

ARTICLE OPEN



Uncovering material deformations via machine learning combined with four-dimensional scanning transmission electron microscopy

Chuqiao Shi¹, Michael C. Cao^{1,2}, Sarah M. Rehn³, Sang-Hoon Bae^{4,5}, Jeewan Kim⁶, Matthew R. Jones^{1,3}, David A. Muller^{1,2,7} and Yimo Han¹✉

Understanding lattice deformations is crucial in determining the properties of nanomaterials, which can become more prominent in future applications ranging from energy harvesting to electronic devices. However, it remains challenging to reveal unexpected deformations that crucially affect material properties across a large sample area. Here, we demonstrate a rapid and semi-automated unsupervised machine learning approach to uncover lattice deformations in materials. Our method utilizes divisive hierarchical clustering to automatically unveil multi-scale deformations in the entire sample flake from the diffraction data using four-dimensional scanning transmission electron microscopy (4D-STEM). Our approach overcomes the current barriers of large 4D data analysis without a priori knowledge of the sample. Using this purely data-driven analysis, we have uncovered different types of material deformations, such as strain, lattice distortion, bending contour, etc., which can significantly impact the band structure and subsequent performance of nanomaterials-based devices. We envision that this data-driven procedure will provide insight into materials' intrinsic structures and accelerate the discovery of materials.

npj Computational Materials (2022)8:114; <https://doi.org/10.1038/s41524-022-00793-9>

INTRODUCTION

Recent advances in the synthesis of materials have led to well-behaved material structures at the nanometer scale. Atomic structure and deformations in these nanomaterials determine their chemical, electronic, and optical properties, affecting their efficiency and performance in their targeted applications. For example, epitaxial growth of optically tunable heterostructures in III–V semiconductors^{1–3} and two-dimensional (2D) materials^{4–8} can lead to local strain and dislocations, which greatly affect the electronic and optical properties due to changes in the local band structure. Even within nanomaterials comprised of a single crystal, such as structurally designed anisotropic metallic nano-plates or prisms^{9–11}, the local bending or deformations also play an important role in determining their optical responses and catalytic behaviors. To study these fine features in the materials, conventional high-resolution transmission electron microscopy and annular dark field scanning transmission electron microscopy (ADF-STEM) have been utilized to reveal the local atomic structure^{2,3,7}. While the limits of the electron microscope resolution are constantly being pushed¹², a common restriction on these imaging techniques is the limited field of view in the sample. As studying the overall structural information of the materials is crucial for mass production and large-scale processing for applications in next-generation devices, techniques such as nanobeam electron diffraction^{13–15} that can map large sample areas have attracted tremendous attention for their potential in determining the sample structure on a large scale.

Although nanobeam electron diffraction has been used for decades to acquire electron diffraction patterns from a large sample area, this technique was limited by conventional charge-

coupled device detectors, which are too slow to collect detailed structural information across the entire sample. However, the development of fast direct electron detectors¹⁶ now allows the collection of a momentum-resolved nanobeam diffraction pattern at each scanning position in a STEM experiment (Fig. 1a), thus generating four-dimensional (4D) data (Fig. 1b). Therefore, the technique enabled by these detectors is often referred to as 4D-STEM¹⁷. Among the direct electron detectors, the electron microscope pixel array detector (EMPAD)¹⁸ has high dynamic range and sensitivity, capable of recording quantitative diffraction patterns without saturating the center beam or cutting off the weak diffracted spots. Using the EMPAD in scanning nanobeam diffraction mode, strain profiles in 2D materials have been mapped across a micrometer scale with sub-picometer precision¹⁹. Although many works^{14,20–23} have been reported to further improve accuracy and precision, this approach relies heavily on prior knowledge of the sample structure. For instance, before applying any strain or phase mapping approach, one needs to determine where and how to place masks on the diffraction patterns. This process varies from sample to sample. Acquiring this information from a large quantity of data from diverse samples makes it difficult to generalize 4D-STEM for materials with unexpected lattice deformations, which usually have a large impact on material properties and device performance.

Recently, machine learning has emerged as a promising method applied in microscopy^{24–30} due to its capability in analyzing complex patterns in large datasets. Specifically, unsupervised learning, which does not require training data, has been utilized to identify the stacking order³¹ and twin boundaries³² in materials. To further extend unsupervised learning for deformation and fine

¹Department of Materials Science and NanoEngineering, Rice University, Houston, TX, USA. ²School of Applied and Engineering Physics, Cornell University, Ithaca, NY, USA.

³Department of Chemistry, Rice University, Houston, TX, USA. ⁴Department of Mechanical Engineering and Materials Science, Washington University in Saint Louis, Saint Louis, MO, USA. ⁵Institute of Materials Science and Engineering, Washington University in Saint Louis, Saint Louis, MO, USA. ⁶Department of Mechanical Engineering, Massachusetts Institute of Technology, Cambridge, MA, USA. ⁷Kavli Institute for Nanoscale Science, Cornell University, Ithaca, NY, USA. ✉email: yimo.han@rice.edu

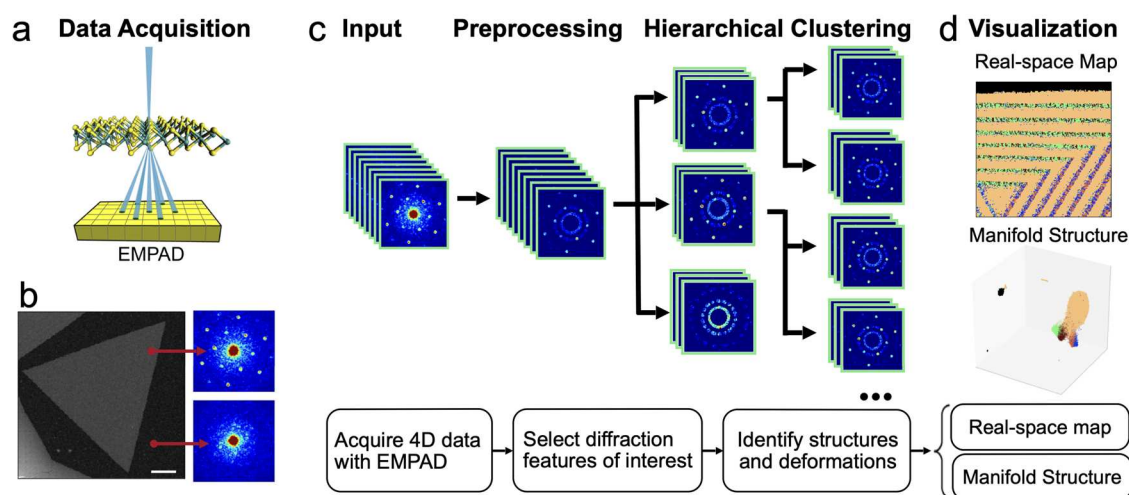


Fig. 1 Schematic of unsupervised learning of 4D-STEM datasets. **a** Schematic of the EMPAD operation, where a diffraction pattern is recorded at each scanning position. **b** A 4D dataset that contains a full diffraction pattern at each pixel in the real-space map. Scale bar, 500 nm. **c** Schematic of the divisive hierarchical clustering workflow on diffraction patterns in the 4D data. **d** Visualization of the clustering results, where the features are mapped back to real space (top) to determine their distribution, and the diffraction patterns are projected onto a lower-dimensional manifold space (bottom) for visualization.

