

Computational Materials Science

Multi Defect Detection and Analysis of Electron Microscopy Images with Deep Learning

--Manuscript Draft--

Manuscript Number:	COMMAT-D-21-00154R1
Article Type:	Full Length Article
Section/Category:	Atomic Description
Keywords:	Deep Learning, Defect Detection and Analysis, Electron Microscopy, Faster R-CNN
Corresponding Author:	Dane Morgan University of Wisconsin-Madison Madison, WI UNITED STATES
First Author:	Mingren Shen
Order of Authors:	Mingren Shen Guanzhao Li Dongxia Wu Yuhan Liu Hima Bharathi Adusumilli Jacob Greaves Wei Hao Nathaniel J. Krakauer Leah Krudy Jacob Perez Varun Sreenivasan Bryan Sanchez Oigimer Torres Wei Li Kevin.G. Field Dane Morgan
Abstract:	Electron microscopy is widely used to explore defects in crystal structures, but human detecting of defects is often time-consuming, error-prone, and unreliable, and is not scalable to large numbers of images or real-time analysis. In this work, we discuss the application of machine learning approaches to find the location and geometry of different defect clusters in irradiated steels. We show that a deep learning based Faster R-CNN analysis system has a performance comparable to human analysis with relatively small training data sets. This study proves the promising ability to apply deep learning to assist the development of automated analysis microscopy data even when multiple features are present and paves the way for fast, scalable, and reliable analysis systems for massive amounts of modern electron microscopy data.
Response to Reviewers:	Prof. Susan Sinnott, Editor Computational Materials Science Dear Prof. Susan Sinnott, Thank you for your positive response and for allowing us to revise our manuscript titled "Multi Defect Detection and Analysis of Electron Microscopy Images with Deep Learning" (manuscript ID COMMAT-D-21-00154). We have carefully examined the reviewer's comments and revised the manuscript accordingly as explained in the reply letter below. We hope that our response is sufficient to address all comments.

Listed below is a summary of changes made to the manuscript in response to the comments by the reviewer, as well as some additional edits made for clarification. The text changes are shown as highlighted text in the revised manuscript.

- We have added more information to explain the reason why we need new tools for defect analysis on Page 4 Line 50-55 of the revised manuscript.
- We have added the scale bar into Figure 1.
- We have corrected the link to the Figure 3.
- We have added discussion of error analysis of diameter difference and its relationship with defect density on Page 19-20 Line 290-304 of the revised manuscript.
- We have modified the text of model reproduction to clarify and include the aspects mentioned by the reviewer. The new text on the Page 20-21 line 311-323 of the revised manuscript.
- We have added one sentence about the default parameter used for CLAHE and Gaussian on Page 25 Line 406-408 of the revised manuscript.
- We have clarified the metric used for grid search of best parameter on Page 28 Line 464-465 of the revised manuscript.
- We have modified text to clarify what was done for watershed and ellipse fitting on Page 29 Line 481-488 of the revised manuscript.
- We have added typo of missing word "twice" for the black dot diameter definition on Page 30 Line 492 of the revised manuscript.
- We have proofread the text and equations and done our best to correct all the typos in the manuscript.

Best regards,

Mingren Shen, Guanzhao Li, Dongxia Wu, Yuhua Liu, Hima Bharathi Adusumilli, Jacob Greaves, Wei Hao, Nathaniel J. Krakauer, Leah Krudy, Jacob Perez, Varun Sreenivasan, Bryan Sanchez, Oigimer Torres, Wei Li, Kevin.G. Field, Dane Morgan

Corresponding author: Dane Morgan, E-mail: ddmorgan@wisc.edu

We thank the reviewers for their comments. We incorporated all the points raised by the reviewers in the revision, which are shown as highlighted text in the resubmitted manuscript document. Below we provide point-by-point responses (in black) to all questions of the reviewers (in blue) and descriptions of any changes to the manuscript. The page numbers refer to the revised manuscript.

Reviewer #2

Multi Defect Detection and Analysis of Electron Microscopy Images with Deep Learning

The authors present results to illustrate the quality of defect identification and property measurement in TEM micrograph images using a machine learning approach.

The work itself is well explained, with good choices made for comparison metrics, and the authors should be commended for making their data and models available.

I think the paper could benefit from tightening up the definitions of the metrics so that a quality comparison can be made between this method and others.

In some places, the enthusiasm for supervised learning shines through without some of the more cautious scientific caveats on reproducibility.

I have detailed below some places the authors might like to revisit.

We thank the reviewer for their supportive comments and detailed suggestions. Based on the comments below, we have revised the text related to metrics for clarity and tried to be more balanced in our discussion of the promise of supervised learning. Details are given in the point-by-point responses below.

Q1. Page 2

The authors state:

"many accurate and standardized data analysis tools [can] be effectively deployed, [none] can replace standard human analysis for identifying and locating defects"

This seems to be self-contradictory. Automated image analyses are routinely used. What are the cases that the authors have in mind where the current methods fail?

We agree with the reviewer that the text here is not clear and what we want to express here is that due to sample and imaging complexity, common automated tools may fail in edge cases or new material systems, so human efforts and hand-tunings are still needed which makes the data volume problem still hard to solve.

We have changed our description of Page 4 Line 50-55 into the following:

Not surprisingly, the EM community has developed many accurate image data analysis tools that can be effectively deployed to accommodate the large volume of EM data^{3,4}. However, due to the complexity of images of material systems, these tools generally still need significant hand-tuning, and in some cases (like counting defects of irradiated materials), human identification of each defect is still the norm.

Q2. Page 6 image caption:

The scale in pixels might be useful here too given the context.

We have added the scale bar into Figure 1.

Q3. Page 9, description of the assessments

These assessments seem a very good set to compare the machine-labeling to human assessment, but as humans are fallible the set by itself does not give an objective quality measure.

Have the authors performed any tests against simulated data, where the ground truth is exactly and unambiguously known (for some geometrical quantities at least)?

Such tests would lend enormous weight to the power of the method, as well as making it easier for other groups to compare different methods in the future.

We agree with the reviewer that human labeling is not as objective as would be desired. We have not performed tests against simulated data for this type of loop defect data as no such simulated data exists, nor is there any easy path to obtain it. Simplified simulated data could be created (e.g., simple circles on a uniform background) but it would be impossible to assess how relevant any model results would be for real data. We are exploring methods to generate simulated data (e.g., using multi-slice methods and generative adversarial networks) but these are significant additional projects and there is always uncertainty about the impact of the necessary approximations. Therefore, we believe simulated data studies are best left to a later work where all the issues can be addressed and are beyond the scope of the present effort.

Q4. Page 9 para 2 and a few other places

There are some broken links "As shown in .., the red circles..."

We have corrected the link to the Figure, and it should be "As shown in Figure 3".

Q5. Page 10, the cutoff for IoU

"As a compromise [] we used cutoff IoU = 0.4"

This value "feels" good for image 2 where the defects are densely packed, but poor for image 3 where the defects are well separated.

Q5a)

Is there a simple scaling for "good" performance in IoU based on e.g. homogeneous distribution density?

While there may be some theoretical justification for IoU values that can be derived from some kind of assumption of homogeneous density it is simpler and more common

in the community to base the result on actual performance vs. IoU, as in Figure 4 of the manuscript. A large IoU is ideal as it assures the best match between true and predicted bounding boxes. However, at some level, this requirement becomes unreasonable, and so we try to take the largest value that does not hurt performance significantly. We find this value as 0.4, which is similar to the standard 0.5 used in the computer vision community (For example, both ImageNet , COCO , and PASCAL , the three largest and widely used image datasets, evaluate object detection performance at a cut-off IoU = 0.5). Also, we note that the cut-off IoU 0.4 corresponds to a 57% overlapping area for two unit-size bounding boxes ($x/(2-x) = 0.4 \Rightarrow x = 4/7 \sim 0.57$ as shown in the figure below), which is quite significant alignment.

Q5b)

Can the authors show examples of defect mismatches with $\text{IoU} < 0.4$, > 0.4 ? This would help readers less familiar with using this metric.

We are not totally sure what the reviewer means here but below we show three pairs of machine learning predicted (red) and ground truth (green) bounding box pairs with IOUs of (a) 0.221 (well below 0.4), (b) 0.402 (near 0.4), and (c) 0.624 (well above 0.4) to illustrate how they look. With a cutoff of 0.4 we would identify (a) as a false positive but (b) and (c) as successfully finding the defect. These results seem intuitively reasonable.

(a) (b) (c)

We are not sure how much these results would inform the general reader, so we have not added them to the manuscript.

Q6. Page 12, description of assessment 3.

Different human observers will categorize differently. Is there any research to indicate a confusion matrix comparable to table 3 for experienced human researchers? As it stands, I have no way of telling if the accuracies given (76%,87%,94%) are excellent, adequate, or hopeless.

The accuracy of human labeling is quite subtle to determine. It can depend sensitively on the system, the specific people, how much they communicate, how much time they spend, how many people check each result, etc. We have therefore not made any attempt to assess this value and view it as an important problem but outside the scope of this paper. Furthermore, we think the question of assessing the model accuracies should not hinge on the relative accuracies of different people. We think it is more important to assess the model in terms of whether the accuracies are adequate to yield physically useful information assuming the ground truth labeling is correct. We are clearly achieving this level of accuracy given the excellent mean size and areal densities we obtain. This paper is focused on assessing the accuracy of the model for reproducing our ground truth labeling. The challenge of obtaining accurate labeling in the community for training such models is an important topic but outside our scope.

Q7. Page 14, description of diameters

How is diameter defined? Particularly for the black-spot damage we could argue in favor of full-width-half maximum, or 95% rule, or...

The "mean" also needs defining - for an ellipse it could be the arithmetic or geometric average diameter, or those values averaged over the sample.

Perhaps the reviewer did not see the definition, but the diameter is defined on Page 29-30 Line 490-493 of the METHODS section as the following text shown

The diameters and areas of defects are defined as follows, where a and b are half the lengths of major and minor axes of the ellipse. The diameter of the $a/2<111>$ and

$a<100>$ defects are defined as 2a. The diameter of the black dot is defined as twice the square root of (ab) .

We agree with the reviewers there might be different ways to define defects, but we think our approach is a reasonable one since all defects are 3D objects and the STEM images are just 2D projection of these 3D objects, so for more elliptical shape objects like $a/2<111>$ and $a<100>$ defects, 2a is a good characterization of its size since b is just a shortened by the projection of the loop being on an inclined plane and for more spherical objects like black dots, twice the square root of (ab) is a good approximation for its real size. And we thank the author since we found we missed the word twice (shown in red in original text) for black dots. This is corrected now, and all calculations are correct since we always comparing the diameter of the defect.

The mean size means the arithmetic mean of defect sizes, and we have added "arithmetic" to all the mean defect diameter discussion occurrences all over the manuscript.

Q8. Page 14, comparison of diameter performance

I don't think the case is really proved here.

Consider a ring with an intensity profile peaking at radius R, but with a line profile through the peak itself having a finite width r.

If $r \ll R$, we should be able to make a really accurate estimate of the radius R.

If $r \sim R$, the task is much harder. In a micrograph image we would be lucky to get a perfect circle with constant line profile, further complicating the issue.

Incautious fitting might easily lead to a case where the shape of the line profile inside and outside the maximum or noise leads to a systematic over or underestimate.

Q8a) How well does eg a Hough transform perform in the cases shown?

It is quite possible that a hand-tuned approach, e.g., a Hough transform could provide better performance on this particular data set and we do not mean to claim otherwise. However, such approaches are unable to tackle general defect detection in TEM images and would require impractically extensive hand tuning of hyperparameters for every dataset – in fact, this exact use case was the motivation of previous work in Ref. 5 of the revised manuscript. This fact cannot really be proven in any mathematical way, but we think that the ubiquitous use of hand labeling and the absence of any widely adopted automated defect detection tool in the radiation damage community is clear evidence that it is the case. The goal of this paper is to help develop a general tool for multiple class defect analysis that does not require extensive hand tuning for every new problem.

Q8b) How big is 2-9 pixels compared to the line profile?

This number of pixels is comparable (and in some cases less) to the thickness of defect loops. The ambiguity associated with the boundary region of the defects may create some of the observed noise.

We agree with the reviewer that finding a proper profile of the defects might help the detection of loops, but the case is more complicated than the ideal settings. Below we show one line profile of pixel intensity of an $a<100>$ loop and this line is considered as the diameter of the loop and for this defect 2-9 is much smaller than the defect diameter.

Q9. Page 16: "2 pixels [error] is considered negligible in terms of impact on material properties."

This is not a valid statement.

We agree that this statement needed some qualification as it will depend on the defect sizes. To further clarify that the error is negligible for our defect population we have done a sensitivity analysis based on previous studies of hardening from loops. As discussed in Field et al. analysis shows for one of our alloys (Fe–Cr–Al alloys) the

hardening under irradiation from loops is

Now consider an error in diameter d defined as e . The fractional error in $\Delta\sigma_y$ due to the error e is $\Delta\sigma_y d + \epsilon - \Delta\sigma_y d \Delta\sigma_y d \approx \epsilon 2d$, where the approximate equality holds for $\epsilon \ll d$. For $e = 1.7$ nm (which is 2 pixels for our largest pixel sizes, see below) and $d = 21.4$ nm (our average sizes of $a/2<111>$ and $a<100>$ defects), we get the fractional error in $\Delta\sigma_y$ as 1.7 nm / $(2 * 21.4$ nm) ≈ 0.04 , which is well within the uncertainty of such microstructure-based analysis. However, for smaller defects this percentage error could clearly become larger. The error of diameter between ML results and human results appears to be approximately symmetrically distributed in positive and negative directions and independent of defect density, as shown in detail in SI section 4 and 5. However, for smaller defects this percentage error could clearly become larger. To address this issue we have added the following text to the manuscript on Page 19-20 Line 290-304:

To further clarify that the error is negligible for our defect population we have done a sensitivity analysis based on previous studies of hardening from loops. As discussed in Field et al. 31 simple dispersed barrier hardening models suggest that the hardening under irradiation from loops is of the form $\Delta\sigma_y = A d$ where A is a constant and d is the diameter of the defect. Now consider an error in diameter d defined as . The fractional error in $\Delta\sigma_y$ due to the error is $\Delta\sigma_y d + \epsilon - \Delta\sigma_y d \Delta\sigma_y d \approx \epsilon 2d$, where the approximate equality holds for $\epsilon \ll d$. For $\epsilon = 1.7$ nm (which is 2 pixels for our largest pixel sizes, see below) and $d = 21.4$ nm (our average sizes of $a/2<111>$ and $a<100>$ defects), we get the fractional error in $\Delta\sigma_y$ as 1.7 nm / $(2 * 21.4$ nm) ≈ 0.04 , which is well within the uncertainty of such microstructure-based analysis. However, for smaller defects this percentage error could clearly become larger. The error of diameter between ML results and human results appears to be approximately symmetrically distributed in positive and negative directions and independent of defect density, as shown in detail in SI section 4 and 5.

Q9a How big is a pixel?

The pixel to nanometer ratio is from 0.14nm/pixel to 0.87nm/pixel as already discussed on Page 17 Line 258 of the manuscript.

Q9b Is the error random or systematic?

The error is nearly random in the sense that it is not always high or low, nor does it show some unexpected groupings. As an illustration, below we show the diameter difference between ML and prediction for each class and shows that forms relatively symmetric and almost Gaussian looking distributions. Here we only considered those paired ground-truth labeling and ML prediction that are both correct in position prediction and defect type classification.

We have put the error analysis into SI section 4.

Q9c Is the error dependent on defect density?

We have examined the arithmetic mean of the error (ME), arithmetic mean of absolute value of the error (MAE), fractional ME (ME/ (arithmetic mean diameter of ground truth labeling), or FME) and fractional MAE (MAE/ (arithmetic mean diameter of ground truth labeling), or FMAE) in diameter for each of our 12 test images and plotted the values vs. the ground truth defect densities. The errors vs. density plots are shown below.

