. Only the scale-factor, contraction/expansion-factor, and ADP are refined, see section B in the Supplementary Information.

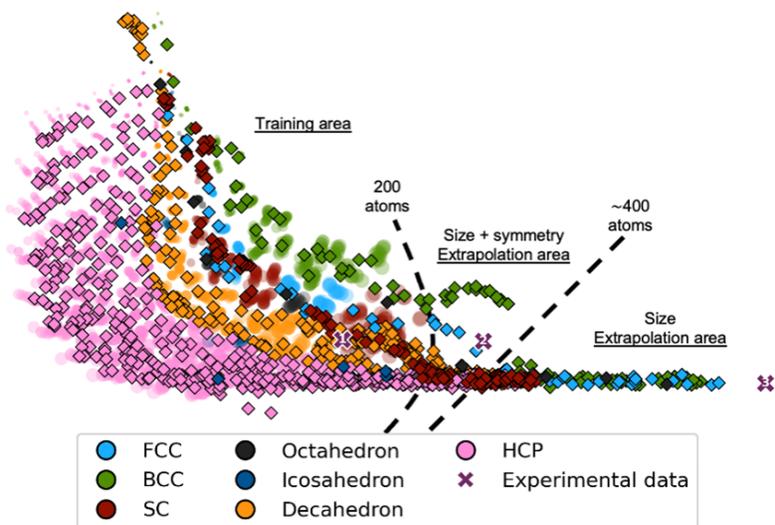


Figure 4: DeepStruc applied on PDFs of structures up to 1000 atoms. Each point is coloured after its structure type, i.e. *fcc* (light blue), octahedral (dark grey), decahedral (orange), *bcc* (green), icosahedral (dark blue), *hcp* (pink), and *sc* (red). Each point in the latent space corresponds to a structure based on its simulated PDF. Test PDFs from structures up to 1000 atoms are plotted as diamond markers on top of the training and validation data which are made semi-transparent. Note that the training set latent space is identical to that plotted in Fig. 2. DeepStruc has only been trained on structures up to 200 atoms. Three experimental PDFs (shown in section C in the Supplementary Information) obtained from differently sized *fcc* nanocrystals estimated to contain 203 (cross marker 1), 371 (cross marker 2), and 1368 (cross marker 3) atoms are illustrated as purple cross markers in the latent space.

125 previously published experimental PDFs from MMNPs. Fig. 5a shows the latent space with the  
 126 predicted location of structures from three experimental PDFs. Here, the location in the latent space  
 127 is represented as distributions rather than as discrete points, and multiple structures are sampled from  
 128 each distribution and compared to the experimental PDF to select the best candidate. The mean of  
 129 the experimental PDF distributions is represented as a black diamond with three ellipsoids indicating  
 130 different confidence intervals with  $\sigma$ : 3, 5 and 7, where  $\sigma$  is the standard deviation of the normal  
 131 distribution.

132 The first experimental dataset that we evaluate was published by Jensen et al.,[26] who identified  
 133 a decahedral structure as the core motif of  $\text{Au}_{144}(\text{p-MBA})_{60}$  nanoparticles. DeepStruc locates the  
 134  $\text{Au}_{144}(\text{p-MBA})_{60}$  PDF (Fig. 5b) in a decahedral region (orange distributions in Fig. 5a) in the latent  
 135 space. Given the generative capabilities of DeepStruc, in theory, we can sample an unlimited number  
 136 of structures for a given PDF. As described in section D of the Supplementary Information, we here  
 137 sampled up to 1000 structures from the three normal distributions ( $\sigma$ : 3, 5, and 7), and compared their  
 138 fit to the experimental PDF. Fig. 5b shows the fit of the best structural prediction, which was among  
 139 the structures sampled from the  $\sigma$ : 3 distributions. DeepStruc predicts a decahedral structure, which  
 140 agrees well with the literature.[26] Other structures sampled from the three distributions are shown  
 141 in Section E of the Supplementary Information, where we also compare the DeepStruc analysis to  
 142 baseline methods. We first consider a brute-force structure-mining method inspired by Banerjee et  
 143 al.,[27], but also compare the DeepStruc results to two simpler ML-algorithms, namely a tree-based  
 144 ML classifier and a regular CVAE without a graph-based input.