structure study, we utilized a divisive hierarchical unsupervised clustering architecture for rapid and semi-automatic 4D data analysis and feature mapping according to intrinsic characteristics and similarity (Fig. 1c). This approach complements the existing strain mapping approaches by providing a rapid and automatic initial analysis of the 4D data (usually < 10 min), as well as uncovering unexpected but significant fine structures and deformations in the materials.

RESULTS

Architecture of divisive hierarchical clustering

The process's overall scheme is described in Fig. 1. After the data acquisition, three main steps are used to process the 4D data: preprocessing the diffraction patterns, hierarchical clustering of the data, and visualizing the results.

Preprocessing involves aligning and masking the diffraction patterns. Alignment corrects the drift in the diffraction patterns caused by slightly misaligned beam tilt when scanning large areas (Supplementary Fig. 1). To do that, we align the center of mass (CoM) of the center beam. We added a circle mask in the center of each diffraction pattern to select the center beam (Supplementary Fig. 1b). Then, the CoM of the center beam was calculated with sub-pixel precision, followed by moving the diffraction pattern towards the center of the detector (Supplementary Fig. 1a). We repeated these three steps until the standard deviation of the CoM in all scanning positions, as well as the error between the CoM and the center of the detector, becomes smaller than 0.01 pixel. This step avoids the confusion caused by the translational shift of the diffraction patterns due to reasons other than the intrinsic sample structure.

Masking contains two parts, where the first part uses a ring mask to mask out the low-angle scattering as well as the zero-padded regions caused by the alignment. We determined the inner and outer radius through plotting the standard deviation (STD) of diffraction pattern. In the example data (Supplementary Fig. 2a, e), the bright center beams always show highest STD due to the high intensity. From the rotational STD plots (Supplementary Fig. 2b, f), we were able to identify the diffraction area and set the inner/outer radius to block the background and the noise (Supplementary Fig. 2c, g). The second part selects the diffraction regions

where the crystal information is stored. In this part, we select the features based on the STD of the intensities in the remaining area. The regions with high STD (30% of the highest) among the ring-masked dataset are selected in each diffraction pattern (Supplementary Fig. 2d, h) and flattened as a feature vector, which is used in the following cluster analysis. The optimal percentage (30%) was determined by experience, as shown in Supplementary Fig. 3. Thus far, this threshold applies for both thin 2D materials (Supplementary Fig. 3b) and thicker nanoprisms (Supplementary Fig. 3e). As a result, the preprocessing step helps organize the 4D datasets so that our method will mainly extract structural information from the datasets, while avoiding any ambiguity caused by the microscope misalignment and background.

The hierarchical clustering architecture is illustrated in Fig. 1c. A single round of clustering is insufficient to determine all the features in the 4D dataset due to features at varying length scales. Initial clustering separates large scale features, typically different materials in the sample, since they cause more noticeable differences in the diffraction pattern. However, diffraction patterns within these initial clusters then have more subtle differences caused by small-scale features, which cannot be separated through a single clustering round. To extract these features at different scales, we employ a divisive hierarchical clustering architecture to cluster the diffraction patterns in the 4D dataset. The divisive architecture is a top-down approach, which starts from the whole dataset, and then clustering is performed recursively when moving down the hierarchy.

In a single round of clustering, we compared different unsupervised learning methods and identified that K-means³³ exhibited the optimal performance considering combined accuracy and speed (Supplementary Table 1 and Supplementary Fig. 4). The K-means algorithm is a common clustering method that divides our real-space points, each represented by a feature vector, into K clusters through minimizing the variance (squared Euclidean distance) within each cluster. To automatically determine the number of clusters, or the K , especially when little prior knowledge about the datasets was presented, we utilize the *elbow method*³⁴ to determine the K number in each round clustering (Supplementary Fig. 5). In the *elbow method*, we calculated the total *within-cluster sum of squares* (WSS) and plotted the error curves according to the number of clusters (or K). Since we do not

need precise WSS values, we adopted the faster but negligibly less precise Mini-Batch K-means method (Supplementary Fig. 4e) to determine the *elbow point*.

Finally, since each diffraction pattern represents one pixel in real space, the labels from the clustering results can be mapped back to real space for a better understanding of the material structures. As shown in the top panel of Fig. 1d, the clustering result is visualized in a real-space color-coded map, with each color representing a single structural feature characterized in a cluster of similar diffraction patterns. Alternatively, we can also visualize the data distribution in a low-dimensional manifold structure (Fig. 1d, bottom panel). The manifold structure is the projection of the high dimensional diffraction patterns to 3D space for visualization, which is achieved by the recently developed uniform manifold approximation and projection (UMAP)³⁵. The lower dimensional manifold structures can assist in the visualization of the distribution and variance of the clusters, which can be used to determine the differences in major physical parameters. Through real-space and manifold structure visualization after clustering, we can map the fine deformations to better understand the distribution of structural features in the sample.

Deformations in WS₂–WSe₂ lateral heterojunction

To test the accuracy of our clustering architecture, we applied our approach to 4D datasets of 2D epitaxial lateral heterojunctions, which contain strain-engineered structures that enable tunable optical properties¹⁹. The sample is composed of an outer and inner region of WS₂ and a middle region of WSe₂, which is not apparent in its ADF-STEM image (Fig. 2a) as the tungsten atoms dominate the ADF contrast. During the clustering, structural features at different scales within the entire flake have been uncovered hierarchically. After three rounds, as shown in the dendrogram (Fig. 2b), fine deformations have been uncovered in the sample.