Figure 1. arithmetic mean of the error changes with the defect density

Figure 2. arithmetic mean of absolute value of the error changes with the defect density

Figure 3. (arithmetic mean of the error) / (arithmetic mean diameter of ground truth labeling) changes with the defect density

Figure 4. (arithmetic mean of absolute value of the error) / (arithmetic mean diameter of ground truth labeling) changes with the defect density

Figures 1 and 2 may show some trend of decreasing value with higher density, but we believe this is almost entirely due to the fact the low densities are correlated with larger loops, which lead to larger absolute errors. To correct for this aspect, we measure fractional (or percentage) errors, shown in Figure 3 and 4. These show no discernable trend with density. We therefore believe that there is no trend of error with density when the confounding effect of defect size is removed.

We have put the error dependence of density discussion into SI Section 5.

We have added the following text on Page 20 Line 301-304 of the revised manuscript.

The errors of diameter between ML results and human results appear to be approximately symmetrically distributed in positive and negative directions and independent of defect density, as shown in detail in SI section 4 and 5.

Q10. Page 16: "Once the model is properly trained, it will yield a unique and reproducible labelling for every image"

This is a difficult statement for me. It is trivially true that the analysis run a second time will give the same result.

But if the training set is added to, or is prepared by a different researcher, then it is not necessarily true.

The labelling is therefore not an algorithmic property of the images, and we should not expect other groups to be able to reproduce it starting from the images alone.

As the authors themselves point out "even the same person may label the same defects differently even after a short break."

We agree with the reviewers' points here and that the meaning of this statement is not clear. We have modified the text somewhat to clarify and include the aspects mentioned by the reviewer. The new text on the Page 20-21 Line 311-323 reads as the following,

While the exact performance of the present automated approach compared to different human researchers is difficult to determine rigorously there is no doubt that the present approach is much more consistent. Previous studies have shown that different labelers tend to label defects in different ways and even the same person may label the same defects differently even after a short break^{2,5}. Such issues can make any given data analysis somewhat unreliable and make it difficult to integrate results across different teams and or time periods in larger analysis efforts. However, once a machine learning model is properly trained, it will yield a unique and reproducible labeling for every image. If the community could converge on a single or small number of models this could greatly increase the reproducibility in labeling of STEM experiments. That said, models trained on different data and/or different human labeling could give different predictions, so establishing community accepted models is an important part of using these approaches to obtain more consistent results.

Q11. Page 19: "The automated analysis [takes] about 0.1 s/image"

How does this value scale with pixel count or defect density?

The Faster R-CNN model will resize each image into 1024 pixel x 1024 pixel, which is also the pixel size of our TEM images, so the time to run object detection is not influenced by a higher number of pixels. However, downsizing an image with a higher

number of pixels will require extra time. We don't know the exact timing for this operation, and it could depend on many factors (e.g., machine speed, algorithms used, exact pixel count, etc.), but it is likely just tens of seconds at most for TEM images on a single processor. We, therefore, do not think it will be a significant issue. Assuming we use the same Faster RCNN hyperparameters the algorithm runs at approximately the same speed regardless of the number of defects in the image. If we have to alter the algorithm to use more proposal regions to detect more defects this could slow it down, but it is difficult to speculate on what would be necessary without detailed studies of the specific images.

Q12. Page 20: Data set preparation

Q12a. Any necessary parameters for CLAHE and Gaussian blur should be given explicitly.

The parameters used for CLAHE and Gaussian blur are all from the default parameter setting of scikit-images and the link to the detailed documentation of parameters used is put in the footnote.

We have added the text with two references of these parameters used on Page 25 Line 406-408 of the revised manuscript.

The parameters used for CLAHE45 and Gaussian blur46 are all from the default parameter setting of scikit-images and details can be found in the references given here for these methods.

Q12b. The use of additional filtered images in the RGB channels is a very neat idea. By how much does it actually help?

We do not have a direct assessment of the impact of adding different filtered images in the RGB channels. It is difficult to assess the model with just one channel as this would involve rewriting the basic Faster RCNN tools which take RGB as their standard input. We could provide the same images to each RGB channel, but this would require retraining the entire model, is a significant effort, and it is not clear how the separate channels would couple. Overall, we think it would be challenging to develop a rigorous and robust quantitative test to determine the impact of this particular choice of the training, and such testing is somewhat outside the main focus of the paper. We are therefore not adding any quantitative assessment of this approach relative to others and we would ask for the reviewer's understanding that there is limited time for such assessment. That said, we think the choice we have made of using different filtered images in the RGB channels is likely beneficial for the working of Faster R-CNN as it provides an additional way for data augmentation to make the model more robust to noise.

Q13. Page 24: "We used grid search [to find] the best choice of these two values."

How is "best" defined? Is this in terms of table 1,2,3 or a combination?

We selected the best values based on maximizing the F1 scores of the testing images, which is the F1 score shown in row 3 in Table 2. We have clarified this in the text on Page 28 Line 464-465.

Q14. Page 25: "Fitting the contour with an ellipse function"

How is the fitting done?

The overall fitting of the defect geometry is done in two steps. First, we use the watershed algorithm provided by OpenCV. The algorithm is applied to the cropped enlarged region of bounding boxes that contains the defect. We followed the official tutorial from OpenCV, and more details can be found there. Then the boundary from the Watershed algorithm was fit to an ellipse. Fitting with an ellipse was useful to match the approach used by the community, obtain a well-defined shape with simple geometric descriptors, and smooth out the otherwise rather rough boundaries found by the Watershed algorithm. OpenCV's fitEllipse() function was called to perform the fit.

We have modified the mentioned text in the paper on Page 29 Line 481-488 to clarify what done for watershed and ellipse fitting,

Watershed methods were applied to find the boundary between defect pixels and background pixels. We followed the official tutorial from OpenCV for performing the watershed and details of the approach can be found there⁵⁰. We then fit the boundaries found from the Watershed algorithm to an ellipse. This fitting was done to match the approach used by the radiation defect analysis community, obtain a well-defined shape with simple geometric descriptors, and smooth out the otherwise rather rough boundaries found by the Watershed algorithm. The fitting was done with OpenCV's fitEllipse() function⁵¹.



Dear Editor,

We would appreciate your consideration of the enclosed manuscript:

Multi Defect Detection and Analysis of Electron Microscopy Images with Deep Learning

By Mingren Shen, Guanzhao Li, Dongxia Wu, Yuhan Liu, Hima Bharathi Adusumilli, Jacob Greaves, Wei Hao, Nathaniel J. Krakauer, Leah Krudy, Jacob Perez, Varun Sreenivasan, Bryan Sanchez, Oigimer Torres, Wei Li, Kevin.G. Field, Dane Morgan

for publication in *Computational Materials Science*.

Electron microscopy is one of the most widely used methods to explore and analyze defects in crystal structures, but when extensive defect identification and analysis is required the existing human-based analysis methods are time-consuming, error-prone, unreliable, and not scalable. It is therefore necessary to automate such analysis to make full use of the data available from modern electron microscopes.

In this work, we demonstrate a practical deep learning based model for automatic STEM image defect detection in irradiated nuclear steels. The model includes the Faster R-CNN module for detection and the watershed flood module for geometry fitting. The model improves significantly over similar previous similar efforts as it is simpler to train and is extended to multiple defect types. Our model performance is quite encouraging. Specifically, on three different defect types we achieved a F1 score of 0.78, and the predicted sizes and areal densities were very similar to the results obtained from previous human researchers.

Overall, we provide an accurate, efficient, reproducible, scalable, and extensible model to support counting of defects in STEM images of irradiated materials, and potentially many other systems. The model can be used to detect multiple types of defects (and potentially other features of the material) simultaneously. The model provides an important step in developing more automated STEM analysis. We believe this work will be of interest to the broad readership of Computational Materials Science and we sincerely hope you consider our paper for publication.

Best regards,

Mingren Shen, Guanzhao Li, Dongxia Wu, Yuhan Liu, Hima Bharathi Adusumilli, Jacob Greaves, Wei Hao, Nathaniel J. Krakauer, Leah Krudy, Jacob Perez, Varun Sreenivasan, Bryan Sanchez, Oigimer Torres, Wei Li, Kevin.G. Field, Dane Morgan

Corresponding author: Dane Morgan, E-mail: ddmorgan@wisc.edu



Prof. Susan Sinnott, Editor
Computational Materials Science

Dear Prof. Susan Sinnott,

Thank you for your positive response and for allowing us to revise our manuscript titled “Multi Defect Detection and Analysis of Electron Microscopy Images with Deep Learning” (manuscript ID COMMAT-D-21-00154). We have carefully examined the reviewer’s comments and revised the manuscript accordingly as explained in the reply letter below. We hope that our response is sufficient to address all comments.

Listed below is a summary of changes made to the manuscript in response to the comments by the reviewer, as well as some additional edits made for clarification. The text changes are shown as highlighted text in the revised manuscript.

- *We have added more information to explain the reason why we need new tools for defect analysis on Page 4 Line 50-55 of the revised manuscript.*
- *We have added the scale bar into Figure 1.*
- *We have corrected the link to the Figure 3.*
- *We have added discussion of error analysis of diameter difference and its relationship with defect density on Page 19-20 Line 290-304 of the revised manuscript.*
- *We have modified the text of model reproduction to clarify and include the aspects mentioned by the reviewer. The new text on the Page 20-21 line 311-323 of the revised manuscript.*
- *We have added one sentence about the default parameter used for CLAHE and Gaussian on Page 25 Line 406-408 of the revised manuscript.*
- *We have clarified the metric used for grid search of best parameter on Page 28 Line 464-465 of the revised manuscript.*
- *We have modified text to clarify what was done for watershed and ellipse fitting on Page 29 Line 481-488 of the revised manuscript.*
- *We have added typo of missing word “twice” for the black dot diameter definition on Page 30 Line 492 of the revised manuscript.*
- *We have proofread the text and equations and done our best to correct all the typos in the manuscript.*

Best regards,

Mingren Shen, Guanzhao Li, Dongxia Wu, Yuhan Liu, Hima Bharathi Adusumilli, Jacob Greaves, Wei Hao, Nathaniel J. Krakauer, Leah Krudy, Jacob Perez, Varun Sreenivasan, Bryan Sanchez, Oigimer Torres, Wei Li, Kevin.G. Field, Dane Morgan

Corresponding author: Dane Morgan, E-mail: ddmorgan@wisc.edu

We thank the reviewers for their comments. We incorporated all the points raised by the reviewers in the revision, which are shown as highlighted text in the resubmitted manuscript document. Below we provide point-by-point responses (in black) to all questions of the reviewers (in blue) and descriptions of any changes to the manuscript. The page numbers refer to the revised manuscript.

Reviewer #2

Multi Defect Detection and Analysis of Electron Microscopy Images with Deep Learning

The authors present results to illustrate the quality of defect identification and property measurement in TEM micrograph images using a machine learning approach.

The work itself is well explained, with good choices made for comparison metrics, and the authors should be commended for making their data and models available.

I think the paper could benefit from tightening up the definitions of the metrics so that a quality comparison can be made between this method and others.

In some places, the enthusiasm for supervised learning shines through without some of the more cautious scientific caveats on reproducibility.

I have detailed below some places the authors might like to revisit.

We thank the reviewer for their supportive comments and detailed suggestions. Based on the comments below, we have revised the text related to metrics for clarity and tried to be more balanced in our discussion of the promise of supervised learning. Details are given in the point-by-point responses below.

Q1. Page 2

The authors state:

"many accurate and standardized data analysis tools [can] be effectively deployed, [none] can replace standard human analysis for identifying and locating defects"

This seems to be self-contradictory. Automated image analyses are routinely used. What are the cases that the authors have in mind where the current methods fail?

We agree with the reviewer that the text here is not clear and what we want to express here is that due to sample and imaging complexity, common automated tools may fail in edge cases or new material systems, so human efforts and hand-tunings are still needed which makes the data volume problem still hard to solve.

We have changed our description of Page 4 Line 50-55 into the following:

Not surprisingly, the EM community has developed many accurate image data analysis tools that can be effectively deployed to accommodate the large volume of EM data^{3,4}. However, due to the complexity of images of material systems, these tools generally still

need significant hand-tuning, and in some cases (like counting defects of irradiated materials), human identification of each defect is still the norm.

Q2. Page 6 image caption:

The scale in pixels might be useful here too given the context.

We have added the scale bar into Figure 1.

Q3. Page 9, description of the assessments

These assessments seem a very good set to compare the machine-labeling to human assessment, but as humans are fallible the set by itself does not give an objective quality measure.

Have the authors performed any tests against simulated data, where the ground truth is exactly and unambiguously known (for some geometrical quantities at least)?

Such tests would lend enormous weight to the power of the method, as well as making it easier for other groups to compare different methods in the future.

We agree with the reviewer that human labeling is not as objective as would be desired. We have not performed tests against simulated data for this type of loop defect data as no such simulated data exists, nor is there any easy path to obtain it. Simplified simulated data could be created (e.g., simple circles on a uniform background) but it would be impossible to assess how relevant any model results would be for real data. We are exploring methods to generate simulated data (e.g., using multi-slice methods and generative adversarial networks) but these are significant additional projects and there is always uncertainty about the impact of the necessary approximations. Therefore, we believe simulated data studies are best left to a later work where all the issues can be addressed and are beyond the scope of the present effort.

Q4. Page 9 para 2 and a few other places

There are some broken links "As shown in .., the red circles..."

We have corrected the link to the Figure, and it should be "As shown in Figure 3".

Q5. Page 10, the cutoff for IoU

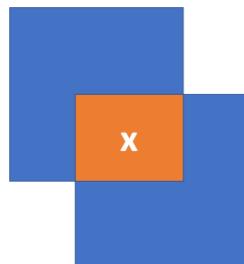
"As a compromise [] we used cutoff IoU = 0.4"

This value "feels" good for image 2 where the defects are densely packed, but poor for image 3 where the defects are well separated.

Q5a)

Is there a simple scaling for "good" performance in IoU based on e.g. homogeneous distribution density?

While there may be some theoretical justification for IoU values that can be derived from some kind of assumption of homogeneous density it is simpler and more common in the community to base the result on actual performance vs. IoU, as in Figure 4 of the manuscript. A large IoU is ideal as it assures the best match between true and predicted bounding boxes. However, at some level, this requirement becomes unreasonable, and so we try to take the largest value that does not hurt performance significantly. We find this value as 0.4, which is similar to the standard 0.5 used in the computer vision community (For example, both ImageNet¹, COCO², and PASCAL³, the three largest and widely used image datasets, evaluate object detection performance at a cut-off IoU = 0.5). Also, we note that the cut-off IoU 0.4 corresponds to a 57% overlapping area for two unit-size bounding boxes ($x / (2-x) = 0.4 \Rightarrow x = 4/7 \sim 0.57$ as shown in the figure below), which is quite significant alignment.



Q5b)

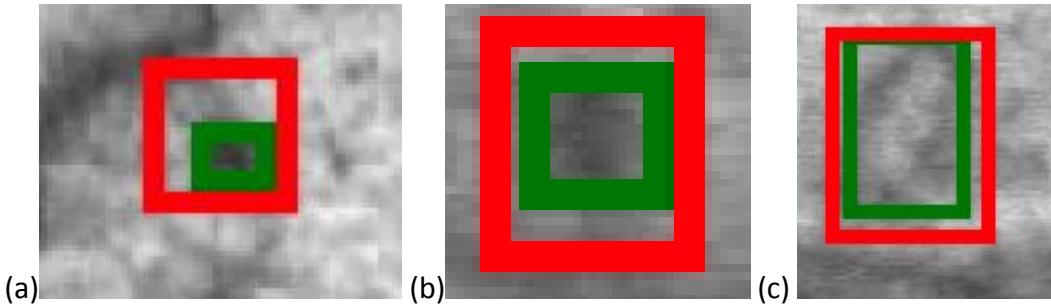
Can the authors show examples of defect mismatches with $\text{IoU} < 0.4$, > 0.4 ? This would help readers less familiar with using this metric.

We are not totally sure what the reviewer means here but below we show three pairs of machine learning predicted (red) and ground truth (green) bounding box pairs with IOUs of (a) 0.221 (well below 0.4), (b) 0.402 (near 0.4), and (c) 0.624 (well above 0.4) to illustrate how they look. With a cutoff of 0.4 we would identify (a) as a false positive but (b) and (c) as successfully finding the defect. These results seem intuitively reasonable.

¹ Russakovsky, Olga, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang et al. "Imagenet large scale visual recognition challenge." *International journal of computer vision* 115, no. 3 (2015): 211-252.