145 The second dataset that we evaluate, published by Quinson et al.,[28] are from 1.8 nm Pt nanoparticles  
 146 with the *fcc* structure (described further in Section F in the Supplementary Information). This size  
 147 corresponds to ca. 203 atoms, i.e. the number of atoms in the particle goes slightly beyond the *fcc*  
 148 structures in the training set that contain only 165 atoms.[28] The location of the predicted mean is  
 149 again shown in Fig. 5a, enclosed by three blue ellipsoids illustrating different  
 150 magnitudes of standard deviation. The mean of the predicted structure is placed near the largest *sc*

151 structures. If DeepStruc only favoured symmetry it would be placed directly on the *fcc* structures.  
152 Interestingly, DeepStruc does not purely favour size either, as it does not position the PDF near the  
153 largest structures which are *hcp* structures of 200 atoms. Instead, we observe that DeepStruc takes  
154 both symmetry and size into account by placing the mean predicted structure adjacent to the largest  
155 *sc* structures containing 185 atoms. To identify the structure from the experimental PDF, we again  
156 sample 1000 structures from the  $\sigma$ : 3, 5 and 7 distributions. When fitting these sampled structures to  
157 the dataset, we obtain the best fit from an *fcc* structure of 146 atoms that is visualized in Fig. 5c and  
158 which agrees with the baseline models (section E in the Supplementary Information). DeepStruc thus  
159 identifies an *fcc* structure even though the size of the MMNP is outside the training set distribution.

160 We also attempted to input PDFs from even larger *fcc* nanoparticles, estimated to have diameters of  
161 2.2 and 3.4 nm, corresponding to 371 and 1368 atoms, respectively (section F in the Supplementary  
162 Information).[28] Their positions in the latent space are shown in Fig. 4 along with the 1.8 nm *fcc*  
163 nanoparticles using cross markers labelled 1, 2, and 3 for increasing size. We observe that they follow  
164 the trend of the simulated *fcc* structures discussed above: while it is possible to estimate both size and  
165 symmetry for the 2.2 nm particles through extrapolation, DeepStruc can only estimate size for the 3.4  
166 nm particle. We note that the size can be read from a PDF directly without any modelling. However,  
167 the ability of DeepStruc to predict structures on experimental data beyond those in the training set is  
168 promising for future structure solution from PDF.

169 While DeepStruc only has been trained on simple MMNPs, we finally evaluate it on a PDF from  
170  $\text{Au}_{144}(\text{PET})_{60}$  nanoparticles, consisting of an icosahedral core of 54 atoms surrounded by a rhombi-  
171 cosidodecahedron shell of 60 atoms (Fig. 5d and e).[26,29] We show the predicted mean position of  
172 the structure with a black diamond enclosed by pink ellipsoids. DeepStruc positions the PDF in the  
173 *hcp* region of the latent space, and when sampling 1000 structures from the distribution with  $\sigma$ : 7,  
174 the best fitting structures is an *hcp* structure with 40 atoms for the  $\text{Au}_{144}(\text{PET})_{60}$  nanoparticle (Fig.  
175 5d). Similar structures are found when sampling from the  $\sigma$ : 3 and  $\sigma$ : 5 distributions. However, the  
176 PDF fit reveals that the reconstructed structure does not capture all peaks in the experimental PDF.  
177 When considering further the latent space, icosahedral structures are strongly underrepresented in our  
178 dataset (section A in the Supplementary Information) which results in an inconsistency when placing  
179 icosahedral structures in the latent space. DeepStruc is thus challenged when solving the icosahedral  
180 core structure of the nanoparticle. However, we observe that one of the test icosahedral structures is  
181 placed near the experimental PDF in latent space within the  $\sigma$ : 5 distribution. Therefore, we again  
182 try to sample 1000 structures by moving the mean of the  $\sigma$ : 3 distribution to the nearest cluster  
183 of icosahedral structures in the latent space, which are located right outside the  $\sigma$ : 7 distribution.  
184 The best fitting structure (Fig. 5e) captures all main peaks of the experimental PDF. Strategies  
185 for sampling of underrepresented structures is discussed further in section D in the Supplementary  
186 Information.