The real space map from the first two rounds of clustering is shown in Fig. 2c, where the method effectively separated the background, WS₂, and WSe₂ (Supplementary Fig. 6a–c). The crystal-line samples are distinguished from the amorphous substrate in the first round of clustering, which is intuitively obvious. In the second round, the WSe₂ and WS₂ are differentiated by the diffraction spots spacing caused by the lattice mismatch between WS₂ and WSe₂ (Supplementary Fig. 6d). The real space map from unsupervised learning provides a precise interface between WS₂ and WSe₂ with defects (Supplementary Fig. 7), which were not recognized in the ADF images. The two separated clusters of WS₂ and WSe₂ in the manifold structure (Fig. 2f) indicate that the lattice constant changes across the junction are discrete. We measured the lattice constant of each diffraction pattern in the junction sample, and the histogram (Fig. 2i) confirms the discrete change. Using the measured lattice constant as the ground truth, our clustering results provide a 99.9% accuracy for the discrete feature.

In the third round of clustering, the sub-clusters of WS₂ and WSe₂ from the second round are analyzed to reveal finer features. Our method separates the WS₂ cluster into two sub-clusters, where the real space map shows a rotational periodicity of the material with graded interfaces (Fig. 2d). The mean diffraction patterns (Supplementary Fig. 6e, f) of each sub-cluster display a slight rotation in the reciprocal lattice (Supplementary Fig. 6g), which corresponds to different lattice rotations in real space. Unlike the discretely separated WS₂ and WSe₂ clusters from the previous round, these two clusters are mixed in the manifold representation (Fig. 2g), indicating that the rotation angle changes continuously. We measured the rotation angle of each diffraction pattern in WS₂, and the histogram (Fig. 2j) confirms the capability of our method to identify continuous structural distortions in materials. The measured accuracy of clustering this continuous feature is 84%.

Meanwhile, the WSe₂ cluster is separated into four clusters, and the real space map is shown in Fig. 2e. The averaged diffraction

patterns of each cluster differ in the intensity of the second order spots (Supplementary Fig. 6h–k), caused by lattice tilts (or ripples) in the sample^{36,37}. Specifically, three separate regions in the corners show different directional ripples since the second-order diffraction spots along one direction display a much stronger intensity than the other two directions (Supplementary Fig. 6i–k). In contrast, the area close to the center is a flat region where all the second-order spots have similar intensity (Supplementary Fig. 6h). The 3D manifold structure of the WSe₂ cluster (Fig. 2h) shows that the tilts in the ripple area are continuous. The results we achieved here are consistent with literature where the ripples form to relax the strain induced by WS₂–WSe₂ lattice mismatch¹⁹.

Uncovering minor ripples in WS₂–WSe₂ superlattices

To test how sensitive our method is for minor deformations in materials, we utilized a coherent 2D superlattice sample which contains a much smaller ripple structure with an aspect ratio of ~1/30 in narrow WSe₂ stripes⁸. Using our method on the entire superlattice sample (Fig. 3a), structural features at different scales were uncovered hierarchically. Our method identified different flakes in the first round of clustering (Fig. 3b). Then lattice differences between WS₂ and WSe₂ were illustrated in the second round (Fig. 3c). The next few rounds unveiled directional uniaxial strain (Fig. 3d and Supplementary Fig. 8a–e), as well as minor ripples hidden in the WSe₂ (Fig. 3e and Supplementary Fig. 8f–k). Due to the coherency in the superlattice, the ripples presented here are much smaller than what we observed in the WS₂–WSe₂ lateral heterojunction sample, indicating the capability of our hierarchical approach for identifying minor deformations in materials. Finally, the manifold structures of the zoomed-in dataset from each round of clustering provide a view of each cluster and subcluster in data space (Supplementary Fig. 9).

Uncovering bending contours in silver nanoprisms

Unlike 2D materials, thicker samples usually present more complex deformations and lattice distortions, which result in more challenges in understanding their local structure. We use our method to investigate the deformations in silver nanoprisms, which have been widely studied due to their optical properties that show great potential in many applications^{9–11}. The silver nanoprisms were drop casted on an amorphous carbon supporting film followed by air drying, which introduced internal deformations due to the surface tension. The ADF-STEM (Fig. 4a) shows an entire flake of a micrometer-sized silver nanoprism, which is known to be single crystal with the [111] zone axis perpendicular to the flat surfaces. Here, we show that our machine learning approach can elucidate deformations in nanoprism flakes in a facile and semi-automated way (Supplementary Fig. 10).

First, the same preprocessing of 4D data was performed on the nanoprism datasets (Supplementary Fig. 11). Focusing on a corner of the nanoprism (Fig. 4b), the first round of clustering distinguishes the sample from the amorphous supporting film (Supplementary Fig. 12a) due to the obvious difference between the amorphous film (Supplementary Fig. 12b) and the nanoprism (Supplementary Fig. 12c) in the diffraction patterns. The following round of clustering identifies the contours in the nanoprism flake (Supplementary Fig. 12d–f). To further investigate the deformation, we clustered the contour data into smaller sub-clusters and identified the two sides of the contour (yellow and orange in Fig. 4c), which is consistent with the general structure of a bending contour. From the averaged diffraction patterns (Supplementary Fig. 12h, i), the difference between the two sides of the bending contour in diffraction space is the intensity variation in the conjugate diffraction spots, which can be explained as the incident electron beam approaching the atomic planes at different angles at two different bending sides. On each side, the Bragg condition in one direction is fulfilled, and the diffraction of the