² Lin, Tsung-Yi, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. "Microsoft coco: Common objects in context." In *European conference on computer vision*, pp. 740-755. Springer, Cham, 2014.

³ Everingham, Mark, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. "The pascal visual object classes challenge: A retrospective." *International journal of computer vision* 111, no. 1 (2015): 98-136.



We are not sure how much these results would inform the general reader, so we have not added them to the manuscript.

Q6. Page 12, description of assessment 3.

Different human observers will categorize differently. Is there any research to indicate a confusion matrix comparable to table 3 for experienced human researchers? As it stands, I have no way of telling if the accuracies given (76%, 87%, 94%) are excellent, adequate, or hopeless.

The accuracy of human labeling is quite subtle to determine. It can depend sensitively on the system, the specific people, how much they communicate, how much time they spend, how many people check each result, etc. We have therefore not made any attempt to assess this value and view it as an important problem but outside the scope of this paper. Furthermore, we think the question of assessing the model accuracies should not hinge on the relative accuracies of different people. We think it is more important to assess the model in terms of whether the accuracies are adequate to yield physically useful information assuming the ground truth labeling is correct. We are clearly achieving this level of accuracy given the excellent mean size and areal densities we obtain. This paper is focused on assessing the accuracy of the model for reproducing our ground truth labeling. The challenge of obtaining accurate labeling in the community for training such models is an important topic but outside our scope.

Q7. Page 14, description of diameters

How is diameter defined? Particularly for the black-spot damage we could argue in favor of full-width-half maximum, or 95% rule, or...

The "mean" also needs defining - for an ellipse it could be the arithmetic or geometric average diameter, or those values averaged over the sample.

Perhaps the reviewer did not see the definition, but the diameter is defined on Page 29-30 Line 490-493 of the **METHODS section** as the following text shown

The diameters and areas of defects are defined as follows, where a and b are half the lengths of major and minor axes of the ellipse. The diameter of the $a/2<111>$ and $a<100>$

*defects are defined as $2a$. The diameter of the black dot is defined as **twice** the square root of (ab) .*

We agree with the reviewers there might be different ways to define defects, but we think our approach is a reasonable one since all defects are 3D objects and the STEM images are just 2D projection of these 3D objects, so for more elliptical shape objects like $a/2<111>$ and $a<100>$ defects, $2a$ is a good characterization of its size since b is just a shortened by the projection of the loop being on an inclined plane⁴ and for more spherical objects like *black dots*, twice the square root of (ab) is a good approximation for its real size. And we thank the author since we found we missed the word twice (shown in red in original text) for black dots. This is corrected now, and all calculations are correct since we always comparing the diameter of the defect.

The mean size means the arithmetic mean of defect sizes, and we have added “arithmetic” to all the mean defect diameter discussion occurrences all over the manuscript.

Q8. Page 14, comparison of diameter performance

I don't think the case is really proved here.

Consider a ring with an intensity profile peaking at radius R , but with a line profile through the peak itself having a finite width r .

If $r \ll R$, we should be able to make a really accurate estimate of the radius R .

If $r \sim R$, the task is much harder. In a micrograph image we would be lucky to get a perfect circle with constant line profile, further complicating the issue.

Incautious fitting might easily lead to a case where the shape of the line profile inside and outside the maximum or noise leads to a systematic over or underestimate.

Q8a) How well does eg a Hough transform perform in the cases shown?

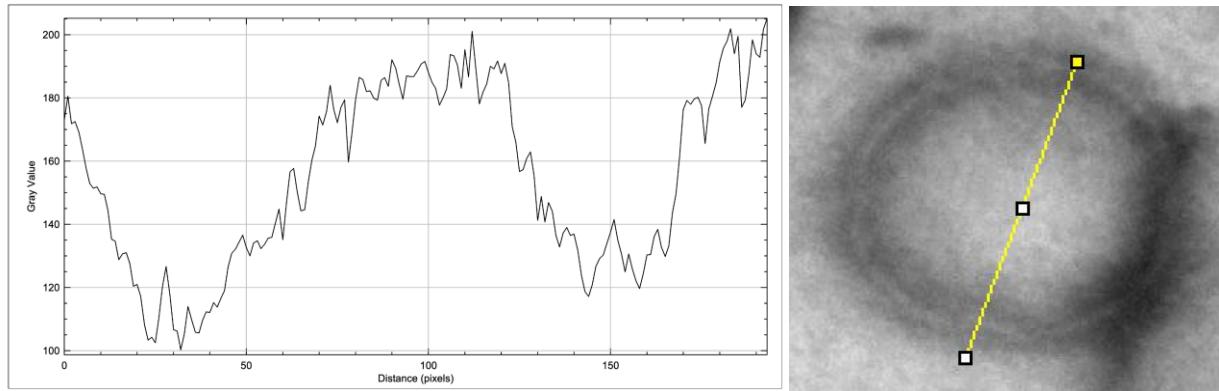
It is quite possible that a hand-tuned approach, e.g., a Hough transform could provide better performance on this particular data set and we do not mean to claim otherwise. However, such approaches are unable to tackle general defect detection in TEM images and would require impractically extensive hand tuning of hyperparameters for every dataset – in fact, this exact use case was the motivation of previous work in Ref. 5 of the revised manuscript. This fact cannot really be proven in any mathematical way, but we think that the ubiquitous use of hand labeling and the absence of any widely adopted automated defect detection tool in the radiation damage community is clear evidence that it is the case. The goal of this paper is to help develop a general tool for multiple class defect analysis that does not require extensive hand tuning for every new problem.

Q8b) How big is 2-9 pixels compared to the line profile?

⁴ Yao, Bo, Danny J. Edwards, and Richard J. Kurtz. "TEM characterization of dislocation loops in irradiated bcc Fe-based steels." Journal of Nuclear Materials 434, no. 1-3 (2013): 402-410.

This number of pixels is comparable (and in some cases less) to the thickness of defect loops. The ambiguity associated with the boundary region of the defects may create some of the observed noise.

We agree with the reviewer that finding a proper profile of the defects might help the detection of loops, but the case is more complicated than the ideal settings. Below we show one line profile of pixel intensity of an $a<100>$ loop and this line is considered as the diameter of the loop and for this defect 2-9 is much smaller than the defect diameter.



Q9. Page 16: "2 pixels [error] is considered negligible in terms of impact on material properties."

This is not a valid statement.

We agree that this statement needed some qualification as it will depend on the defect sizes. To further clarify that the error is negligible for our defect population we have done a sensitivity analysis based on previous studies of hardening from loops. As discussed in Field et al. analysis shows for one of our alloys (Fe–Cr–Al alloys⁵) the hardening under irradiation from loops is

$$\Delta\sigma_y = M\alpha\mu b\sqrt{\rho d} = A\sqrt{d}$$

Now consider an error in diameter d defined as e . The fractional error in $\Delta\sigma_y$ due to the error e is $(\Delta\sigma_y(d + \epsilon) - \Delta\sigma_y(d))/\Delta\sigma_y(d) \approx \epsilon/(2d)$, where the approximate equality holds for $\epsilon \ll d$. For $e = 1.7$ nm (which is 2 pixels for our largest pixel sizes, see below) and $d = 21.4$ nm (our average sizes of $a/2<111>$ and $a<100>$ defects), we get the fractional error in $\Delta\sigma_y$ as 1.7 nm / $(2 * 21.4$ nm) ≈ 0.04 , which is well within the uncertainty of such microstructure-based analysis. However, for smaller defects this percentage error could clearly become larger. The error of diameter between ML results and human results appears to be approximately symmetrically distributed in positive and negative directions and independent of defect density, as shown in detail in SI section 4 and 5. However, for smaller defects this percentage error could clearly

⁵ Field, Kevin G., Xunxiang Hu, Kenneth C. Littrell, Yukinori Yamamoto, and Lance L. Snead. "Radiation tolerance of neutron-irradiated model Fe–Cr–Al alloys." Journal of Nuclear Materials 465 (2015): 746-755.

become larger. To address this issue we have added the following text to the manuscript on Page 19-20 Line 290-304:

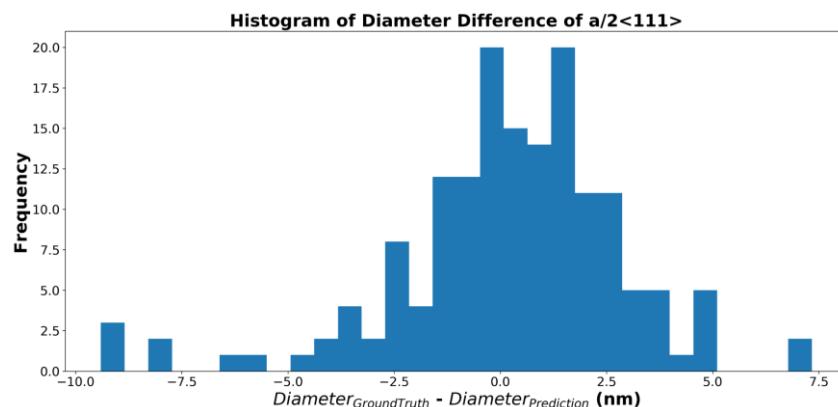
To further clarify that the error is negligible for our defect population we have done a sensitivity analysis based on previous studies of hardening from loops. As discussed in Field et al.³¹ simple dispersed barrier hardening models suggest that the hardening under irradiation from loops is of the form $\Delta\sigma_y = A\sqrt{d}$ where A is a constant and d is the diameter of the defect. Now consider an error in diameter d defined as ε . The fractional error in $\Delta\sigma_y$ due to the error ε is $(\Delta\sigma_y(d + \varepsilon) - \Delta\sigma_y(d)) / \Delta\sigma_y(d) \approx \varepsilon/(2d)$, where the approximate equality holds for $\varepsilon \ll d$. For $\varepsilon = 1.7\text{ nm}$ (which is 2 pixels for our largest pixel sizes, see below) and $d = 21.4\text{ nm}$ (our average sizes of $a/2<111>$ and $a<100>$ defects), we get the fractional error in $\Delta\sigma_y$ as $1.7\text{ nm} / (2 * 21.4\text{ nm}) \approx 0.04$, which is well within the uncertainty of such microstructure-based analysis. However, for smaller defects this percentage error could clearly become larger. The error of diameter between ML results and human results appears to be approximately symmetrically distributed in positive and negative directions and independent of defect density, as shown in detail in SI section 4 and 5.

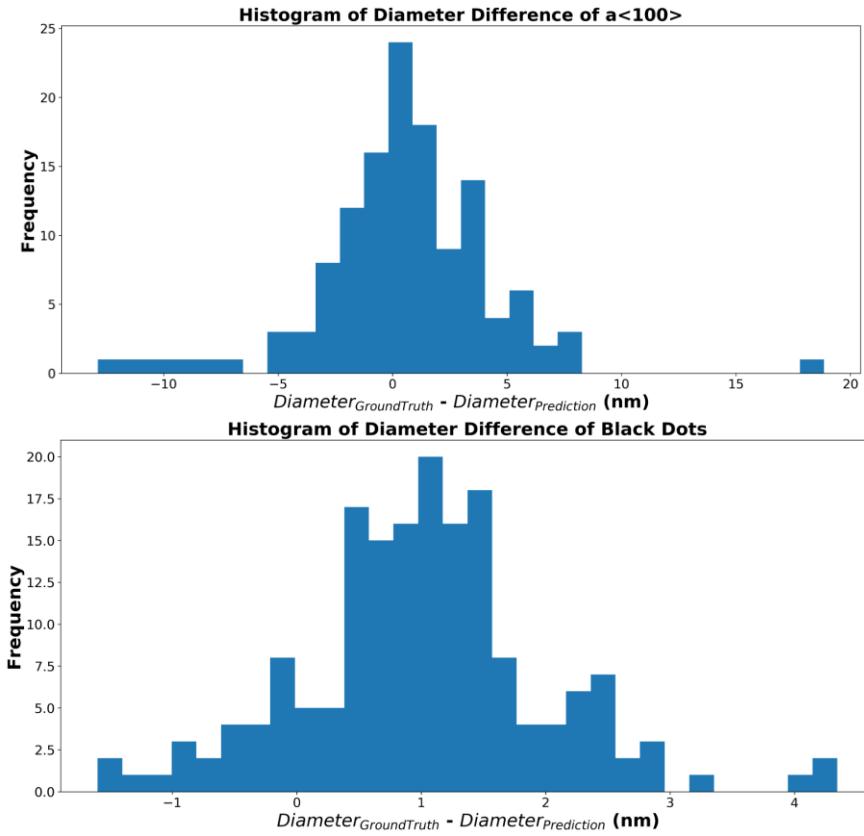
Q9a How big is a pixel?

The pixel to nanometer ratio is from 0.14 nm/pixel to 0.87 nm/pixel as already discussed on Page 17 Line 258 of the manuscript.

Q9b Is the error random or systematic?

The error is nearly random in the sense that it is not always high or low, nor does it show some unexpected groupings. As an illustration, below we show the diameter difference between ML and prediction for each class and shows that forms relatively symmetric and almost Gaussian looking distributions. Here we only considered those paired ground-truth labeling and ML prediction that are both correct in position prediction and defect type classification.





We have put the error analysis into SI section 4.

[Q9c Is the error dependent on defect density?](#)

We have examined the arithmetic mean of the error (ME), arithmetic mean of absolute value of the error (MAE), fractional ME (ME/ (arithmetic mean diameter of ground truth labeling), or FME) and fractional MAE (MAE/ (arithmetic mean diameter of ground truth labeling), or FMAE) in diameter for each of our 12 test images and plotted the values vs. the ground truth defect densities. The errors vs. density plots are shown below.