## 187 2.4 Structure determination from PDF: *fcc*, *hcp*, and stacking faulted nanoparticles

188 To obtain a deeper understanding of the latent space’s behaviour, we investigate a dataset only  
189 containing *fcc*, *hcp*, and stacking faulted structures. *Fcc* and *hcp* structures are distinguished by the  
190 stacking sequence of closed packed layers in their structures: while *fcc* structures can be described  
191 by ABCABC stacking, *hcp* structures have ABABAB stacking. Structures with other sequences are  
192 stacking faulted structures. We hypothesize that stacking faulted structures can be considered an  
193 ‘interpolation’ in the discrete space between the *fcc* and *hcp* structure type.[30]

194 Examples of reconstructed *fcc* (blue), *hcp* (pink), and different stacking faulted structures (purple)  
195 and their position in the new latent space are illustrated in Fig. S8a. The MMNPs cluster in size,  
196 whilst we also observe that *fcc* and *hcp* structures separate in the latent space. It is evident that the  
197 stacking faulted structures are located in between the *fcc* and *hcp* structures in the latent space as  
198 hypothesized. It is chemically reasonable that they are positioned in this exact order based on their  
199 similarity to *fcc* and *hcp*. For example, the structure with ABCABA layers, shown in Fig. S8 with a  
200 purple star is structurally close *fcc*. We see that it is also located closer to the *fcc* structures in the  
201 latent space. On the other hand, the structure with ABCBCB layers (marked as a purple diamond in  
202 Fig. S8) can be considered structurally more closely related to *hcp* than *fcc*. DeepStruc places this  
203 structure adjacent to *hcp* structures of the same size in the latent space. DeepStruc can thus insert  
204 stacking faulted structures between *fcc* and *hcp* into the latent space in a chemically meaningful way.