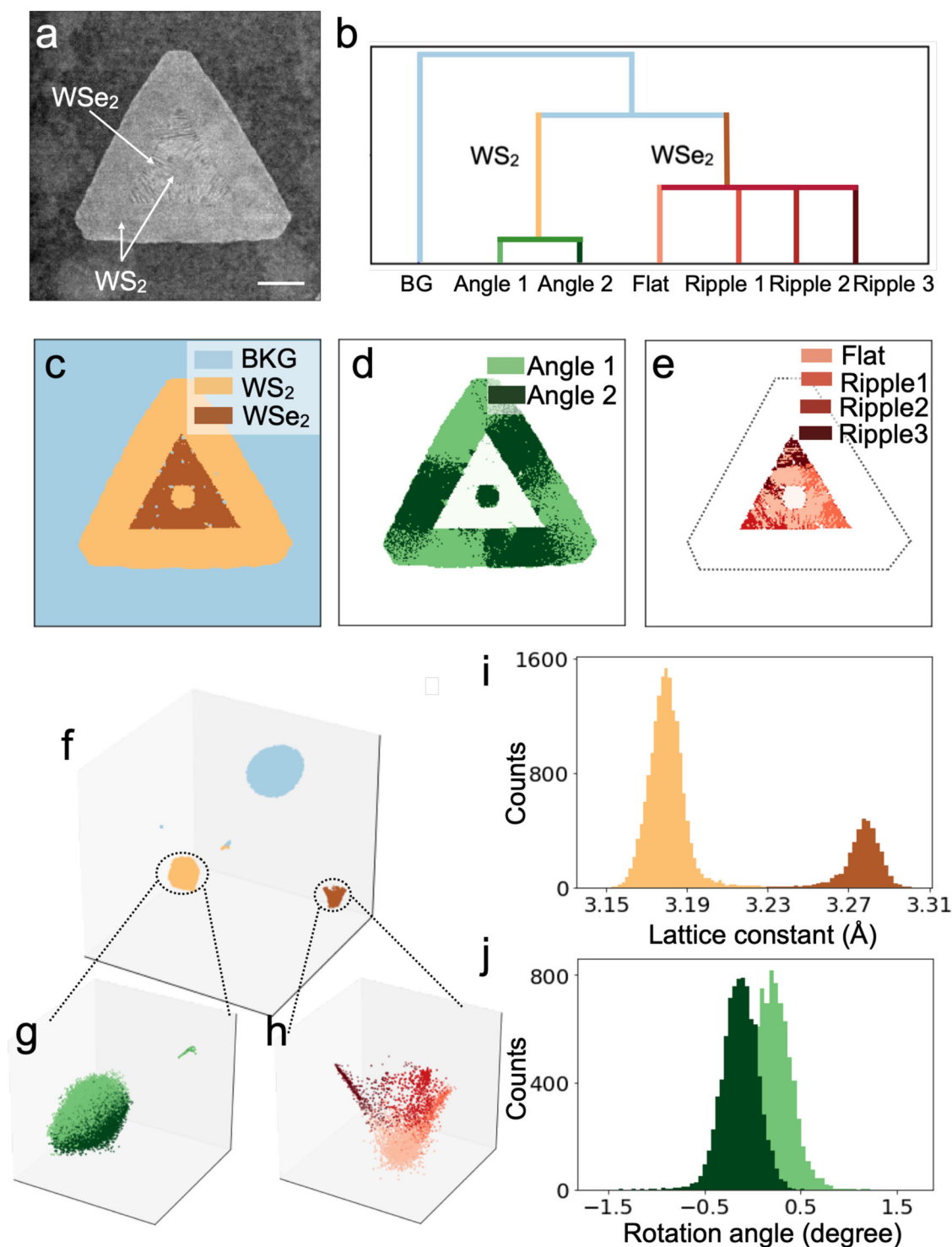


Fig. 2 Clustering results on WS_2 - WSe_2 lateral heterojunction. **a** ADF-STEM image of the WS_2 - WSe_2 junction. Scale bar, 500 nm. **b** Dendrogram of the unsupervised learning workflow. **c–e** Real space maps from the unsupervised learning results showing different compositions and deformations in the sample. **f–h** 3D manifold structure displaying the data distribution in each round of the clustering. **i** Histogram of measured lattice constants of WS_2 and WSe_2 color-coded with results in the second-round clustering. **j** Histogram of measured rotation angles of WS_2 color-coded with results in the third-round clustering of WS_2 .

incident beam would result in strong corresponding diffraction peaks (Supplementary Fig. 13a). By placing masks on the conjugate diffraction spots to form virtual dark-field (DF) images (Supplementary Fig. 13b), the high contrast contours in the virtual

DF images are consistent with the previously reported experimental DF-TEM images in silver nanoprisms. In addition, our simulated diffraction patterns (Supplementary Fig. 14) suggest that the thickness variations can be detected due to the distinct

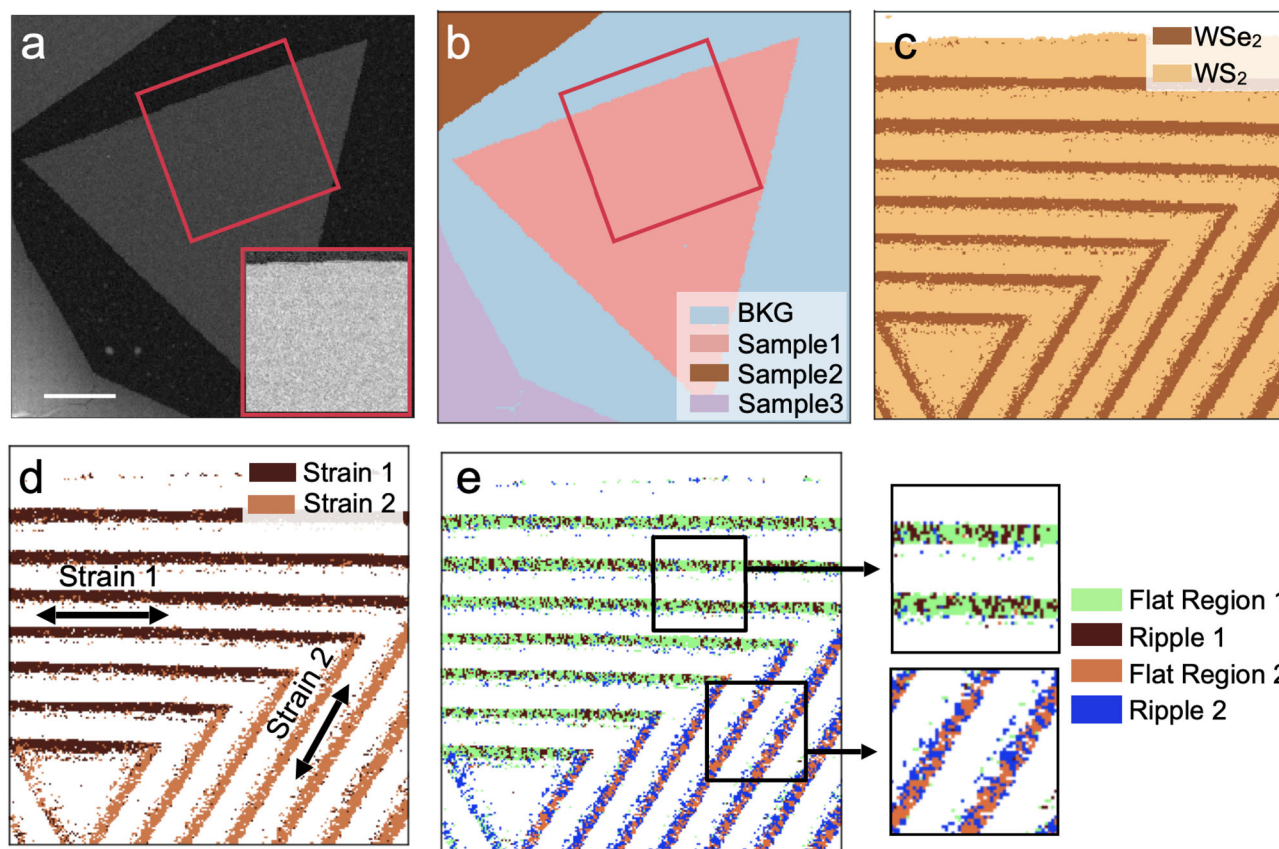


Fig. 3 Clustering results on WS_2 - WSe_2 superlattices. **a** ADF-STEM images of the 2D multi-junction superlattice with a zoomed-in area shown in the inset. The ADF contrast is dominated by heavy tungsten atoms, thus no superlattice structure appears. Scale bar, 500 nm. **b** Real-space map of the clustering results from the 4D dataset of the WS_2 - WSe_2 multi-junction sample, where different flakes are identified from the substrate. **c** Zoomed-in real-space maps of the multi-junction sample area shown in the red box in **(a)**. **d** Real-space map of the sub-clustering results in WSe_2 with two colors indicating different strain profiles. **e** Real-space map of another clustering round on each subcluster in **(d)**, providing more structural details of lattice ripples in WSe_2 due to the strain.

diffraction intensity, confirming the uniformity of the thickness in our Ag nanoprism (10 nm from previous publication¹¹).