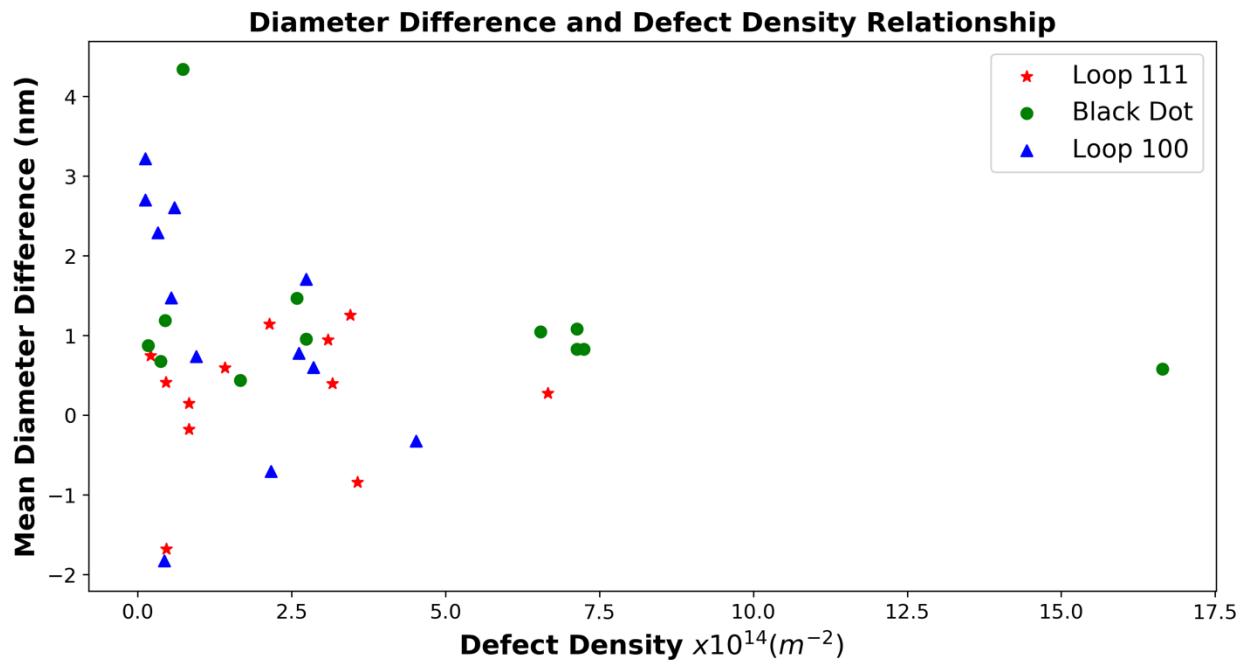
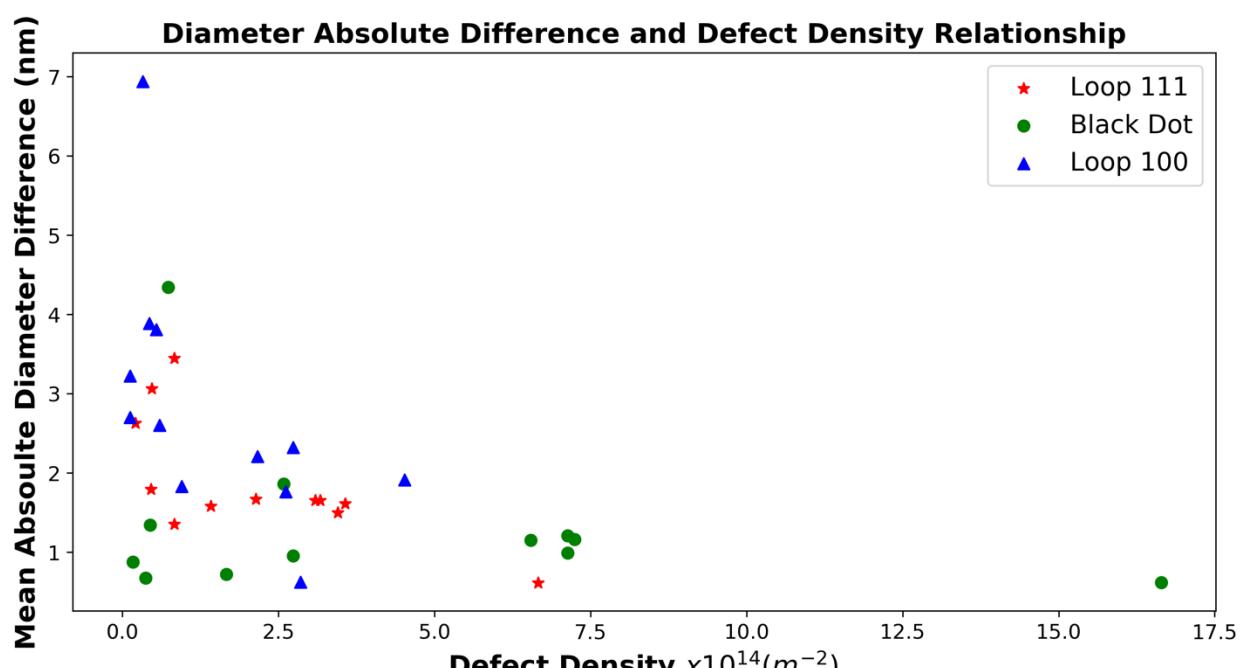


Figure 1. arithmetic mean of the error changes with the defect density



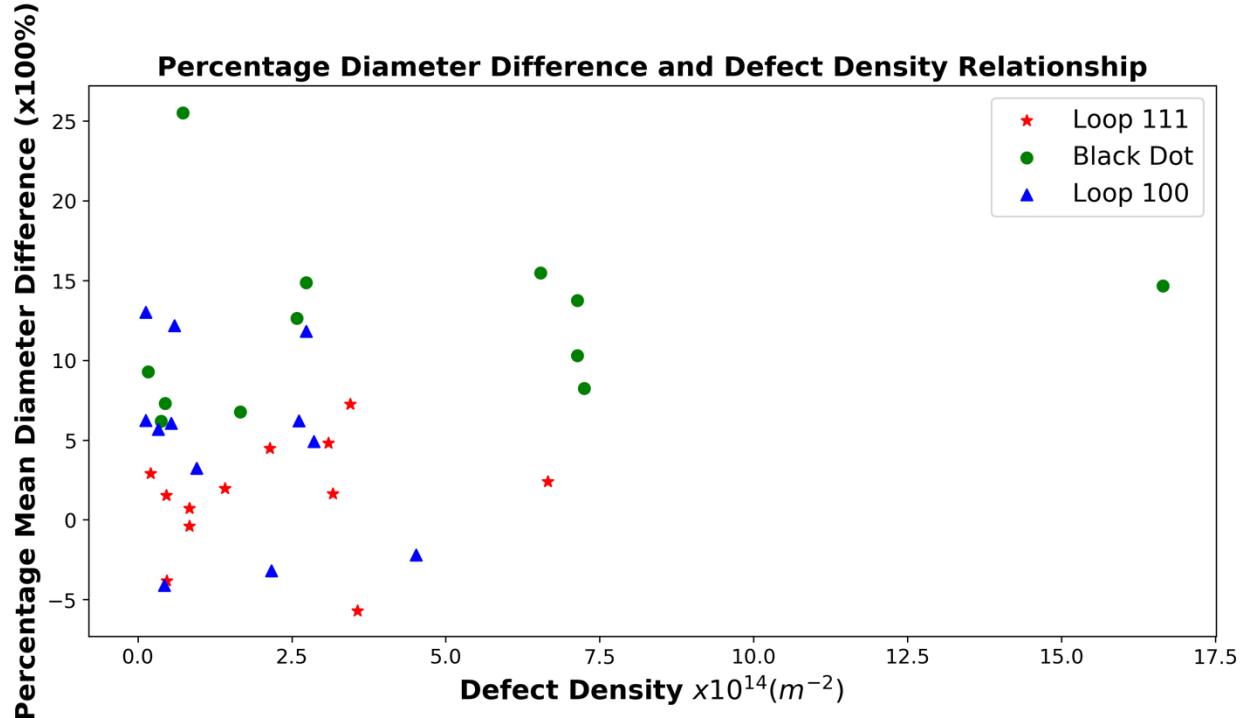


Figure 3. (arithmetic mean of the error) / (arithmetic mean diameter of ground truth labeling) changes with the defect density

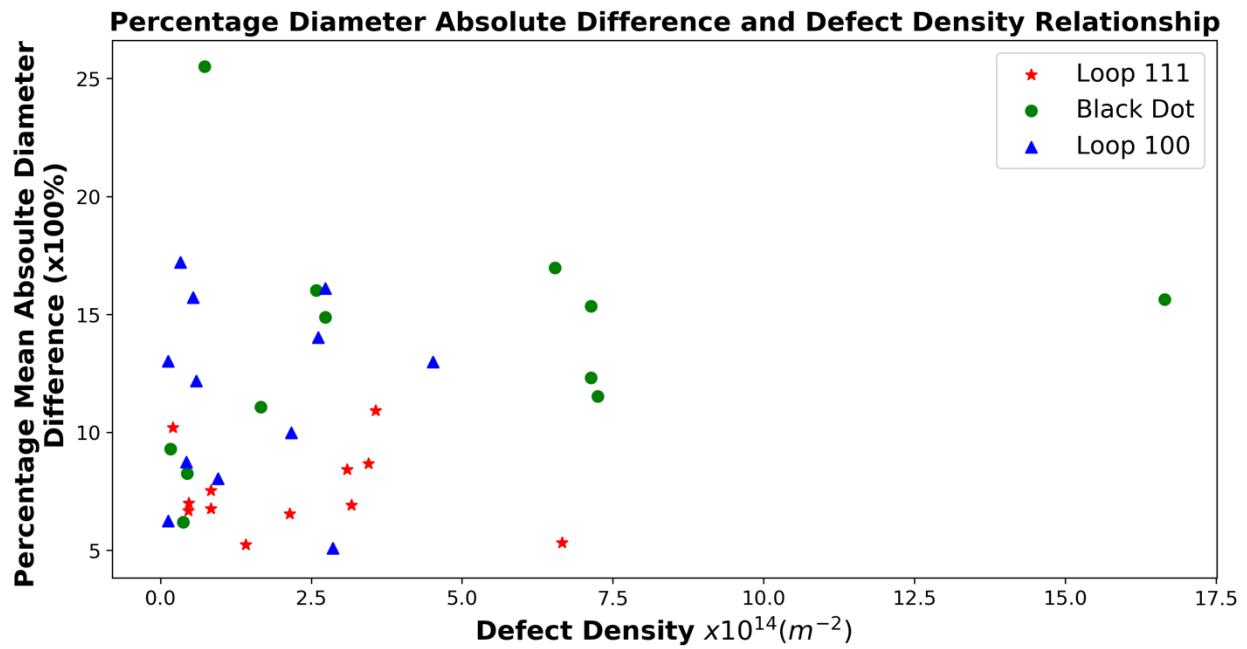


Figure 4. (arithmetic mean of absolute value of the error) / (arithmetic mean diameter of ground truth labeling) changes with the defect density

Figures 1 and 2 may show some trend of decreasing value with higher density, but we believe this is almost entirely due to the fact the low densities are correlated with larger loops, which

lead to larger absolute errors. To correct for this aspect, we measure fractional (or percentage) errors, shown in Figure 3 and 4. These show no discernable trend with density. We therefore believe that there is no trend of error with density when the confounding effect of defect size is removed.

We have put the error dependence of density discussion into SI Section 5.

We have added the following text on Page 20 Line 301-304 of the revised manuscript.

The errors of diameter between ML results and human results appear to be approximately symmetrically distributed in positive and negative directions and independent of defect density, as shown in detail in SI section 4 and 5.

Q10. Page 16: "Once the model is properly trained, it will yield a unique and reproducible labelling for every image"

This is a difficult statement for me. It is trivially true that the analysis run a second time will give the same result.

But if the training set is added to, or is prepared by a different researcher, then it is not necessarily true.

The labelling is therefore not an algorithmic property of the images, and we should not expect other groups to be able to reproduce it starting from the images alone.

As the authors themselves point out "even the same person may label the same defects differently even after a short break."

We agree with the reviewers' points here and that the meaning of this statement is not clear. We have modified the text somewhat to clarify and include the aspects mentioned by the reviewer. The new text on the Page 20-21 Line 311-323 reads as the following,

While the exact performance of the present automated approach compared to different human researchers is difficult to determine rigorously there is no doubt that the present approach is much more consistent. Previous studies have shown that different labelers tend to label defects in different ways and even the same person may label the same defects differently even after a short break^{2,5}. Such issues can make any given data analysis somewhat unreliable and make it difficult to integrate results across different teams and/or time periods in larger analysis efforts. However, once a machine learning model is properly trained, it will yield a unique and reproducible labeling for every image. If the community could converge on a single or small number of models this could greatly increase the reproducibility in labeling of STEM experiments. That said, models trained on different data and/or different human labeling could give different predictions, so establishing community accepted models is an important part of using these approaches to obtain more consistent results.

Q11. Page 19: "The automated analysis [takes] about 0.1 s/image"

How does this value scale with pixel count or defect density?

The Faster R-CNN model will resize each image into 1024 pixel x 1024 pixel, which is also the pixel size of our TEM images, so the time to run object detection is not influenced by a higher number of pixels. However, downsizing an image with a higher number of pixels will require extra time. We don't know the exact timing for this operation, and it could depend on many factors (e.g., machine speed, algorithms used, exact pixel count, etc.), but it is likely just tens of seconds at most for TEM images on a single processor. We, therefore, do not think it will be a significant issue. Assuming we use the same Faster RCNN hyperparameters the algorithm runs at approximately the same speed regardless of the number of defects in the image. If we have to alter the algorithm to use more proposal regions to detect more defects this could slow it down, but it is difficult to speculate on what would be necessary without detailed studies of the specific images.

Q12. Page 20: Data set preparation

Q12a. Any necessary parameters for CLAHE and Gaussian blur should be given explicitly.

The parameters used for CLAHE⁶ and Gaussian blur⁷ are all from the default parameter setting of scikit-images and the link to the detailed documentation of parameters used is put in the footnote.

We have added the text with two references of these parameters used on Page 25 Line 406-408 of the revised manuscript.

The parameters used for CLAHE⁴⁵ and Gaussian blur⁴⁶ are all from the default parameter setting of scikit-images and details can be found in the references given here for these methods.

Q12b. The use of additional filtered images in the RGB channels is a very neat idea. By how much does it actually help?

We do not have a direct assessment of the impact of adding different filtered images in the RGB channels. It is difficult to assess the model with just one channel as this would involve rewriting the basic Faster RCNN tools which take RGB as their standard input. We could provide the same images to each RGB channel, but this would require retraining the entire model, is a significant effort, and it is not clear how the separate channels would couple. Overall, we think it would be challenging to develop a rigorous and robust quantitative test to determine the impact of this particular choice of the training, and such testing is somewhat outside the main focus of the paper. We are therefore not adding any quantitative assessment of this approach relative to

⁶ https://scikit-image.org/docs/stable/api/skimage.exposure.html#skimage.exposure.equalize_adapthist

⁷ <https://scikit-image.org/docs/dev/api/skimage.filters.html#skimage.filters.gaussian>

others and we would ask for the reviewer's understanding that there is limited time for such assessment. That said, we think the choice we have made of using different filtered images in the RGB channels is likely beneficial for the working of Faster R-CNN as it provides an additional way for data augmentation to make the model more robust to noise.

Q13. Page 24: "We used grid search [to find] the best choice of these two values."

How is "best" defined? Is this in terms of table 1,2,3 or a combination?

We selected the best values based on maximizing the F1 scores of the testing images, which is the F1 score shown in row 3 in Table 2. We have clarified this in the text on Page 28 Line 464-465.

Q14. Page 25: "Fitting the contour with an ellipse function"

How is the fitting done?

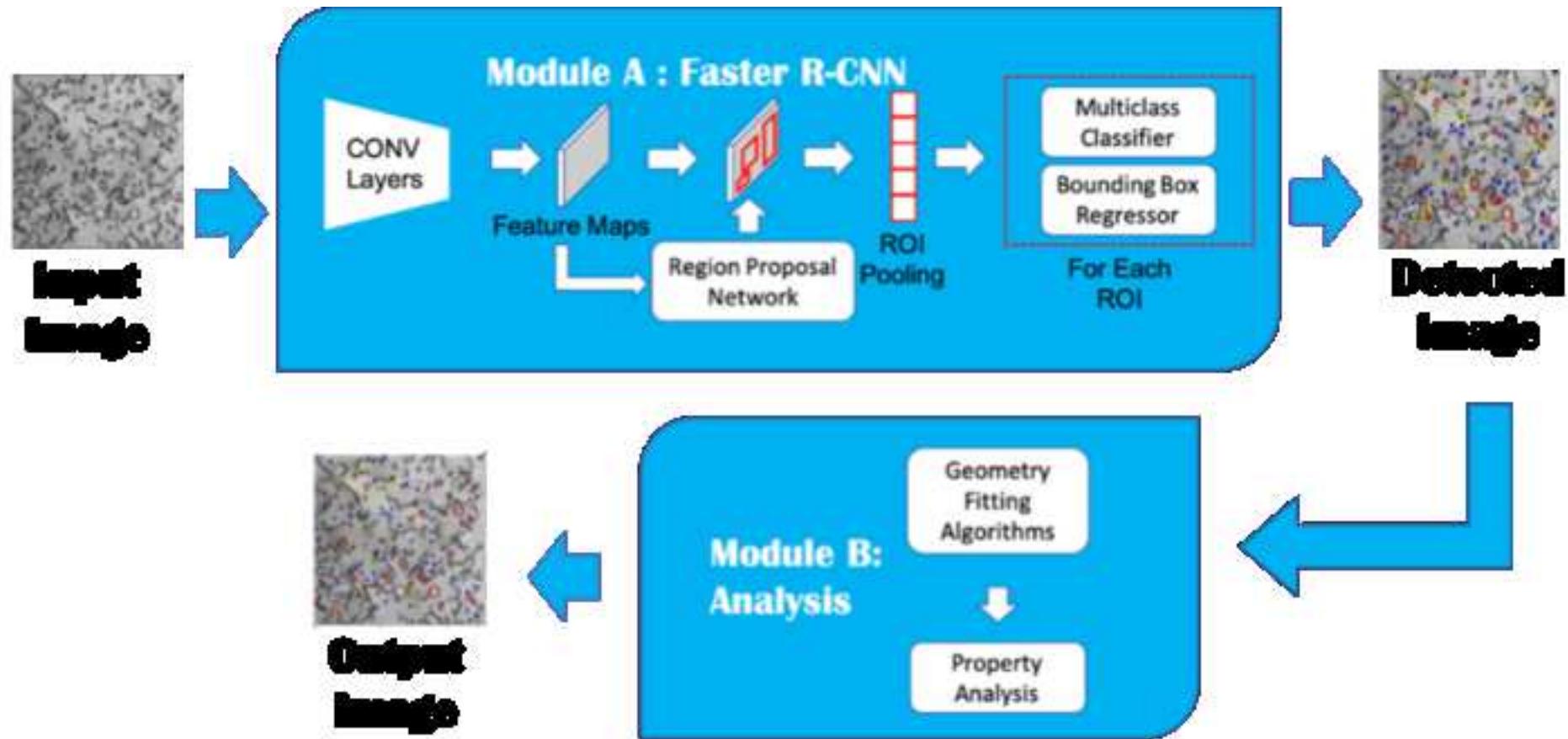
The overall fitting of the defect geometry is done in two steps. First, we use the watershed algorithm provided by OpenCV. The algorithm is applied to the cropped enlarged region of bounding boxes that contains the defect. We followed the official tutorial from OpenCV, and more details can be found there⁸. Then the boundary from the Watershed algorithm was fit to an ellipse. Fitting with an ellipse was useful to match the approach used by the community, obtain a well-defined shape with simple geometric descriptors, and smooth out the otherwise rather rough boundaries found by the Watershed algorithm. OpenCV's `fitEllipse()` function⁹ was called to perform the fit.