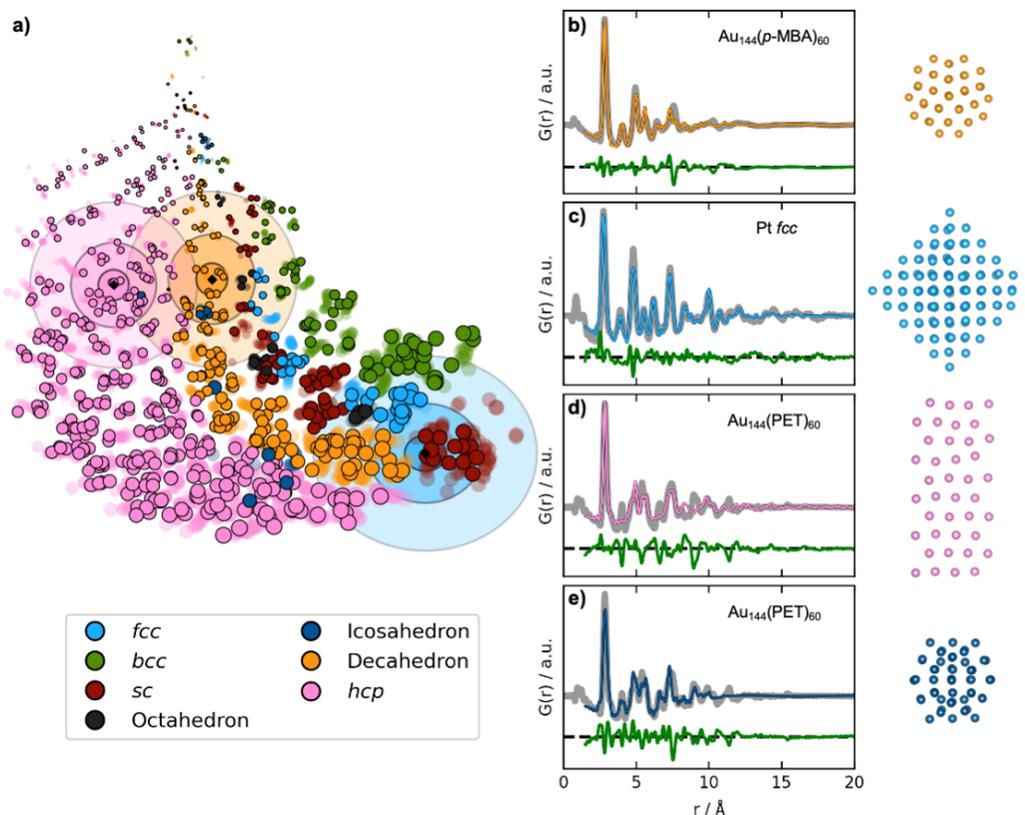


Figure 5: Fitting experimental PDFs with structures obtained by DeepStruc. a) The DeepStruc latent space showing predicted latent space positions for structures from three experimental PDFs. The predicted means are shown as diamond markers, which are enclosed by three rings, indicating the sampling regions for  $\sigma$ : 3, 5, and 7. b) PDF fit of the reconstructed structure from the Au<sub>144</sub>(*p*-MBA)<sub>60</sub> PDF[26] c) PDF fit of the reconstructed structure from the 1.8 nm Pt nanoparticle PDF from Quinson et al.[28], d) PDF fit of the reconstructed structure from the Au<sub>144</sub>(*p*-MBA)<sub>60</sub> PDF[26] using a *hcp* structure. e) PDF fit of the reconstructed structure from the Au<sub>144</sub>(*p*-MBA)<sub>60</sub> PDF[26] using an icosahedral structure. Note that the test set structures shown here are the predicted structures from DeepStruc obtained during inference on experimental PDFs.

205 Fig. S8b illustrates the fits of the reconstructed structures to the PDF data. The difference curves  
 206 indicate that the predicted and true structures are very close to being identical, which is supported by  
 207 the MAE of the atomic positions on  $0.030 \pm 0.019$  Å (section E in the Supplementary Information).  
 208 While disorder causes a broadening of the peaks, the disorder in the generated structures is minor and  
 209 structures with distinct difference between the layers and in the correct sequence can be reconstructed  
 210 to a satisfying degree. This is a promising result, showing that a graph-based CVAE can be used as a  
 211 tool to determine the structure of stacking faulted nanoparticles from PDFs,[31,32] which is a topic  
 212 of significant current interest.[33-37]

### 213 3 Discussion

214 We have shown the potential of using a DGM for structure determination from simulated and  
 215 experimental PDFs. Our graph-based CVAE algorithm, DeepStruc, provides valuable information  
 216 through its latent space, as the MMNP structures cluster based on symmetry and size in agreement  
 217 with their structural chemistry. Using experimental data, the Au<sub>144</sub>(*p*-MBA)<sub>60</sub> nanoparticle was  
 218 determined to be decahedral, Pt nanoparticles were determined to be *fcc* and the Au<sub>144</sub>(PET)<sub>60</sub>  
 219 was determined to have an icosahedral core structure, all in agreement with previous literature. While

220 these systems are relatively simple MMNPs, we recognise that there are more complex materials  
221 where the measured PDF would not contain sufficient information to solve the structure. DeepStruc  
222 would then be limited to provide a distribution of starting models which can aid in the further structure  
223 analysis.

224 Our approach is only restricted by the distribution of the structural training set. When DeepStruc  
225 is trained on *fcc*, *hcp*, and stacking faulted structures, it will locate the stacking faulted structures  
226 in between the *fcc* and *hcp* structures. This suggests a strategy for training DeepStruc models on  
227 different chemical systems that also ‘interpolate’ from one to another when this can be identified.  
228 DeepStruc does not yet provide a completely general structure solution approach, but gives critical  
229 insight into how DGMs can interact with structural and diffraction information to yield candidate  
230 structures and ultimately structure solutions.

231 We suggest to implement DeepStruc as part of PDF-in-the-cloud (PDFitc.org),[38] where the training  
232 data can gradually be expanded over time. So far, the structures investigated are fairly ordered and  
233 contain some symmetry, but in the future, we plan to expand DeepStruc to chemical systems with  
234 more atoms and higher complexity such as metal oxide nanoparticles and alloys. Combining the  
235 PDF conditioning with data from complimentary techniques could prove important for structure  
236 determination of more complex systems. Such studies would both enable structure determination from  
237 a combined modelling perspective, but it would also reveal fundamental aspects of the information  
238 content of the different datasets for solving structure problems.