In the manifold structure (Fig. 4d–f), the data is first segmented into two main clusters (Fig. 4d), amorphous background (black) and nanoprism (dark purple). Further sub-clustering on the nanoprism cluster separated the ring-like appendage (bending contour) from the original sample manifold as the new-subcluster (Fig. 4e). The thin asymmetric shape of the bending contour manifold (Fig. 4f) indicates a continuous deformation, which is shown in the diffraction pattern as an intensity change in the two conjugate diffraction spots.

Clustering real space images for virtual imaging

Conventionally, we represent 4D-STEM datasets in real-space major order (x, y, k_x, k_y) . However, we can also view the data in momentum-space major order (k_x, k_y, x, y) . Instead of a 2D array of diffraction patterns, we can visualize the 4D dataset as a 2D array of real-space images generated from diffraction pattern intensities at a single momentum-space coordinate (Fig. 5a). Due to this property of 4D datasets, our hierarchical method could also be extended to clustering these real-space images (Fig. 5b), aimed at analyzing the data from a different dimension. However, the quantitative intensity of real-space images varies dramatically across the diffraction space. Consequently, the clustering results are dominated by the intensity effect, overpowering the actual feature information shown in the contrast (Supplementary Fig. 15a). To uncover the actual features, we normalized the real space

image as an additional preprocessing step before clustering (Supplementary Fig. 15b), which provides similar intensity in the real space images for better recognition of sample features using our method. Followed by the hierarchical method, the real-space features in different momentum-space pixels can be uncovered.

We re-investigate 4D data of the WS_2 - WSe_2 superlattices using this clustering approach to segment the diffraction space into multiple clusters (Fig. 5c). This method works by effectively placing virtual objective apertures in the diffraction space, which can accurately select the diffraction coordinates based on their generated real-space image similarity. In each cluster, the averaged real-space image shows results equivalent to different modes in STEM, including virtual bright-field (BF) (Fig. 5d) and DF (Fig. 5e, f) images. The DF images of different flakes are separated (Fig. 5e, f) according to their different lattice orientations (cyan and yellow in Fig. 5c). Compared with the hierarchical clustering on diffraction patterns (Fig. 3), the results from clustering real-space images provide crystal information that is stored in momentum-space. However, we observed striping in the BF image (Fig. 5d) generated from pixels in the center beam. Upon further investigation of the data (Supplementary Fig. 15c, d), we found the stripes are caused by the center beam alignment in the pre-processing of the 4D data, which indicates the importance of more accurate alignment of the beam tilt for 4D nanobeam electron diffraction mode in STEM.

The virtual imaging through unsupervised learning can not only be used on 2D materials but also bulk materials. We investigated

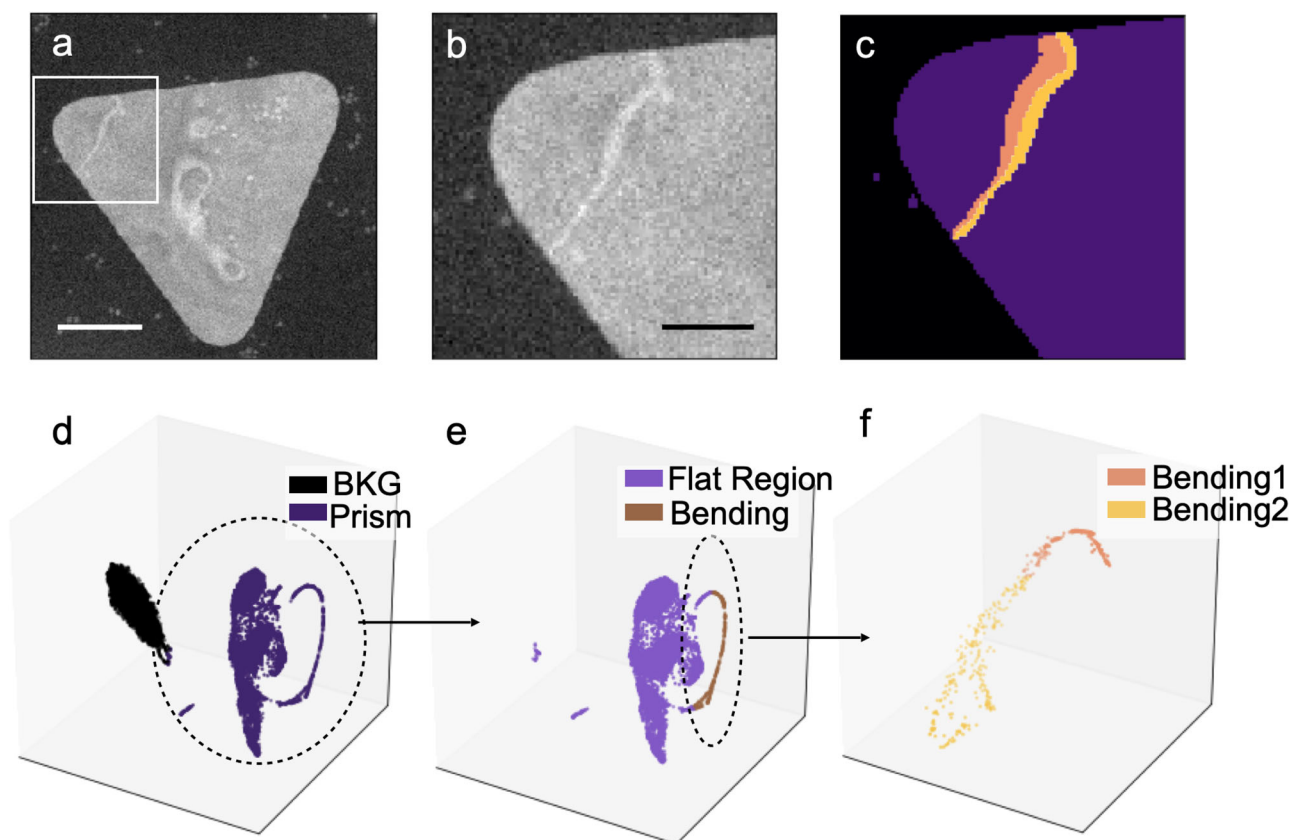


Fig. 4 Clustering results of silver nanoprisms. **a** ADF-STEM image of a silver nanoprism placed on amorphous SiN_x supporting film. Scale bar, 500 nm. **b** Zoomed-in ADF image from the white box in (a). Scale bar, 200 nm. **c** Real space map of the final cluster and sub-cluster results. **d–f** Manifold structure of the data in each hierarchical clustering process.