We have modified the mentioned text in the paper on Page 29 Line 481-488 to clarify what done for watershed and ellipse fitting,

Watershed methods were applied to find the boundary between defect pixels and background pixels. We followed the official tutorial from OpenCV for performing the watershed and details of the approach can be found there⁵⁰. We then fit the boundaries found from the Watershed algorithm to an ellipse. This fitting was done to match the approach used by the radiation defect analysis community, obtain a well-defined shape with simple geometric descriptors, and smooth out the otherwise rather rough boundaries found by the Watershed algorithm. The fitting was done with OpenCV's `fitEllipse()` function⁵¹.

⁸ https://docs.opencv.org/master/d3/db4/tutorial_py_watershed.html

⁹ https://docs.opencv.org/master/d3/dc0/group__imgproc__shape.html#gaf259efaad93098103d6c27b9e4900ffa



[Click here to view linked References](#)

1
2
3
4
5 **Multi Defect Detection and Analysis of Electron Microscopy**
6
7
8 **Images with Deep Learning**
9
10 3
11
12
13
14 4 Mingren Shen¹, Guanzhao Li¹, Dongxia Wu⁹, Yuhua Liu⁴, Hima Bharathi Adusumilli⁸,
15
16 5 Jacob Greaves¹, Wei Hao⁴, Nathaniel J. Krakauer⁴, Leah Krudy⁵, Jacob Perez⁴, Varun
18
19 6 Sreenivasan⁴, Bryan Sanchez⁶, Oigimer Torres⁷, Wei Li³, Kevin.G. Field², Dane Morgan¹
21
22 7
23
24
25 8 ¹ Department of Materials Science and Engineering, University of Wisconsin-
26
27
28 9 Madison, Madison, Wisconsin, 53706, USA
29
30
31 10 ² Materials Science and Technology Division, Oak Ridge National Laboratory;
32
33
34 11 Current address: Nuclear Engineering and Radiological Sciences Department,
35
36
37 12 University of Michigan-Ann Arbor
38
39
40 13 ³ Google, Mountain View, California, 94043, USA
41
42
43 14 ⁴ Department of Computer Sciences, University of Wisconsin-Madison, Madison,
44
45
46 15 Wisconsin, 53706, USA
47
48
49 16 ⁵ Department of Mathematics, Hope College, 141 E 12th Street, Holland Michigan,
50
51
52 17 49423, USA
53
54
55 18 ⁶ Department of Computer Science and Engineering, University of Puerto Rico at
56
57 19 Mayagüez, Mayaguez, 00682, Puerto Rico
58
59
60
61
62
63
64
65

1
2
3
4
5 20 ⁷ Electrical and Computer Engineering Department, University of Puerto Rico at
6
7
8 21 Mayagüez, Mayaguez, 00681, Puerto Rico
9
10 22 ⁸ Microsoft Corporation, Redmond, WA, 98052, USA
11
12
13 23 ⁹ Department of Mathematics, University of Wisconsin-Madison, Madison,
14
15
16 24 Wisconsin, 53706, USA
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5 **Abstract**
6
7
8 Electron microscopy is widely used to explore defects in crystal structures, but human
9
10 detecting of defects is often time-consuming, error-prone, and unreliable, and is not
11
12 scalable to large numbers of images or real-time analysis. In this work, we discuss the
13
14 application of machine learning approaches to find the location and geometry of different
15
16 defect clusters in irradiated steels. We show that a deep learning based Faster R-CNN
17
18
19
20
21
22 analysis system has a performance comparable to human analysis with relatively small
23
24
25 training data sets. This study proves the promising ability to apply deep learning to assist
26
27
28 the development of automated analysis microscopy data even when multiple features are
29
30
31 present and paves the way for fast, scalable, and reliable analysis systems for massive
32
33
34 amounts of modern electron microscopy data.
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5 36 INTRODUCTION
6
7
8 37 Electron microscopy (EM) is one of the most powerful tools for researchers to
9
10 38 extract and collect micrometer down to angstrom scale structural and morphological
11
12 39 properties of materials, including repeated structural units (e.g., unit cells of crystals) and
13
14 40 defected regions (e.g., grain boundary, impurities, defect clusters). Traditionally,
15
16 41 researchers have to manually label defects and repeatedly measure the relevant properties
17
18 42 to obtain statistically meaningful values, which is time-consuming, error-prone,
19
20 43 inconsistent, and hard to scale¹. The issue of scaling has become pressing as increasing
21
22 44 usage and advancement of EM techniques, such as high-speed detector and automated
23
24 45 sample exploration in EM, now generate massive amounts of image data (e.g., up to
25
26
27 46 thousands of images from a single experiment or condition can be generated in minutes)
28
29
30 47 which will keep increasing in the near future². Data on this scale cannot be practically
31
32 48 examined by humans, and automated approaches are therefore now necessary to utilize
33
34 49 the full power of modern EM. Not surprisingly, the EM community has developed many
35
36 50 accurate image data analysis tools that can be effectively deployed to accommodate the
37
38 51 large volume of EM data^{3,4}. However, due to the complexity of images of material
39
40 52 systems, these tools generally still need significant hand-tuning, and in some cases (like
41
42 53 counting defects of irradiated materials), human identification of each defect is still the
43
44 54 norm.

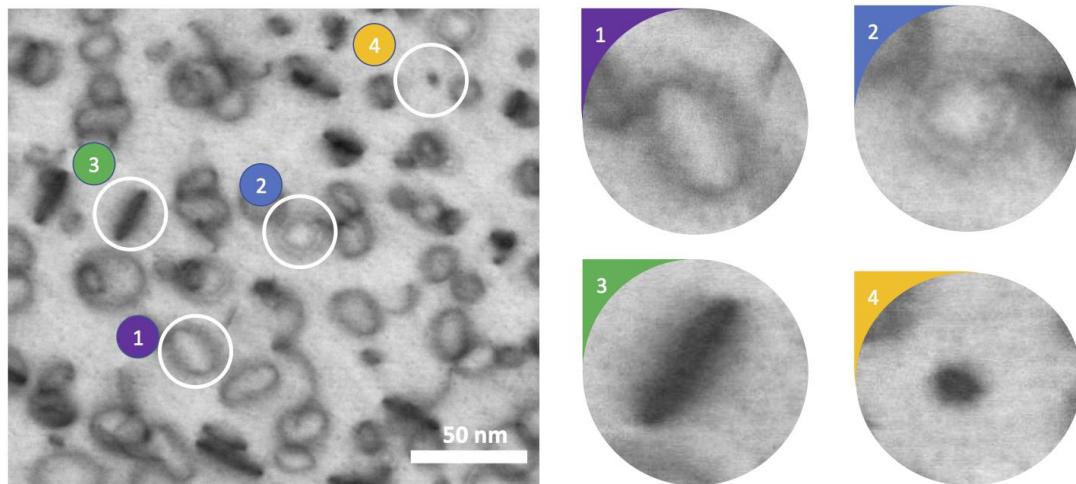
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5 Well-designed and well-tested automated analysis has proven to be significantly
6
7
8 more efficient, repeatable, and standardized than human analysis for discrete cases^{2,4,56}.
9
10 57 Developing automatic methods for EM image processing has drawn a great deal of
11
12 58 interest in the material science community. Automation efforts typically rely on
13
14 59 traditional computer vision technology, such as variance hybridized mean local
15
16 60 thresholding⁷, texture representation, and template matching methods like bag of visual
17
18 61 words (BoW)^{8,9}, key-point matching methods^{10,11}, Hessian-based Blob boundary
19
20 62 detection methods¹² and sometimes obtain better performance^{13,14} by utilizing tools from
21
22 63 broader areas, such as incorporating synthetic image dataset¹⁵ and using machine learning
23
24 64 methods like support vector machine¹⁶ and k-means clustering². However, these efforts
25
26 65 typically require extensive human tuning and/or are limited to specific tasks.
27
28
29
30
31
32
33
34
35
36
37 66 Recent developments (<10 years) in deep learning methods^{17,18} have
38
39 67 demonstrated that object detection in images can be automated with minimal
40
41 68 hyperparameters and yield human or even better than human levels of performance.
42
43
44
45 69 These frameworks are now being adapted towards finding defects in metals, and
46
47 70 particularly nuclear materials, including the automated detection of dislocation loops,
48
49
50 71 cavities, precipitates, and line dislocations^{5,19–21}. Generally, there are three different
51
52 72 approaches for applying deep learning frameworks to defect detection in microscopy
53
54
55 73 images. The first is using the combination of both traditional techniques and deep
56
57
58 74 learning tools⁵, for example, Li et al. develops an analysis model that includes a local

1
2
3
4
5 75 visual content descriptor widely used in computer vision called Local Binary Patterns
6
7
8 76 (LBP)²² descriptor, feature selecting methods called AdaBoost²³, and Convolutional
9
10 77 Neural Network (CNN) module to screen candidate bounding boxes to obtain the best
11
12 78 performance. This approach is more like an intermediate stage of applying deep learning
13
14 79 since it does not follow the complete end-to-end pattern of deep learning practice²⁴ but it
15
16 80 helped show that value of new deep learning methods for this class of problems. The
17
18 81 second approach relies on the encoder-decoder framework to find features²⁵ in EM
19
20 82 images. Examples include using weakly supervised learning methods of encoder-decoder
21
22 83 to study the local atom movements²⁶, U-Net to study nanoparticle segmentations²⁷, and a
23
24 84 modified U-Net framework to segment defects in STEM images of steels¹⁹. The encoder-
25
26 85 decoder framework can extract the most relevant information in images and use the
27
28 86 extracted inner state to do other tasks, but the performance of encoder-decoder
29
30 87 framework relies on the extracted inner state and good performance requires careful
31
32 88 training²⁵. The third category of methods are using mature object-detection frameworks,
33
34 89 for example, Chen et al. used Mask R-CNN to study the microstructural segmentation of
35
36 90 aluminum alloy²⁸ and Anderson, et al. used Faster R-CNN to study helium bubbles in
37
38 91 irradiated X-750 alloy²⁰. Here we explore the first use of similar mature object-detection
39
40 92 deep learning methods as this third category, specifically Faster R-CNN, to obtain
41
42 93 properties of dislocation loops with varying Burgers vector and habit plane in neutron-

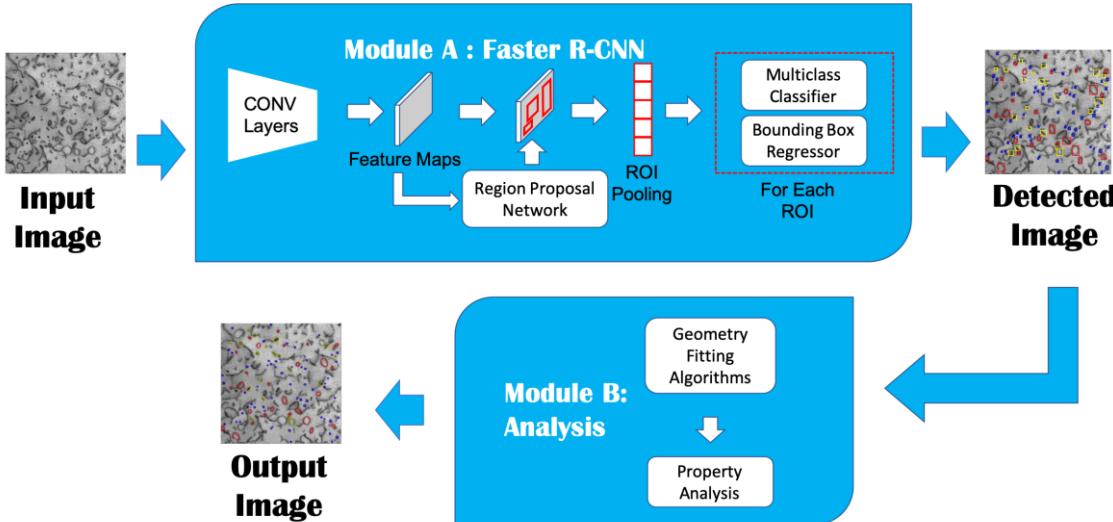
1
2
3
4
5 94 irradiation related iron-chromium-aluminum (FeCrAl) materials, These materials are
6
7 95 important for the development of next generation nuclear reactors²⁹⁻³¹.
8
9
10 96 Analyzing the locations and sizes of defects in materials that have undergone
11
12
13 97 irradiation is a widely used application of electron microscopy. In such studies, the key
14
15
16 98 properties are the total number and distribution of each type of defect. Typical defects of
17
18
19 99 interest include grain boundaries, precipitates, dislocation lines, dislocation loops,
20
21
22 100 stacking fault **tetrahedra**, cavities (voids, bubbles), and co-called “black-spot” defects,
23
24
25 101 which are small defect clusters of interstitials and sometimes vacancies^{1,32}. For this study,
26
27
28 102 we focus on the dislocation loops formed within a ferritic alloy, where the loops exist on
29
30
31 103 specific habit planes that manifest themselves with different morphologies due to the
32
33
34 104 projection of a 3D volume imaged using EM³³. Typical microstructural images of
35
36
37 105 irradiated ferritic steels contain four prominent types of defects: (1) open ellipse loops
38
39
40 106 (single ring edge), (2) open ellipse loops (double ring edges), (3) closed solid elliptical
41
42
43 107 loops, (4) closed circular solid dots³³. Figure 1 shows a sample STEM image containing
44
45
46 108 all four morphologies of loops obtained from a ferritic alloy irradiated in a materials test
47
48
49 109 reactor. In this paper, we used a modern deep learning-based object detection model
50
51
52 110 called Faster Regional CNN (Faster R-CNN)³⁴, a widely used deep learning based object
53
54
55 111 detection model¹⁷. We use the Faster R-CNN to develop an automatic defect detection
56
57
58 112 system for all four morphologies commonly observed in irradiated steels with a body-
59
60
61 113 centered cubic structure and then additional post-processing to analyze their geometrical
62
63
64
65

1
2
3
4
5 information (specifically, size and areal density). This paper serves to demonstrate the
6
7 power of deep learning-based computer vision models for material image studies and
8
9 suggests the possibility that most aspects of defect analysis may soon be practically
10
11 automated, and many, if not all, handcrafted feature-based methods may be replaced by
12
13 deep learning methods.
14
15
16
17
18
19
20



120
121 Figure 1. Selected bright field scanning transmission electron microscopy (STEM) image of an irradiated ferritic alloy
122 showing four common morphologies of dislocation loops: (1) open ellipse loops (single ring edge), (2) open ellipse
123 loops (double ring edges), (3) closed elliptical solid loops, (4) closed circular solid dots. Open single edge ellipse loops
124 (1) are dislocation loops with a Burgers vector of $a_0/2 \langle 111 \rangle$. Open double edge ellipse loops (2) and closed elliptical
125 solid loops (3) are dislocation loops with a Burgers vector of $a_0 \langle 100 \rangle$. Closed circular solid dots (4) are black dot
126 defects with a Burgers vector of either $a_0/2 \langle 111 \rangle$ or $a_0 \langle 100 \rangle$. Image size: Primary image is 290 × 290 nm; inset scales
arbitrary.