## 239 **4 Data availability**

240 Code for DeepStruc and the baseline models are available at:

241 <https://github.com/EmilSkaaning/DeepStruc>

242 <https://github.com/AndyNano/Brute-force-PDF-modelling>

243 <https://github.com/AndyNano/MetalFinder>

244 <https://github.com/AndyNano/CVAE>

## 245 **5 Acknowledgements**

246 Acknowledgements This work is part of a project that has received funding from the European  
247 Research Council (ERC) under the European Union’s Horizon 2020 Research and Innovation Pro-  
248 gramme (grant agreement No. 804066). We are grateful to the Villum Foundation for financial  
249 support through a Villum Young Investigator grant (VKR00015416). Funding from the Danish Min-  
250 istry of Higher Education and Science through the SMART Lighthouse is gratefully acknowledged.  
251 We acknowledge support from the Danish National Research Foundation Center for High Entropy  
252 Alloy Catalysis (DNRF 149). Work in the Billinge group was supported by the U.S. National Science  
253 Foundation through grant DMREF-1922234.

## 254 **6 Author contributions**

255 ETSK and ASA contributed to all aspects of the paper. MNW wrote the code associated to the  
256 tree-based classifier. SJLB, RS and KMØJ supervised the project. All authors contributed to the  
257 writing of the manuscript.

## 258 **7 Competing interests**

259 Competing interests The authors declare no competing interests.

## 260 **References**

261 [1] David, W. I. F. & Shankland, K. Structure determination from powder diffraction data. *Acta Crystallogr. A*  
262 **64**, 52-64 (2008).

- 263 [2] Cheetham, A. K. & Goodwin, A. L. Crystallography with powders. *Nat. Mater.* **13**, 760-762 (2014).
- 264 [3] Billinge, S. J. L. & Kanatzidis, M. G. Beyond crystallography: the study of disorder, nanocrystallinity and  
265 crystallographically challenged materials with pair distribution functions. *Chem. Commun.*, 749-760 (2004).
- 266 [4] Young, C. A. & Goodwin, A. L. Applications of pair distribution function methods to contemporary problems  
267 in materials chemistry. *J. Mater. Chem.* **21**, 6464-6476 (2011).
- 268 [5] Christiansen, T. L., Cooper, S. R. & Jensen, K. M. Ø. There's no place like real-space: elucidating size-  
269 dependent atomic structure of nanomaterials using pair distribution function analysis. *Nanoscale Adv.* **2**,  
270 2234-2254 (2020).
- 271 [6] Zhu, H., Huang, Y., Ren, J., Zhang, B., Ke, Y., Jen, A. K.-Y., Zhang, Q., Wang, X.-L. & Liu, Q. Bridg-  
272 ing Structural Inhomogeneity to Functionality: Pair Distribution Function Methods for Functional Materials  
273 Development. *Adv. Sci.* **8**, 2003534 (2021).
- 274 [7] Billinge, S. J. L. & Levin, I. The problem with determining atomic structure at the nanoscale. *Science* **316**,  
275 561-565 (2007).
- 276 [8] Yang, L., Juhas, P., Terban, M. W., Tucker, M. G. & Billinge, S. J. L. Structure-mining: screening structure  
277 models by automated fitting to the atomic pair distribution function over large numbers of models. *Acta*  
278 *Crystallogr. A* **76**, 395-409 (2020).
- 279 [9] Banerjee, S., Liu, C.-H., Jensen, K. M. Ø., Juhas, P., Lee, J. D., Tofanelli, M., Ackerson, C. J., Murray, C. B.  
280 & Billinge, S. J. L. Cluster-mining: an approach for determining core structures of metallic nanoparticles from  