the epitaxial InGaP grown on GaAs, where the complexity of strain and dislocations in the InGaP layer create challenges to quantitatively understand the material structure from conventional methods (ADF-STEM image in inset of Fig. 5h). As the strain profile in the material dramatically affects its properties and performance, uncovering the complicated crystalline structure in the film becomes crucial. For such complex materials, our real-space clustering approach segments the reciprocal space into six parts. A BF image (Fig. 5i) is generated from the segment where the center beam is located (green part in Fig. 5h). A low-angle ADF (LAADF) image (Fig. 5j) is reconstructed from low scattering angles between the center beam and first order diffraction spots (dark blue part in Fig. 5h), which highlights the amorphous carbon protecting layer on top of the film due to the amorphous scattering ring formed at this low angle. In addition, higher angle areas that exclude the diffraction spots (light blue part of Fig. 5h) display elastically scattered electrons that roughly represent the thickness of the sample (Fig. 5k). The remaining three clusters indicate the crystallinity of the thin film, providing virtual DF images of three lattice orientations (Fig. 5l–n). From the diffraction map (Fig. 5h), we conclude that Fig. 5l, m show small tilts possibly caused by the lattice strain in InGaP, while Fig. 5n displays a twin domain formed in the film.

DISCUSSION

We intend to compare different dimension reduction methods and determine the most efficient one for our application. The EMPAD records 128×128 diffraction patterns so each pattern has 16834 pixels. However, the useful structural information is

concentrated only in the diffraction spots, which is a small part in the whole pattern. A condensed representation of the data significantly improves computational speed and feature selection. In the preprocessing step, we chose to place STD masks on diffraction patterns to reduce the data size. However, there are more elegant approaches to reduce the dimension of the data such as Principal Component Analysis (PCA) and UMAP, but they are not desirable for our hierarchical clustering method.

PCA computes the principal components based on matrix factorization but since it is applied to the whole dataset, only the main difference is captured while the minor differences are ignored. For example, when applying the hierarchical clustering on the PCA components of the WSe₂–WS₂ junction in Fig. 2, the WSe₂ and WS₂ samples can be separated in the first two rounds (Supplementary Fig. 16b). However, in the next round of clustering, the rotations cannot be captured (Supplementary Fig. 16c) unless we re-apply PCA on each sub-cluster before each round clustering, which takes more time and memory. The need to re-apply dimension reduction on each sub-cluster makes PCA undesirable for our clustering method.

The other feature extraction method we tested was UMAP, which approximates the features in low-dimensional manifold space. Compared with PCA, UMAP can capture both the major and minor features in the sample. However, the clustering results show striping artifacts (Supplementary Fig. 16d), while the real deformation features are hidden. Due to this artifact, we do not use UMAP as a dimension reduction method but only as a visualization method.

In addition, there are alternative clustering methods besides K-Means, including Agglomerative clustering, BIRCH clustering, Spectral clustering, etc. To quantitatively compare the efficiency and performance of different clustering methods, we created the ground truth