1
2
3
4
5 129 Faster R-CNN is a CNN based end-to-end deep learning object detection model
6
7
8 130 that outputs both the object position and its class³⁴. As shown in Figure 2, Faster R-CNN
9
10 131 is a two-stage detector where the region proposal network (RPN) proposes Region of
11
12 132 Interest (ROI), and the following ROI regressor and classifier will fine tune the final
13
14 133 output results including the size and position of the object contained bounding boxes and
15
16 134 the corresponding object label³⁴. Given an image, the shared convolutional layers will
17
18 135 extract a feature map from the input image by performing a series of convolution and
19
20 136 max pooling operations. Then based on the extracted feature map, the RPN will put a set
21
22 137 of predefined anchor boxes on the feature map and output the probability of whether the
23
24 138 anchor box belongs to an object of interest or plain background. It worth mentioning that
25
26 139 RPN ignores the specific object class of each bounding box and the following ROI
27
28 140 regressor and classifier are responsible for the specific class and refined location of the
29
30 141 objects. The refining network predicts certain object labels and refines the size and
31
32 142 position of each bounding box based on the feature map generated by the ROI-pooling
33
34 143 layers³⁵. The RPN and ROI components are trained jointly to minimize the loss function
35
36 144 sums from both of them³⁴. After the Faster R-CNN module A, those images with detected
37
38 145 defects are sent to module B to extract geometric information such as defect diameters, as
39
40 146 shown in Figure 2.
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57 147



148 Figure 2. Schematic flow chart of proposed deep learning based automated detection approach. Input micrograph
 149 images go through the pipeline of Module A—Faster R-CNN Detector, Module B—Image Property Analysis. After
 150 Module A, the loop locations and bounding boxes are identified and then for each identified bounding box, geometry
 151 fitting algorithms are called to determine the defect shape and size in Module B.

152
 153
 154 For details on the data sets, training approaches, and methods of defect
 155 identification and analysis, please see Methods section. Please see Data section for
 156 summary of all data made available.

158 RESULTS

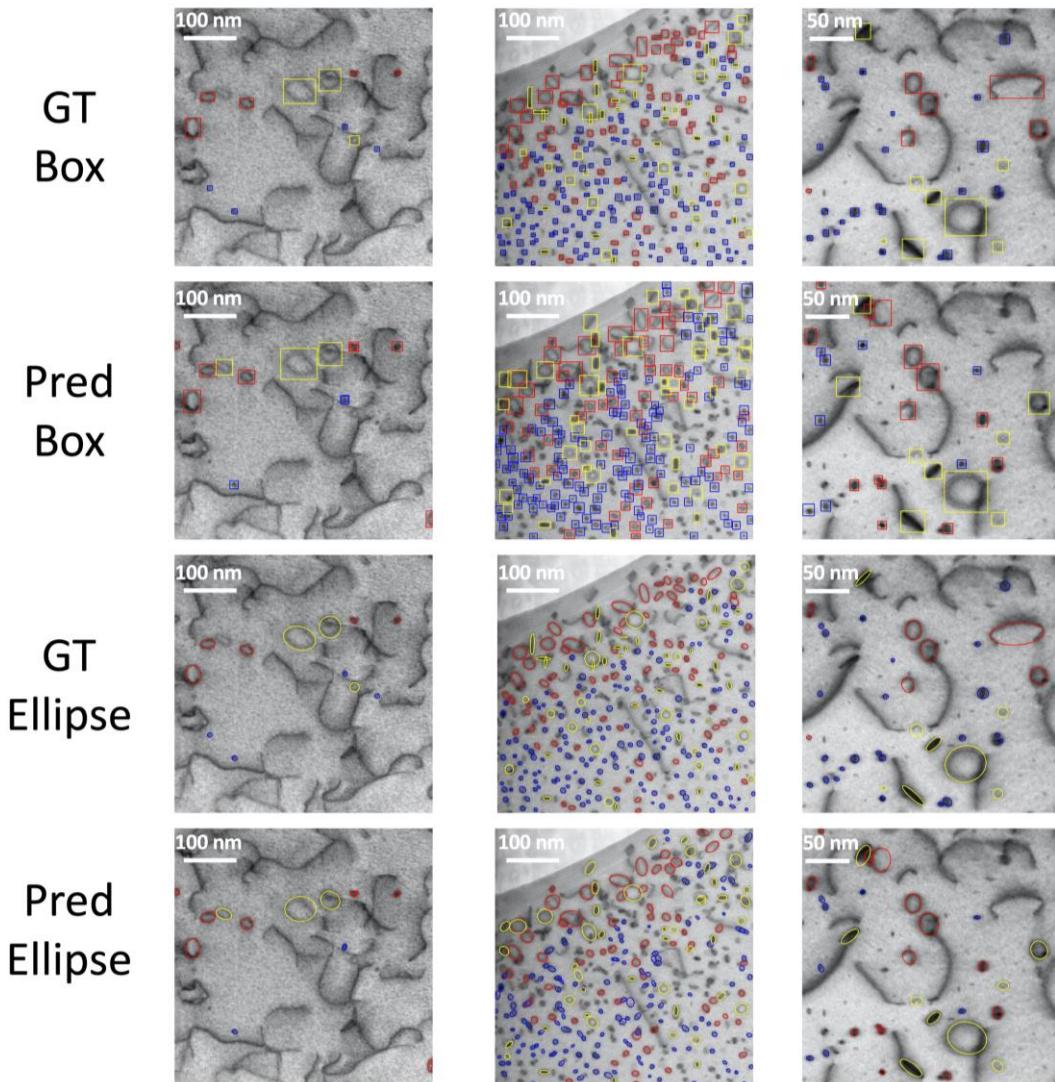
159 To assess the machine predictions, four types of approaches were taken. The first
 160 approach was a qualitative comparison of machine to human labeled images, where we
 161 looked for large fractions of errors, e.g., more than 40%, and for trends in errors that

1
2
3
4
5 162 might indicate a major issue but made no attempt to quantify agreement. This assessment
6
7 163 tests all aspects of the model as it compares to the ground truth human results, which
8
9 164 include the bounding box predictions (the defect detection part of Module A in Figure 2),
10
11 165 the defect type identifications (the categorization part of Module A in Figure 2), and the
12
13 166 geometric shape determination (Module B in Figure 2). The second assessment approach
14
15 167 was a quantitative assessment of the ability to identify a defect, regardless of defect types.
16
17 168 This assessment tested the defect detection part of Module A (see Figure 2). This
18
19 169 assessment was a binary categorization problem and success was quantified with
20
21 170 precision, recall, and F1 score. The third assessment was a quantitative assessment of the
22
23 171 ability to identify a defect type once a defect had been correctly identified and tested the
24
25 172 categorization part of Module A (see Figure 2). This assessment was a three-category
26
27 173 categorization problem and was quantified using the confusion matrix with precision,
28
29 174 recall and F1 calculated for each class. Finally, the fourth assessment was a quantitative
30
31 175 assessment of the ability to quantify the geometric properties of defects. This assessment
32
33 176 tested the geometric analysis of Module B (see Figure 2) and compared machine and
34
35 177 human predictions of average and standard deviations in size and areal density for each
36
37 178 defect type. We discuss each of the four assessments below and label them assessment 1-
38
39 179 4 for clarity. In all cases the comparisons are made on the test data set described in
40
41
42 180 Methods section.

1
2
3
4
5 **Assessment 1.** After feeding the images into the Faster R-CNN detectors, the
6
7 resulting detections were plotted on the original images. As shown in Figure 3, the red
8
9 circles represent the dislocation loops with a Burgers vector of $a_0/2\langle 111 \rangle$ (Type 1 in
10
11 Figure 1), while the yellow and blue circles represent $a_0\langle 100 \rangle$ direction loops (Type 2
12
13 in Figure 1) and “black dot” defects (Type 4 in Figure 1) respectively. The data
14
15 from both human-labeled and machine detected results are plotted in the same manner.
16
17
18 More comparisons can be found in Supplement Information Section 1. To a human
19
20 observer the machine results show strong correlation of bounding box location, defect
21
22 type identification, and defect shape with the ground truth human labeling which
23
24 indicates the effectiveness of the proposed automatic defect detection system.
25
26
27
28
29
30
31 **Assessment 2.** The performance of the detection part of Module A (see Figure 2)
32
33
34 of the trained model was evaluated in terms of precision, recall, and F1 score by
35
36 comparing the detected result with the human labeled result of the 12-image testing set,
37
38 as shown in Figure 4. The precision describes the percentage of all machine predicted
39
40 bounding boxes that are judged to have correct positions, and the recall value describes
41
42 the percentage of all human labeled defects that are identified as in a bounding box by the
43
44 machine algorithm. F1 is the harmonic mean of the precision and recall which can be
45
46 used to assess the overall performance of the defect location task³⁵. The IoU (Intersection
47
48 over Union) method was used to determine if a given defect was identified by a bounding
49
50 box and is described within the provided Methods section. The cutoff IoU, which must be
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5 201 exceeded to consider the bounding box to have identified the defect, is a hyperparameter
6
7 202 that can be fine-tuned based on the purpose of the object detection task³⁵. Figure 4
8
9 203 showed a drop in performance as the cutoff IoU increased. This trend agreed with
10
11
12 204 expectations as the higher cutoff IoU meant it was harder for the predicted bounding box
13
14 205 to be judged successful. However, setting the cutoff IoU to an extremely small threshold
15
16 206 could lead to the problem that the predicted bounding boxes are associated with defects
17
18 207 for which only a small part of the defect is actually in the bounding box, which will likely
19
20 208 cause problems in the defect identification (Module B) of our model. As a compromise,
21
22 209 for all the further assessments in this paper, we used cutoff IoU = 0.4 to determine when
23
24 210 the machine predictions were considered to match a given defect. This choice kept
25
26 211 nearly optimal performance of the detector (based on Figure 4) and an adequately
27
28 212 demanding standard for predictions.

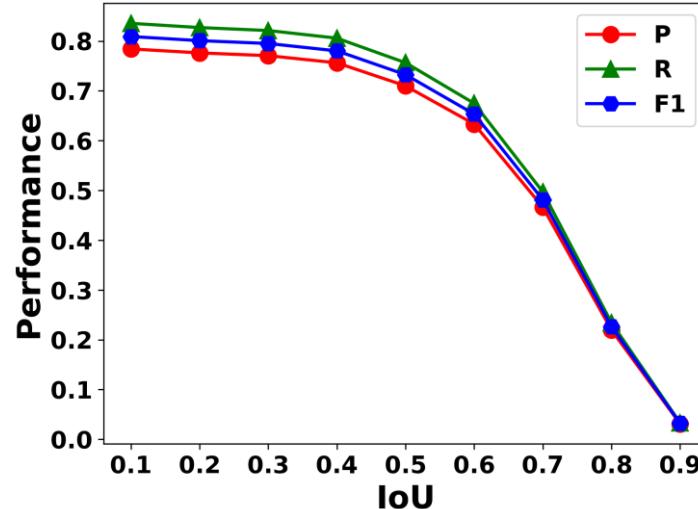
Image 1 Image 2 Image 3



215 Figure 3. Selected data images to show the detector performance and the fitted ellipse of our automatic

216 analysis system. These three test images are selected from the test dataset of 12 images (see Methods). The “Ground
217 Truth (GT)” shows the bounding box and ellipse human labeling (colored by defect type), the “Prediction (Pred) Box”

1
2
3
4
5 218 shows the predicted bounding boxes (colored by defect type), and the “Prediction (Pred) Ellipse” shows the resulting
6
7
8 219 fits to the specific defect geometry (colored by defect type as described in the text).
9
10
11
12 220



33
34 221
35
36 222 Figure 4. Summary of defect location recognition performance of all types of defects evaluated using precision and
37
38 recall metrics, regardless of defect types. The test set contained 12 images, and, for all IoU (Intersection over Union)
39
40
41 223 values and we used a threshold confidence score 0.25 for Faster R-CNN output. (see Method Section).
42
43
44
45 225
46
47 226 **Assessment 3.** Table 1 shows the confusion matrix of the predictions made by
48
49 Faster R-CNN detector evaluating its capability to correctly categorize defects. Each row
50
51 in the confusion matrix represents a class that is predicted by the detector, and each
52
53 column represents a class labeled by human researchers. The diagonal elements of the
54
55 228 table represent the correct classification made by the detector and off diagonal elements
56
57
58 229
59
60
61
62
63
64
65 230

1
 2
 3
 4
 5 231 represent errors of different types. We also show the percentage accuracy of each type of
 6 defect in parentheses. The 76%, 87%, and 94% accuracy indicates that once the Faster R-
 7
 8 232 CNN model locates the defect, it can classify the type of defect based on their
 9
 10 233 morphology within the image with good accuracy, although some improvement of the
 11
 12 234 76% value is likely possible for the $a_0\langle 100 \rangle a<100>$ loops. We also report the
 13
 14 235 classification performance using precision, recall and F1 score in Table 2. Given inherent
 15
 16 236 errors of human performance we take scores for precision, recall and F1 of 0.78 as
 17
 18 237 approximately the upper limit that can be obtained with the present labeling. Table 2
 19
 20 238 shows F1 from about 0.65 to 0.78, which demonstrates significant capabilities but is
 21
 22 239 likely less than can be achieved, suggesting opportunities for further improvements.
 23
 24 230
 25
 26 231
 27
 28 232
 29
 30 233
 31 234
 32
 33 235
 34
 35 236
 36
 37 237 Table 1. Summary of the classification performance for each type of defects at cutoff IoU 0.4. Values in parenthesis
 38
 39 238 give the % each number represents of the total number of defects in that class as determined by the human labeling.
 40
 41
 42 239
 43
 44
 45
 46
 47
 48

	$a_0/2\langle 111 \rangle$ Loop	Black Dot	$a_0\langle 100 \rangle$ Loop
$a/2\langle 111 \rangle$ Loop	239 (87.2%)	21	14
Black Dot	17	416 (94.3%)	8
$a\langle 100 \rangle$ Loop	33	13	166 (78.3%)

49 244
 50
 51 245 Table 2. The performance report for each class.
 52
 53
 54

	$a_0/2\langle 111 \rangle$ Loop	Black Dot	$a_0\langle 100 \rangle$ Loop
Precision	0.73	0.65	0.62
Recall	0.83	0.71	0.72
F1	0.78	0.68	0.67

1
2
3
4
5 246
6
7

8 **Assessment 4.** The second Module provides geometric information for each
9
10 247 defect through fitting ellipses. While the fits can provide a range of detailed information,
11
12 248 we are particularly interested in the arithmetic mean and the associated standard deviation
13
14 249 of the defect diameter as well as the areal density in an image for each type of defect.
15
16 250 These values are commonly quoted values in literature within irradiated materials studies.
17
18 251 Table 3 compares the human labeled arithmetic mean diameters and areal densities to the
19
20 252 ones predicted by the automatic analysis system. The discrepancy of arithmetic mean
21
22 253 diameter between the human labeled ground truth and predictions is within 10% in all
23
24 254 cases, which is considerably less than might be expected for variation among different
25
26 255 humans² and we consider a strong success. Furthermore, the errors in arithmetic mean
27
28 256 diameters are in the range 0.7-1.1 nm, which corresponds to a range of two to nine pixels
29
30
31 257 (based on the range 0.14nm/pixel to 0.87nm/pixel for our test data, see SI Section 1). The
32
33 258 errors of about 1 nm correspond to about 5-10% for our data which is somewhat larger
34
35 259 than might be expected from direct labeling errors on 10-15nm. Thus, it is unlikely that
36
37 260 any human labeling is meaningfully accurate to much below this level. However, the
38
39
40 261 human and machine learning black dot radii do not fall within a 95% confidence interval,
41
42
43 262 suggesting that the algorithm does not yield exactly the same means as the human ground
44
45
46 263 truth. Some errors will come from the machine detection (failures in precision and recall,
47
48
49 264 see Figure 4) and defect type assignment (see Table 1). Additional errors are associated
50
51
52 265 with the defect type assignment (see Table 1). Additional errors are associated