281 atomic pair distribution function data. *Acta Crystallogr. A* **76**, 24-31 (2020).
- 282 [10] Christiansen, T. L., Kjær, E. T. S., Kovyakh, A., Röderen, M. L., Høj, M., Vosch, T. & Jensen, K. M. Ø.  
283 Structure analysis of supported disordered molybdenum oxides using pair distribution function analysis and  
284 automated cluster modelling. *J. Appl. Crystallogr.* **53**, 148-158 (2020).
- 285 [11] Juhás, P., Cherba, D. M., Duxbury, P. M., Punch, W. F. & Billinge, S. J. L. Ab initio determination of  
286 solid-state nanostructure. *Nature* **440**, 655-658 (2006).
- 287 [12] Juhás, P., Granlund, L., Duxbury, P. M., Punch, W. F. & Billinge, S. J. L. The Liga algorithm for ab initio  
288 determination of nanostructure. *Acta Crystallogr. A* **64**, 631-640 (2008).
- 289 [13] Juhas, P., Granlund, L., Gujarathi, S. R., Duxbury, P. M. & Billinge, S. J. L. Crystal structure solution from  
290 experimentally determined atomic pair distribution functions. *J. Appl. Crystallogr.* **43**, 623-629 (2010).
- 291 [14] Cliffe, M. J., Dove, M. T., Drabold, D. & Goodwin, A. L. Structure determination of disordered materials  
292 from diffraction data. *Phys. Rev. Lett.* **104**, 125501 (2010).
- 293 [15] Cliffe, M. J. & Goodwin, A. L. Nanostructure determination from the pair distribution function: a parametric  
294 study of the INVERT approach. *J. Phys.: Condens. Matter* **25**, 454218 (2013).
- 295 [16] Choudhary, K., DeCost, B., Chen, C., Jain, A., Tavazza, F., Cohn, R., WooPark, C., Choudhary, A., Agrawal,  
296 A., Billinge, S. J. L., Holm, E., Ong, S. P. & Wolverton, C. Recent Advances and Applications of Deep Learning  
297 Methods in Materials Science. *arXiv*, 2110.14820 (2021).
- 298 [17] Razavi, A., Van den Oord, A. & Vinyals, O. Generating diverse high-fidelity images with vq-vae-2. *Adv.*  
299 *Neural Inf. Process. Syst.* **32** (2019).
- 300 [18] Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J. & Aila, T. Analyzing and improving the image  
301 quality of stylegan. *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,  
302 8110-8119 (2020).
- 303 [19] Anker, A. S., Kjær, E. T. S., Dam, E. B., Billinge, S. J. L., Jensen, K. M. Ø. & Selvan, R. Characterising the  
304 Atomic Structure of Mono-Metallic Nanoparticles from X-Ray Scattering Data Using Conditional Generative  
305 Models. *Proceedings of the 16th International Workshop on Mining and Learning with Graphs (MLG)* (2020).
- 306 [20] Lim, J., Ryu, S., Kim, J. W. & Kim, W. Y. Molecular generative model based on conditional variational  
307 autoencoder for de novo molecular design. *J. Cheminformatics* **10**, 1-9 (2018).
- 308 [21] Samarakoon, A. M., Barros, K., Li, Y. W., Eisenbach, M., Zhang, Q., Ye, F., Sharma, V., Dun, Z. L., Zhou,  
309 H., Grigera, S. A., Batista, C. D. & Tennant, D. A. Machine-learning-assisted insight into spin ice Dy<sub>2</sub>Ti<sub>2</sub>O<sub>7</sub>.  
310 *Nat. Commun.* **11**, 892 (2020).
- 311 [22] Gómez-Bombarelli, R., Wei, J. N., Duvenaud, D., Hernández-Lobato, J. M., Sánchez-Lengeling, B.,  
312 Sheberla, D., Aguilera-Iparraguirre, J., Hirzel, T. D., Adams, R. P. & Aspuru-Guzik, A. Automatic Chemical  
313 Design Using a Data-Driven Continuous Representation of Molecules. *ACS Cent. Sci.* **4**, 268-276 (2018).
- 314 [23] Marks, L. D. Surface structure and energetics of multiply twinned particles. *Philos. Mag. A* **49**, 81-93  
315 (1984).