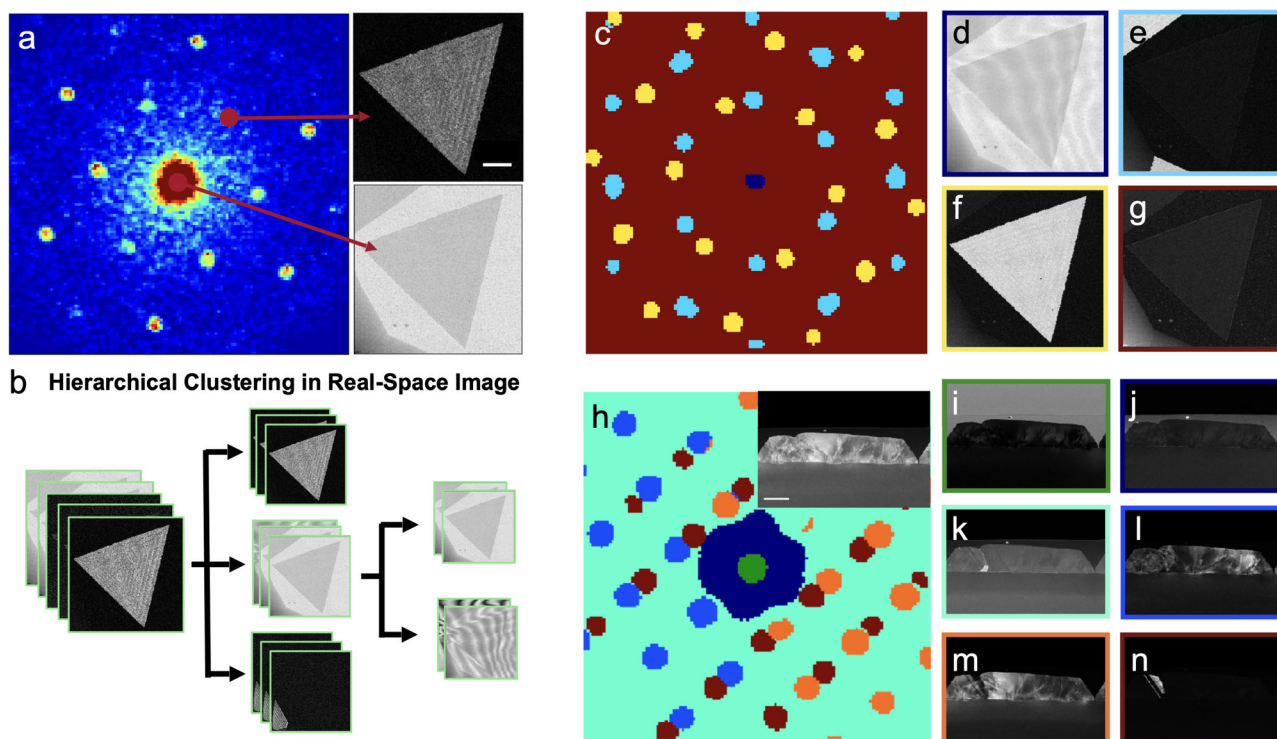


Fig. 5 Clustering results on real space images in 4D-STEM datasets. **a** Visualizing 4D data in a momentum-major order, where each pixel in diffraction pattern can be considered as a real-space image. Scale bar, 500 nm. **b** Schematic of the divisive hierarchical clustering architecture on real space images. **c** Map of hierarchical clustering results in diffraction space on a $\text{WS}_2\text{-WSe}_2$ superlattice. **d–g** Mean real-space images of the superlattice in each cluster. **d** Displays the image from dark blue area in **(c)**, which represents a virtual BF image; **e, f** are the images from light blue and yellow portions in **(c)**, corresponding to DF images for different flakes. **g** Sums all other areas in **(c)**, indicating a thickness variation in the sample. **h** Map of hierarchical clustering results in diffraction space on a cross-sectional InGaP/GaAs crystal, with an ADF image displayed in the inset. Scale bar, 200 nm. **i–n** Mean real-space images of the sample in each cluster. **i** corresponds to the center beam, the green area in **(h)**, and shows a virtual BF image; **j, k** are from the dark blue and cyan area in respect, which are from the amorphous carbon and thickness effect in the sample; **l–n** are virtual DF images for the cross-section, showing strain effects and twin grains.

label on the WS_2 cluster from the manually measured rotation map (Supplementary Fig. 4a). We set the positive rotation angle as label 1 and the negative rotation angle as label 2. Then we cluster this dataset with different methods, record the time and calculate the accuracy. Based on the results shown in Supplementary Fig. 4, the Mini-batch K-means method is the fastest. However, due to the random initialization of this algorithm, the clustering results are not stable. The spectral clustering provides the best accuracy but takes a much longer time. Considering the tradeoff between time and performance, K-means was chosen as the clustering method in each round of the hierarchical clustering architecture.

Here, we summarize the combination of different dimension reduction and clustering methods on the WS_2 dataset. The time and accuracy of each method are shown in Supplementary Table 1, which proved the STD selection and K-means are suitable for our hierarchical clustering architecture.

Finally, Our method demonstrates a quick initial analysis of 4D-STEM data and allows for easier and quicker discovery of unexpected deformations in a crystalline sample. However, our approach does not focus on providing highly accurate maps of such deformations, due to the limitation of clustering data. The method provided high accuracy (>99%) for discrete features, but only showed a reasonable accuracy (>80%) to distinguish a continuous deformation (Supplementary Table 1). Nevertheless, it is sufficient to inform users of the existence of such continuous features and provide a quick initial analysis of the 4D data, but further processing of the 4D data is required to quantitatively map them.

We also evaluate the precision (or consistency) of the K-means approach by running the algorithm 100 times on the same dataset with random initial clustering centers. The testing datasets contain continuous features. The result (Supplementary Fig. 17) shows that 88% of the total runs possess an accuracy above 75% for this continuous feature, indicating the reliability of the approach.

In summary, we have demonstrated a method using divisive hierarchical unsupervised machine learning to perform initial analysis of 4D-STEM datasets, which may accelerate and automate the study of lattice deformations in materials. We have applied this method to extract such features from different 2D lateral heterojunctions, thicker 3D materials, and cross-sectional crystals. The purely data-driven analysis uncovers different types of material deformations in the samples, such as strain, lattice distortion, bending contour, etc. Combined with highly accurate mapping techniques, our method may lead to a crucial step towards a fully autonomous method for the analysis of subtle lattice deformations. In addition, this approach may be potentially expanded to broader material systems or other imaging techniques that generate large and multidimensional datasets, benefiting the development of future materials, techniques, and applications.

METHODS

EMPAD data acquisition

All 4D-STEM datasets were acquired on an aberration-corrected Thermo Fisher Titan Themis. The 4D datasets of the 2D lateral heterojunction were acquired at 80 kV. A 1.6 mrad convergence angle was used, leading to a

~3 nm probe size. More experimental details can be found in our previous work¹⁹. The 4D datasets collected on the silver nanoprism sample and cross-sectional semiconductors were acquired at 300 kV. For 80 kV electron beam, 151 analog-to-digital units (ADUs) correspond to one microscope electron per pixel, while for 300 kV electron beam, 579 ADUs represent one electron per pixel¹⁸. For all the datasets, an exposure time of 1.86 ms (1 ms acquisition time along with 0.86 ms readout time) was employed when acquiring the EMPAD 4D datasets. The estimated dose rate is $\sim 10^5 \text{ e}^- \text{ \AA}^{-2} \text{ s}^{-1}$. The scan size in real space (the number of pixels the beam scan across) can be set from 64×64 to 512×512 . The scan size of the data used in this paper was 256×256 .

DATA AVAILABILITY

The data that supports the findings are available at <https://zenodo.org/communities/hanlab-rice/>.

CODE AVAILABILITY

Codes are available at https://github.com/Chuqiao2333/Hierarchical_Clustering.

Received: 22 November 2021; Accepted: 21 April 2022;