1
 2
 3
 4
 5 266 with intrinsic errors in the machine and human ellipse labeling, where both have some
 6
 7 267 uncertainty due to ambiguity or variances in the morphology of defects in images. In
 8
 9 268 particular, some defects are not well fit by an ellipse (e.g. some have a more rectangular
 10
 11 269 shape, as can be seen in Figure 3), making this form of labeling difficult for both human
 12
 13 270 and machine. Another error to consider is that as the number of pixels per feature goes
 14
 15 271 down, the intrinsic error due to the resolution (pixel/nm) will artificially go up. For
 16
 17 272 instance, a 100 nm loop where the resolution is 1 pixel/nm where the labeling is off by 1
 18
 19 273 pixel will yield a 1% error. If the labeling is off by 1 pixel for a 5 nm loop, the error will
 20
 21 274 be 20% even though the per pixel error is the same. Seeing as the black dots are all of
 22
 23 275 small [arithmetic](#) mean diameter (<10 nm), they will intrinsically have a higher error
 24
 25
 26 276 compared to the other classes where the diameters are 2-3 times larger.
 27
 28
 29
 30
 31
 32
 33
 34 277
 35
 36
 37 278 Table 3. Comparison of [arithmetic](#) mean defect diameter and standard deviation of mean loop diameter between ground
 38
 39
 40 279 truth labeling and our automatic analysis model prediction with an IoU of 0.4. The values in parenthesis are the relative
 41
 42
 43
 44
 45 280 percentage error between ground truth human labelling results and the automatic analysis results.
 46
 47

Defect Type	Ground Truth			Automatic Analysis Model		
	Arithmetic Mean diameter (nm)	Standard Deviation of Mean Diameter (nm)	Areal density (m ⁻²)	Arithmetic Mean diameter (nm)	Standard Deviation of Mean Diameter (nm)	Areal Density (m ⁻²)
a₀/2 (111) Loop	22.4	0.7	1.77×10 ¹⁴	23.1 (3.1%)	0.8	2.21×10 ¹⁴ (24.9%)
Black Dot	8.2	0.1	3.41×10 ¹⁴	9.1 (10.9%)	0.2	4.98×10 ¹⁴ (46.0%)

<u>a₀(100)</u>	20.3	0.8	1.32×10 ¹⁴	22.4 (10.3%)	0.9	1.79×10 ¹⁴ (35.6%)
---------------------------	------	-----	-----------------------	-----------------	-----	----------------------------------

DISCUSSION

The above results demonstrate that that the trained model potentially performs well enough to replace human in a workflow on similar types of data. The precision and recall values for assessing detection in the range 62-83% which are comparable or less than human variation² from previous assessments. The machine defect type misidentifications are at the level of 10-25% (see Table 1), and a significant fraction of this variation may also be due to ground truth ambiguities or errors. The final machine predicted diameters are within a nanometer, approximately 2 pixels in images, which is a level of error that is considered negligible in terms of impact on material properties. To further clarify that the error is negligible for our defect population we have done a sensitivity analysis based on previous studies of hardening from loops. As discussed in Field et al.³¹ simple dispersed barrier hardening models suggest that the hardening under irradiation from loops is of the form $\Delta\sigma_y = A \sqrt{d}$ where A is a constant and d is the diameter of the defect. Now consider an error in diameter d defined as ε. The fractional error in Δσ_y due to the error ε is $(\Delta\sigma_y(d + \varepsilon) - \Delta\sigma_y(d))/\Delta\sigma_y(d) \approx \varepsilon/(2d)$, where the approximate equality holds for ε < d. For ε = 1.7 nm (which is 2 pixels for our largest pixel sizes, see below) and d = 21.4 nm (our average sizes of a/2<111> and a<100> defects), we get the fractional error in Δσ_y as 1.7 nm / (2 * 21.4 nm) ≈ 0.04, which is well

1
2
3
4
5 300 within the uncertainty of such microstructure-based analysis. However, for smaller
6 301 defects this percentage error could clearly become larger. The errors of diameter between
7
8 302 ML results and human results appear to be approximately symmetrically distributed in
9
10 303 positive and negative directions and independent of defect density, as shown in detail in
11
12 304 SI section 4 and 5. Furthermore, previous studies indicate that the differences of
13
14 305 arithmetic mean diameter between different human labelers can be comparable or larger
15
16 306 than values found here between the ML and human results⁵. The discrepancy in areal
17
18 307 densities is somewhat larger than might be intuitively expected just from the percentage
19
20 308 error in the arithmetic mean diameters. However, additional errors are introduced by the
21
22
23 309 exact definition of areal density (see Methods section) and the additional errors
24
25
26 310 introduced by the imperfect precision and recall.

27
28
29
30
31 311 While the exact performance of the present automated approach compared to
32
33 312 different human researchers is difficult to determine rigorously there is no doubt that the
34
35
36 313 present approach is much more consistent. Previous studies have shown that different
37
38
39 314 labelers tend to label defects in different ways and even the same person may label the
40
41
42 315 same defects differently even after a short break^{2,5}. Such issues can make any given data
43
44
45 316 analysis somewhat unreliable and make it difficult to integrate results across different
46
47
48 317 teams and or time periods in larger analysis efforts. However, once a machine learning
49
50
51 318 model is properly trained, it will yield a unique and reproducible labeling for every
52
53
54 319 image. If the community could converge on a single or small number of models this

1
2
3
4
5 320 could greatly increase the reproducibility in labeling of STEM experiments. That said,
6
7 321 models trained on different data and/or different human labeling could give different
8
9 322 predictions, so establishing community accepted models is an important part of using
10
11 323 these approaches to obtain more consistent results.

16 324 The approach applied here is readily scalable to very large data sets. Analyzing a
17
18 325 single image with our model on a reasonablele state of the art GPU (NVIDIA's GeForce
19
20 326 GTX 1080 GPU) takes about 0.1s, so analyzing all the images in a typical experiment can
21
22 327 be done easily in minutes, even less if multiple GPUs are used and as GPU and related
23
24 328 processors (e.g., TPU) continue to get faster. As large scale distributed cloud service
25
26 329 provider like Google, Amazon and Microsoft are providing cloud service for deep
27
28 330 learning applications with GPU machines³⁶, it would be easy to scale to process even
29
30 331 larger amount of data. Furthermore, significant speedup can likely be obtained if desired.
31
32
33 332 We developed the system with the Python code language and the ChainerCV deep
34
35 333 learning framework, both of which were chosen for ease of development not for the high-
36
37 334 performance in deployment. Replacing Python with C/C++ or using high-performance
38
39 335 deep learning frameworks, e.g. Caffe³⁷, could potentially accelerate the prediction speed
40
41 336 of the current model. In particular, the deep learning community is actively designing
42
43 337 new methods to accelerate the running speed of model e.g. model compression, weight
44
45 338 sharing, or parameter pruning³⁸ which could also boost the speed of ours. As an example
46
47 339 of how fast deep learning AI algorithms can be, researchers from Google have recently

1
2
3
4
5 340 applied deep learning models for cancer diagnosis on data during the actual process of
6
7
8 341 conducting an optical microscopy experiment³⁹.
9
10 342 The approach applied here is also readily adapted to new defect types and
11
12 343 systems. The present model was trained with only a relatively small amount of training
13
14 344 data due to the use of transfer learning⁴⁰. With only modest additional data sets (e.g., on
15
16 345 the scale of thousands of defects or possibly fewer) and a few rounds of further training
17
18 346 as described in Section II of SI, researchers could likely extend the present model to
19
20 347 more defects (e.g. separating the two orientations of 111 loops or adding voids,
21
22 348 preexisting dislocations, etc.), different imaging conditions (e.g., changes in microscopes,
23
24 349 imaging modes, orientation, focus, etc.), and different materials (e.g. other metal alloys).
25
26
27
28
29
30
31
32
33
34 350 There are several areas where significant improvements may be obtainable. The
35
36
37 351 first is that the use of real-world data in the study has led to significant time spent
38
39 352 labeling and introducing unavoidable human biases and errors into the deep learning
40
41
42 353 model and its assessment. However, it is possible that simulated images could be both
43
44
45 354 more accurately labeled and generated in large volume, potentially allowing much more
46
47
48 355 accurate models to be trained.
49
50
51 356 The second area where significant improvement is likely is that deep learning
52
53
54 357 methods for object detection continue to evolve rapidly. In particular, deep learning
55
56
57 358 segmentation models¹⁷, which learn a label for every pixel, could be equally or more
58
59
60 359 accurate and remove the step of fitting contours in a bounding box to get geometric
61
62
63
64
65

1
2
3
4
5 360 information. Such an approach was applied recently to automatically detect information
6
7
8 361 about dislocation lines, precipitates and voids in STEM images¹⁹.
9
10
11 362

12
13 363 **Conclusion**
14
15

16 364 This study demonstrated a practical deep learning based automatic STEM image
17
18 365 defect detection system implemented by incorporating Faster R-CNN for detection and
19
20 366 watershed flood algorithm for geometry fitting. Compared with other models proposed
21
22 367 before, our model reduced the training effort by utilizing only one module for detection
23
24 368 and expanded capability to simultaneously recognize multiple classes of defects. The
25
26 369 approach developed here achieved reasonably reliable performance, with an F1 score of
27
28 370 0.78, and predicted sizes and areal densities within the uncertainty of results from human
29
30 371 researchers. The automated analysis on NVIDIA's GeForce GTX 1080 GPU processor is
31
32 372 about 0.1 s/image, hundreds of times faster than human analysis (≥ 1 minute/image), and
33
34 373 trivially parallelizable and scalable on more processors. The model can also be readily
35
36 374 extended to new defects, systems, and conditions with modest training requirements.
37
38 375 Thus, our approach provides an accurate, efficient, reproducible, scalable, and extensible
39
40 376 method which could replace or greatly enhance human analysis in future studies related
41
42 377 to STEM images.
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5 378 We believe that this framework can be used on many defect and other STEM
6
7
8 379 features simultaneously, eventually providing a general tool for automated analysis
9
10 380 across many STEM applications.

11 381
12
13
14 382 **METHODS**
15
16
17

18 383 *Data Set Collection*

19
20 384 Data set collection was completed as part of a large-scale effort to characterize
21
22 385 iron-chromium-aluminum (FeCrAl) materials neutron-irradiated within the High Flux
23
24 386 Isotope Reactor at Oak Ridge National Laboratory. The dataset comprises a series of
25
26
27 387 published^{29,31,41} and unpublished data. [The data](#) collection was completed over 3 years
28
29
30 388 and spaned a range of different FeCrAl alloys, including model, commercial, and
31
32
33
34 389 engineering-grade alloys irradiated to light water reactor-relevant conditions (e.g., <15
35
36
37 390 displacements per atom and temperatures of nominally 285–320°C). Images generation
38
39
40 391 are described in more details in Li et al⁵.
41
42
43 392
44
45
46
47

48 393 *Data Set Preparation*
49
50

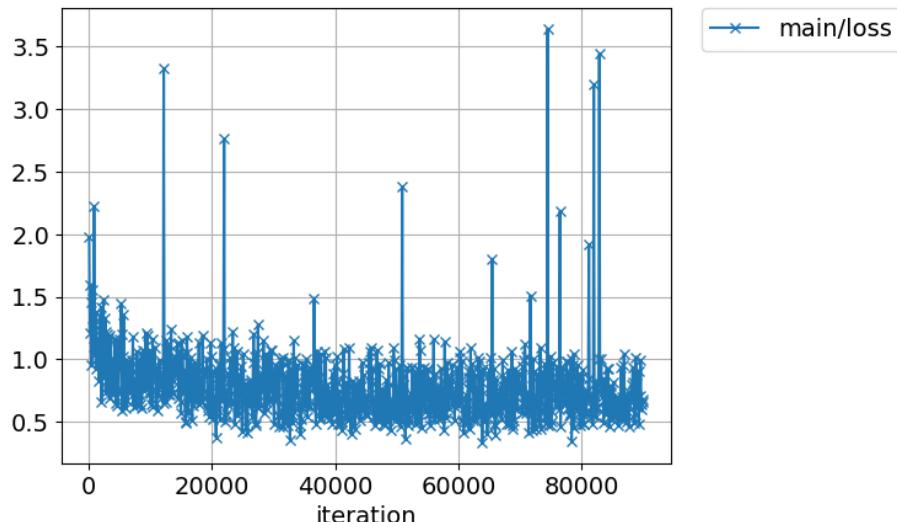
51 394 We used ImageJ^{42,43}, an open-source software for analysis of scientific images, to
52
53 395 manually label all the training and testing data set. And since STEM images are gray
54
55
56 396 scale and ChainerCV⁴⁴ expects input images with RGB channels, some modifications are
57
58
59 397 necessary. We use the direct STEM image gray scale for the R channel. Then we use
60
61
62
63
64

1
2
3
4
5 398 modifications of the original image gray scale for the G and B channels. Specifically,
6
7
8 399 following Li et al.⁵, for G we use a local contrast enhancement of the original gray scale
9
10 400 channel saturated to maximum/minimum and for B we use a Gaussian bluer filter of the
11
12 401 original gray scale STEM images. For the local contrast enhancement in channel G, we
13
14 402 use the Contrast Limited Adaptive Histogram Equalization (CLAHE), a common
15
16 403 algorithm used for local contrast enhancement that makes local detail of STEM image
17
18 404 enhanced even in regions that are darker or lighter than most of the image. The Gaussian
19
20 405 filter used in channel B represents cases where there might be noises or blurring in the
21
22 406 STEM images. The parameters used for CLAHE⁴⁵ and Gaussian blur⁴⁶ are all from the
23
24 407 default parameter setting of scikit-images and details can be found in the references given
25
26
27 408 here for these methods. The purpose of adding two more channels in this way is to
28
29
30 409 improve the model performance and make the model more robust by providing more
31
32
33 410 information about various contrast levels or blurring.
34
35
36
37
38
39
40
41 411 For the training and testing on the Faster R-CNN model, a total of 165 STEM
42
43
44 412 images of irradiated ferritic alloys were collected and labeled. The images were taken at
45
46
47 413 different experimental conditions of temperature and irradiation damage level so that the
48
49
50 414 data includes varying defect sizes, shapes, and areal density. We constructed the ground
51
52
53 415 truth labeling by giving each image in the dataset to at least two groups of at least two
54
55
56 416 researchers per group who together labeled each image in that dataset. In some cases, no
57
58
59 417 absolute consensus could be reached on whether a feature was a defect and/or what type

1
2
3
4
5 418 it had, in which case a best effort was made based on group discussion. Details of the
6
7
8 419 protocol are in the SI.
9
10
11 420 The test dataset was randomly selected from the complete image dataset, so that
12
13 421 the training and test were split by approximately 10:1 ratio. The training dataset was then
14
15 422 augmented to 918 images in total, which could provide more training instances without
16
17
18 423 spending more manpower on labeling. The data is augmented by rotating and/or flipping
19
20 424 each image in the training set, a standard method previously well established to improved
21
22 425 results in some cases⁴⁷.
23
24
25
26
27
28 426
29
30
31 427 ***Model Training***
32
33
34 428 The Faster R-CNN model used VGG-16 as its backbone architecture and we
35
36
37 429 adopted the module provided by ChainerCV⁴⁴ as the Module A in Figure 2 and using
38
39
40 430 watershed function provided by OpenCV⁴⁸ as the second module. The initial weights of
41
42
43 431 Faster R-CNN was loaded from the pre-trained weights from ImageNet which is a
44
45
46 432 common practice in the deep learning training strategy⁴⁰ called transfer learning.
47
48
49 433 Although ImageNet is trained for image classification, not object detection, there are
50
51
52 434 enough similarities in key features to support effective transfer learning of weights.
53
54
55 435 Transfer learning can reduce the amount of data and training time required for good
56
57
58 436 performance⁴⁰. The Faster R-CNN module was optimized with Stochastic Gradient
59
60 437 Descent (SGD) on a single Nvidia GeForce GTX 1080 GPU. The best hyper parameter
61
62
63
64
65

1
2
3
4
5 438 set was found by performing [hyperparameter](#) search of learning rate from 10^{-3} to 10^{-6} and
6
7
8 439 we adjust the needed iteration numbers correspondingly. The best choice of hyper
9
10
11 440 parameter is a decayed learning rate starting from 10^{-4} and each 20000 iterations the
12
13
14 441 learning rate will decay to one tenth of the previous one. In total 90000 iterations were
15
16
17 442 performed, and a learning loss curve is shown in Figure 5. The geometry extraction
18
19
20 443 module needed no training.