- 316 [24] Banerjee, S., Liu, C.-H., Lee, J. D., Kovyakh, A., Grasmik, V., Prymak, O., Koenigsmann, C., Liu, H.,  
317 Wang, L., Abeykoon, A. M. M., Wong, S. S., Epple, M., Murray, C. B. & Billinge, S. J. L. Improved Models for  
318 Metallic Nanoparticle Cores from Atomic Pair Distribution Function (PDF) Analysis. *J. Phys. Chem. C* **122**,  
319 29498-29506 (2018).
- 320 [25] Juhas, P., Farrow, C. L., Yang, X., Knox, K. R. & Billinge, S. J. L. Complex modeling: a strategy and  
321 software program for combining multiple information sources to solve ill posed structure and nanostructure  
322 inverse problems. *Acta Crystallogr. A* **71**, 562-568 (2015).
- 323 [26] Jensen, K. M. Ø., Juhas, P., Tofanelli, M. A., Heinecke, C. L., Vaughan, G., Ackerson, C. J. & Billinge, S. J.  
324 L. Polymorphism in magic-sized Au<sub>144</sub>(SR)<sub>60</sub> clusters. *Nat. Commun.* **7** (2016).
- 325 [27] Banerjee, S., Liu, C.-H., Jensen, K., Juhás, P., Lee, J. D., Tofanelli, M., Ackerson, C. J., Murray, C. B. &  
326 Billinge, S. J. L. Cluster-mining: an approach for determining core structures of metallic nanoparticles from  
327 atomic pair distribution function data. *Acta Crystallogr. A* **76**, 24-31 (2020).
- 328 [28] Quinson, J., Kacenauskaite, L., Christiansen, T. L., Vosch, T., Arenz, M. & Jensen, K. M. Ø. Spatially  
329 Localized Synthesis and Structural Characterization of Platinum Nanocrystals Obtained Using UV Light. *ACS*  
330 *Omega* **3**, 10351-10356 (2018).
- 331 [29] Yan, N., Xia, N., Liao, L., Zhu, M., Jin, F., Jin, R. & Wu, Z. Unraveling the long-pursued Au<sub>144</sub> structure  
332 by x-ray crystallography. *Sci. Adv.* **4**, eaat7259 (2018).
- 333 [30] Bertolotti, F., Moscheni, D., Migliori, A., Zacchini, S., Cervellino, A., Guagliardi, A. & Masciocchi, N. A  
334 total scattering Debye function analysis study of faulted Pt nanocrystals embedded in a porous matrix. *Acta*  
335 *Crystallogr. A* **72**, 632-644 (2016).
- 336 [31] Masadeh, A. S., Bozin, E. S., Farrow, C. L., Paglia, G., Juhas, P., Billinge, S. J. L., Karkamkar, A. &  
337 Kanatzidis, M. G. Quantitative size-dependent structure and strain determination of CdSe nanoparticles using  
338 atomic pair distribution function analysis. *Phys. Rev. B* **76** (2007).
- 339 [32] Yang, X., Masadeh, A. S., McBride, J. R., Božin, E. S., Rosenthal, S. J. & Billinge, S. J. L. Confirmation  
340 of disordered structure of ultrasmall CdSe nanoparticles from X-ray atomic pair distribution function analysis.  
341 *Phys. Chem. Chem. Phys.* **15**, 8480-8486 (2013).
- 342 [33] Cenker, J., Sivakumar, S., Xie, K., Miller, A., Thijssen, P., Liu, Z., Dismukes, A., Fonseca, J., Anderson, E.,  
343 Zhu, X., Roy, X., Xiao, D., Chu, J.-H., Cao, T. & Xu, X. Reversible strain-induced magnetic phase transition in  
344 a van der Waals magnet. *Nat. Nanotechnol.* (2022).
- 345 [34] Rong, X., Liu, J., Hu, E., Liu, Y., Wang, Y., Wu, J., Yu, X., Page, K., Hu, Y.-S., Yang, W., Li, H., Yang,  
346 X.-Q., Chen, L. & Huang, X. Structure-Induced Reversible Anionic Redox Activity in Na Layered Oxide  
347 Cathode. *Joule* **2**, 125-140 (2018).
- 348 [35] Charles, D. S., Feygenson, M., Page, K., Neufeind, J., Xu, W. & Teng, X. Structural water engaged  
349 disordered vanadium oxide nanosheets for high capacity aqueous potassium-ion storage. *Nat. Commun.* **8**,  
350 15520 (2017).
- 351 [36] Gao, P., Metz, P., Hey, T., Gong, Y., Liu, D., Edwards, D. D., Howe, J. Y., Huang, R. & Mixture, S. T. The  
352 critical role of point defects in improving the specific capacitance of -MnO<sub>2</sub> nanosheets. *Nat. Commun.* **8**,  
353 14559 (2017).
- 354 [37] Metz, P. C., Koch, R. & Mixture, S. T. Differential evolution and Markov chain Monte Carlo analyses of  
355 layer disorder in nanosheet ensembles using total scattering. *J. Appl. Crystallogr.* **51**, 1437-1444 (2018).
- 356 [38] Yang, L., Culbertson, E. A., Thomas, N. K., Vuong, H. T., Kjaer, E. T. S., Jensen, K. M. Ø., Tucker, M. G.  
357 & Billinge, S. J. L. A cloud platform for atomic pair distribution function analysis: PDFfitc. *Acta Crystallogr. A*  
358 **77**, 2-6 (2021).