Published online: 18 May 2022

REFERENCES

- Kim, Y. et al. Remote epitaxy through graphene enables two-dimensional material-based layer transfer. *Nature* **544**, 340–343 (2017).
- Bae, S.-H. et al. Graphene-assisted spontaneous relaxation towards dislocation-free heteroepitaxy. *Nat. Nanotechnol.* **15**, 272–276 (2020).
- Yan, R. et al. GaN/NbN epitaxial semiconductor/superconductor heterostructures. *Nature* **555**, 183–189 (2018).
- Duan, X. et al. Lateral epitaxial growth of two-dimensional layered semiconductor heterojunctions. *Nat. Nanotechnol.* **9**, 1024–1030 (2014).
- Li, M.-Y. et al. Epitaxial growth of a monolayer WSe₂–MoS₂ lateral p–n junction with an atomically sharp interface. *Science* **349**, 524–528 (2015).
- Zhang, Z. et al. Robust epitaxial growth of two-dimensional heterostructures, multiheterostructures, and superlattices. *Science* **357**, eaan6814 (2017).
- Han, Y. et al. Sub-nanometre channels embedded in two-dimensional materials. *Nat. Mater.* **17**, 129–133 (2018).
- Xie, S. et al. Coherent, atomically thin transition-metal dichalcogenide superlattices with engineered strain. *Science* **359**, 1131–1136 (2018).
- Rodríguez-González, B., Pastoriza-Santos, I. & Liz-Marzán, L. M. Bending contours in silver nanoprisms. *J. Phys. Chem. B* **110**, 11796–11799 (2006).
- Germain, V., Li, J., Ingert, D., Wang, Z. L. & Pileni, M. P. Stacking faults in formation of silver nanodisks. *J. Phys. Chem. B* **107**, 8717–8720 (2003).
- Rehn, S. M. et al. Mechanical reshaping of inorganic nanostructures with weak nanoscale forces. *Nano Lett.* **21**, 130–135 (2021).
- Jiang, Y. et al. Electron ptychography of 2D materials to deep sub-ångström resolution. *Nature* **559**, 343–349 (2018).
- Zuo, J. M. & Tao, J. *Scanning Electron Nanodiffraction and Diffraction Imaging* 393–427 (Springer New York, 2010).
- Ozdol, V. B. et al. Strain mapping at nanometer resolution using advanced nanobeam electron diffraction. *Appl. Phys. Lett.* **106**, 253107 (2015).
- Baumann, F. H. High precision two-dimensional strain mapping in semiconductor devices using nanobeam electron diffraction in the transmission electron microscope. *Appl. Phys. Lett.* **104**, 262102 (2014).
- Faruqi, A. R. & McMullan, G. Direct imaging detectors for electron microscopy. *Nucl. Instrum. Methods Phys. Res. Sect. Accel. Spectrometers Detect. Assoc. Equip.* **878**, 180–190 (2018).
- Ophus, C. Four-dimensional scanning transmission electron microscopy (4D-STEM): From scanning nanodiffraction to ptychography and beyond. *Microsc. Microanal.* **25**, 563–582 (2019).
- Tate, M. W. et al. High dynamic range pixel array detector for scanning transmission electron microscopy. *Microsc. Microanal.* **22**, 237–249 (2015).
- Han, Y. et al. Strain mapping of two-dimensional heterostructures with sub-picometer precision. *Nano Lett.* **18**, 3746–3751 (2018).
- Yuan, R., Zhang, J. & Zuo, J.-M. Lattice strain mapping using circular Hough transform for electron diffraction disk detection. *Ultramicroscopy* **207**, 112837 (2019).
- Savitzky, B. H. et al. py4DSTEM: A software package for four-dimensional scanning transmission electron microscopy data analysis. *Microsc. Microanal.* **27**, 712–743 (2021).
- Pekin, T. C., Gammer, C., Ciston, J., Minor, A. M. & Ophus, C. Optimizing disk registration algorithms for nanobeam electron diffraction strain mapping. *Ultramicroscopy* **176**, 170–176 (2017).
- Padgett, E. et al. The exit-wave power-spectrum transform for scanning nanobeam electron diffraction: Robust strain mapping at subnanometer resolution and subpicometer precision. *Ultramicroscopy* **214**, 112994 (2020).
- Lee, C.-H. et al. Deep learning enabled strain mapping of single-atom defects in two-dimensional transition metal dichalcogenides with sub-picometer precision. *Nano Lett.* **20**, 3369–3377 (2020).
- Ge, M., Su, F., Zhao, Z. & Su, D. Deep learning analysis on microscopic imaging in materials science. *Mater. Today Nano* **11**, 100087 (2020).
- Li, X. et al. Manifold learning of four-dimensional scanning transmission electron microscopy. *Npj Comput. Mater.* **5**, 5 (2019).
- Zhang, C., Feng, J., DaCosta, L. R. & Voyles, P. M. Atomic resolution convergent beam electron diffraction analysis using convolutional neural networks. *Ultramicroscopy* **210**, 112921 (2019).
- Xu, W. & LeBeau, J. M. A deep convolutional neural network to analyze position averaged convergent beam electron diffraction patterns. *Ultramicroscopy* **188**, 59–69 (2018).
- Jesse, S. et al. Big data analytics for scanning transmission electron microscopy ptychography. *Sci. Rep.* **6**, 26348 (2016).
- Agar, J. C. et al. Machine detection of enhanced electromechanical energy conversion in PbZr_{0.2}Ti_{0.8}O₃ thin films. *Adv. Mater.* **30**, 1800701 (2018).
- Mehta, A. N. et al. Unravelling stacking order in epitaxial bilayer MX₂ using 4D-STEM with unsupervised learning. *Nanotechnology* **31**, 445702 (2020).
- Martineau, B. H., Johnstone, D. N., van Helvoort, A. T. J., Midgley, P. A. & Eggeman, A. S. Unsupervised machine learning applied to scanning precession electron diffraction data. *Adv. Struct. Chem. Imaging* **5**, 3 (2019).
- Hartigan, J. A. & Wong, M. A. A K-means clustering algorithm. *J. R. Stat. Soc. Ser. C. Appl. Stat.* **28**, 100–108 (1979).
- Bholowalia, Purnima, and Arvind Kumar. EBK-means: A clustering technique based on elbow method and k-means in WSN. *Int. J. Comput. Appl.* **105**, 9 (2014).
- McInnes, L., Healy, J., Saul, N. & Großberger, L. UMAP: Uniform manifold approximation and projection. *J. Open Source Softw.* **3**, 861 (2018).
- Meyer, J. C. et al. The structure of suspended graphene sheets. *Nature* **446**, 60–63 (2007).
- Sung, S. H., Schnitzer, N., Brown, L., Park, J. & Hovden, R. Stacking, strain, and twist in 2D materials quantified by 3D electron diffraction. *Phys. Rev. Mater.* **3**, 064003 (2019).

ACKNOWLEDGEMENTS

C.S. and Y.H. are supported by start-up funds provided by Rice University. Y.H. acknowledges the support from the Welch Foundation (C-2065-20210327). M.C. and D.A. M are supported by the NSF MRSEC program (DMR-1719875). S.M.R. would like to acknowledge financial support from a National Science Foundation Graduate Research Fellowship (No. 1842494). This work made use of the Electron Microscopy Center at Rice University. We thank Prof. Ashok Veeraraghavan, Mengnan Zhao, and Guanhuai Gao for helpful discussion. We also thank Prof. Jiwoong Park and Saïen Xie for providing 2D junction samples, which have been previously reported in these references^{8,19}.

AUTHOR CONTRIBUTIONS

Experiments and data analysis were performed by C.S. under the supervision of Y.H.; C.S. and M.C. contributed to the machine learning algorithm under supervision of Y. H.; samples were provided by S.M.R., M.R.J., S.-H.B., J.K.; datasets in Figs. 1–3 were acquired at Cornell by Y.H. under the supervision of D.A.M.

COMPETING INTERESTS

The authors declare no competing interests.

ADDITIONAL INFORMATION

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41524-022-00793-9>.

Correspondence and requests for materials should be addressed to Yimo Han.

Reprints and permission information is available at <http://www.nature.com/reprints>

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022