21
22 444



1
2
3
4
5 452 used to measure the performance of object detection models¹⁷. IoU is calculated from the
6
7 453 ratio of overlap area of a ground truth bounding box and a predicted bounding box to the
8
9 454 area of union of two bounding boxes. The range of IoU is from 0 to 1 where 0 means no
10
11 455 overlap found between two bounding boxes and 1 means the two bounding boxes are
12
13 456 perfectly overlapping. The threshold IoU is the value used to judge the prediction quality
14
15 457 of the overlapping of ground truth bounding boxes and prediction bounding boxes. A
16
17 458 higher threshold IoU requires more accurate location prediction of the bounding box
18
19 459 detector, which will generally reduce performance, but lower the threshold IoU could
20
21 460 lead a predicted bounding box to being assigned to no defect or the wrong defect. And
22
23 461 another important hyperparameter is the threshold confidence score, a value from 0.0 to
24
25 462 1.0 used by Faster R-CNN internally to discard low confidence proposals in the RPN, and
26
27 463 it can change the total number of outputs of Faster R-CNN. We used grid search of the
28
29 464 threshold IoU and confidence score to search the best choice of these two values [based on](#)
30
31 465 [maximizing the F1 scores](#), with confidence score from the list [0.001, 0.005, 0.01, 0.05,
32
33 466 0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5, 0.55, 0.6] and the threshold IoU from the list
34
35 467 [0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9]. We selected 0.25 as the confidence score for
36
37 468 Faster R-CNN and [showed](#) the performance changes with 0.4 threshold IoU in Figure 4.
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57 470 ***Geometry Fitting of Analysis Module***

1
2
3
4
5 471 After the Faster R-CNN module was performed on specific image, the analysis
6
7 472 module was called to obtain shape and size of the defect contained in bounding box. As
8
9 473 shown in the third column in [Figure 3](#), the approach fits the defect with elliptical contours
10
11
12 474 to estimate their actual shapes and diameters. The approach uses the watershed algorithm
13
14
15 475 to identify the pixels that make up the defect contour and then fit those to an ellipse. The
16
17
18 476 watershed algorithm is a widely used technique for image segmentation purposes that
19
20
21 477 views any gray scale image as a topographic surface where the high (e.g. white) pixel
22
23
24 478 values represents peaks while the low (e.g. black) pixel values denotes valleys. The
25
26
27 479 algorithm tries to grow the region areas by flooding the valleys and where different
28
29
30 480 regions meet with each other are the watershed lines needed for image segmentation⁴⁹.
31
32
33
34 481 [Watershed methods were applied to find the boundary between defect pixels and](#)
35
36
37 482 [background pixels. We followed the official tutorial from OpenCV for performing the](#)
38
39
40 483 [watershed and details of the approach can be found there⁵⁰.](#) [We then fit the boundaries](#)
41
42
43 484 [found from the Watershed algorithm to an ellipse. This fitting was done to match the](#)
44
45
46 485 [approach used by the radiation defect analysis community, obtain a well-defined shape](#)
47
48
49 486 [with simple geometric descriptors, and smooth out the otherwise rather rough boundaries](#)
50
51
52 487 [found by the Watershed algorithm. The fitting was done with OpenCV's](#)
53
54
55 488 [fitEllipse \(\) function⁵¹.](#) All codes were based with OpenCV⁴⁸ and by applying the
56
57 489 second module we could get precise information about the defects' position, size, and
58
59
60 490 orientations. The diameters and areas of defects are defined as follows, where a and b are
61
62
63
64
65

1
2
3
4
5 491 half the lengths of major and minor axes of the ellipse. The diameter of the $a/2<111>$ and
6
7 492 $a<100>$ defects are defined as $2a$. The diameter of the black dot is defined as twice the
8
9
10 493 square root of (ab) . The area of all defects is defined as πab . The areal density is the sum
11
12 494 of defect areas in a set of images divided by the total area of the set of images.
13
14 495

15
16 496 ***Data Availability***
17
18
19
20
21 497 We used a subset of published STEM images of the irradiated FeCrAl alloy
22
23
24
25 498 system⁵ (<https://publish.globus.org/jspui/handle/ITEM/997>) which were labeled and used
26
27
28 499 in this study. The data are available at Figshare
29
30
31 500 (<https://doi.org/10.6084/m9.figshare.8266484>) and the source code for the model is
32
33
34 501 available on Github (<https://github.com/uw-cmg/multitype-defect-detection>). The
35
36
37 502 dataset includes both the images and bounding boxes we used for this project. Data on
38
39
40 503 Figshare also includes a CSV file with all data used in plots in this paper.
41
42
43 504

44
45 505 ***Supporting Information***
46
47
48 506 We showed the fitting results of all 12 testing images in section 1 of SI. In Section
49
50
51 507 2 of SI, we presented the labeling process that has been used in a previous study⁵ and
52
53
54 508 prepared a detailed instruction document to record our labeling process, which can be
55
56
57 509 easily used for other defect images. In section 3 of SI, detailed statistics distribution of
58
59
60
61
62
63
64
65

1
2
3
4
5 510 the human labelling and machine predicting results of diameters and areal density were
6
7
8 511 showed.
9
10
11 512

12
13 513 **Competing Interests:**

14
15 514 There are no competing interests in relation to the work described.
16
17
18
19 515

20
21 516 **Acknowledgments:**
22
23
24

25 517 We would like to thank Wisconsin Applied Computing Center (WACC) for
26
27 518 providing access to CPU/GPU cluster, Euler. And special thanks to Colin Vanden Heuvel
28
29 519 for helping us use GPUs and install software needed and sincere thanks to Vanessa
30
31 520 Meschke for helping organize the undergraduate students participate in this research in
32
33 521 summer 2018.
34
35
36
37 522

38
39 523 **Funding¹:**
40
41
42

43 524 Research was sponsored by the Department of Energy (DOE) Office of Nuclear
44
45 525 Energy, Advanced Fuel Campaign of the Nuclear Technology Research and Development
46
47
48
49
50
51
52

53 ¹ Notice: This manuscript has been authored by UT-Battelle, LLC, under contract DE-AC05-00OR22725 with the US Department of
54 Energy (DOE). The US government retains and the publisher, by accepting the article for publication, acknowledges that the US
55 government retains a nonexclusive, paid-up, irrevocable, worldwide license to publish or reproduce the published form of this
56 manuscript, or allow others to do so, for US government purposes. DOE will provide public access to these results of federally
57 sponsored research in accordance with the DOE Public Access Plan (<http://energy.gov/downloads/doe-public-access-plan>).
58
59
60
61
62
63
64
65

1
2
3
4
5 526 program (formerly the Fuel Cycle R&D program). Neutron irradiation of FeCrAl alloys at
6
7
8 527 Oak Ridge National Laboratory's High Flux Isotope Reactor user facility was sponsored
9
10 528 by the Scientific User Facilities Division, Office of Basic Energy Sciences, DOE. Support
11
12 529 for D. M. was provided by the National Science Foundation Cyberinfrastructure for
13
14 530 Sustained Scientific Innovation (CSSI) program, award No. 1931298. Support for select
15
16 531 undergraduate participants over some periods provided by the NSF University of
17
18 532 Wisconsin-Madison Materials Research Science and Engineering Center (DMR 1720415)
19
20 533 and the Schmidt Foundation.
21
22
23
24
25
26
27
28
29 534

1
2
3
4
5 535 **Reference**
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 536 1. Jenkins, M. . & Kirk, M. . *Characterisation of Radiation Damage by Transmission Electron Microscopy. Iop* **20002352**, (Taylor & Francis, 2000).
- 538 2. Jesse, S. *et al.* Big Data Analytics for Scanning Transmission Electron Microscopy Ptychography. *Sci. Rep.* **6**, 26348 (2016).
- 540 3. Kalinin, S. V., Sumpter, B. G. & Archibald, R. K. Big-deep-smart data in imaging for guiding materials design. *Nat. Mater.* **14**, 973–980 (2015).
- 542 4. Duval, L. *et al.* Image processing for materials characterization: Issues, challenges and opportunities. in *2014 IEEE International Conference on Image Processing, ICIP 2014* 4862–4866 (IEEE, 2014). doi:10.1109/ICIP.2014.7025985
- 544 5. Li, W., Field, K. G. & Morgan, D. Automated defect analysis in electron microscopic images. *npj Comput. Mater.* **4**, 1–9 (2018).
- 546 6. Park, C. & Ding, Y. Automating material image analysis for material discovery. *MRS Commun.* **9**, 545–555 (2019).
- 548 7. Groom, D. J. *et al.* Automatic segmentation of inorganic nanoparticles in BF TEM micrographs. *Ultramicroscopy* **194**, 25–34 (2018).
- 550 8. DeCost, B. L., Francis, T. & Holm, E. A. Exploring the microstructure manifold: Image texture representations applied to ultrahigh carbon steel microstructures. *Acta Mater.* **133**, 30–40 (2017).
- 552 9. Decost, B. L. & Holm, E. A. A computer vision approach for automated analysis

- 1
2
3
4
5 555 and classification of microstructural image data. *Comput. Mater. Sci.* **110**, 126–
6
7
8 556 133 (2015).
9
10 557 10. DeCost, B. L., Jain, H., Rollett, A. D. & Holm, E. A. Computer Vision and
11
12
13 558 Machine Learning for Autonomous Characterization of AM Powder Feedstocks.
14
15
16 559 *Jom* **69**, 456–465 (2017).
17
18
19 560 11. DeCost, B. L. & Holm, E. A. Characterizing powder materials using keypoint-
20
21
22 561 based computer vision methods. *Comput. Mater. Sci.* **126**, 438–445 (2017).
23
24
25 562 12. Marsh, B. P., Chada, N., Sanganna Gari, R. R., Sigdel, K. P. & King, G. M. The
26
27
28 563 Hessian Blob Algorithm: Precise Particle Detection in Atomic Force Microscopy
29
30
31 564 Imagery. *Sci. Rep.* **8**, 978 (2018).
32
33
34 565 13. Chowdhury, A., Kautz, E., Yener, B. & Lewis, D. Image driven machine learning
35
36
37 566 methods for microstructure recognition. *Comput. Mater. Sci.* **123**, 176–187 (2016).
38
39
40 567 14. Vlcek, L., Maksov, A., Pan, M., Vasudevan, R. K. & Kalinin, S. V. Knowledge
41
42
43 568 Extraction from Atomically Resolved Images. *ACS Nano* **11**, 10313–10320 (2017).
44
45
46 569 15. DeCost, B. L. & Holm, E. A. A large dataset of synthetic SEM images of powder
47
48
49 570 materials and their ground truth 3D structures. *Data Br.* **9**, 727–731 (2016).
50
51
52 571 16. Gola, J. *et al.* Advanced microstructure classification by data mining methods.
53
54
55 572 *Comput. Mater. Sci.* **148**, 324–335 (2018).
56
57
58 573 17. Liu, L. *et al.* Deep Learning for Generic Object Detection: A Survey. (2018).
59
60
61 574 18. Zou, Z., Shi, Z., Guo, Y. & Ye, J. Object Detection in 20 Years: A Survey. 1–39
62
63
64
65

- 1
2
3
4
5 575 (2019).
6
7
8 576 19. Roberts, G. *et al.* Deep Learning for Semantic Segmentation of Defects in
9 Advanced STEM Images of Steels. *Sci. Rep.* **9**, (2019).
10
11
12
13 578 20. Anderson, C. M., Klein, J., Rajakumar, H., Judge, C. D. & B, L. K. Automated
14
15
16 579 Classification of Helium Ingress in Irradiated X-750. 1–7 (2019).
17
18
19 580 21. Rusanovsky, M. *et al.* Anomaly Detection using Novel Data Mining and Deep
20
21
22 581 Learning Approach.
23
24
25 582 22. Ahonen, T., Hadid, A. & Pietikäinen, M. Face description with local binary
26
27
28 583 patterns: Application to face recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*
29
30
31 584 **28**, 2037–2041 (2006).
32
33
34 585 23. Viola, P. & Jones, M. Rapid object detection using a boosted cascade of simple
35
36
37 586 features. in *Proceedings of the 2001 IEEE Computer Society Conference on*
38
39
40 587 *Computer Vision and Pattern Recognition. CVPR 2001* **1**, I-511-I-518 (IEEE
41
42
43 588 Comput. Soc).
44
45 589 24. Lecun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444 (2015).
46
47
48 590 25. Badrinarayanan, V., Kendall, A. & Cipolla, R. SegNet: A Deep Convolutional
49
50
51 591 Encoder-Decoder Architecture for Image Segmentation. *IEEE Trans. Pattern*
52
53
54 592 *Anal. Mach. Intell.* **39**, 2481–2495 (2017).
55
56
57 593 26. Ziatdinov, M. *et al.* Deep Learning of Atomically Resolved Scanning
58
59
60 594 Transmission Electron Microscopy Images: Chemical Identification and Tracking
61
62
63
64
65

- 1
2
3
4
5 595 Local Transformations. *ACS Nano* **11**, 12742–12752 (2017).
6
7
8 596 27. Zafari, S., Eerola, T., Ferreira, P., Kälviäinen, H. & Bovik, A. Automated
9
10 597 Segmentation of Nanoparticles in BF TEM Images by U-Net Binarization and
11
12
13 598 Branch and Bound. in *Lecture Notes in Computer Science (including subseries*
14
15
16 599 *Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* **11678**
17
18
19 600 **LNCS**, 113–125 (Springer Verlag, 2019).
20
21
22 601 28. Chen, D., Zhang, P., Liu, S., Chen, Y. & Zhao, W. Aluminum alloy
23
24 602 microstructural segmentation in micrograph with hierarchical parameter transfer
25
26
27 603 learning method. *J. Electron. Imaging* **28**, 1 (2019).
28
29
30
31 604 29. Field, K. G., Briggs, S. A., Sridharan, K., Yamamoto, Y. & Howard, R. H.
32
33
34 605 Dislocation loop formation in model FeCrAl alloys after neutron irradiation below
35
36
37 606 1 dpa. *J. Nucl. Mater.* **495**, 20–26 (2017).
38
39
40 607 30. Parish, C. M., Field, K. G., Certain, A. G. & Wharry, J. P. Application of STEM
41
42 608 characterization for investigating radiation effects in BCC Fe-based alloys. *J.*
43
44
45 609 *Mater. Res.* **30**, 1275–1289 (2015).
46
47
48 610 31. Field, K. G., Hu, X., Littrell, K. C., Yamamoto, Y. & Snead, L. L. Radiation
49
50
51 611 tolerance of neutron-irradiated model Fe-Cr-Al alloys. *J. Nucl. Mater.* **465**, 746–
52
53
54 612 755 (2015).
55
56
57 613 32. Zinkle, S. J. & Busby, J. T. Structural materials for fission & fusion energy.
58
59
60 614 *Materials Today* **12**, 12–19 (2009).
61
62
63
64
65

- 1
2
3
4
5 615 33. Yao, B., Edwards, D. J. & Kurtz, R. J. TEM characterization of dislocation loops
6
7
8 616 in irradiated bcc Fe-based steels. *J. Nucl. Mater.* **434**, 402–410 (2013).
9
10
11 617 34. Ren, S., He, K., Girshick, R. & Sun, J. Faster R-CNN: Towards Real-Time Object
12
13
14 618 Detection with Region Proposal Networks. *IEEE Trans. Pattern Anal. Mach.
15
16 Intell.* **39**, 1137–1149 (2015).
17
18
19 620 35. Zhao, Z.-Q., Zheng, P., Xu, S.-T. & Wu, X. Object Detection With Deep Learning:
20
21
22 621 A Review. *IEEE Trans. Neural Networks Learn. Syst.* 1–21 (2019).
23
24
25 622 doi:10.1109/TNNLS.2018.2876865
26
27
28 623 36. Dean, J. *et al.* Large Scale Distributed Deep Networks. in *Advances in Neural
29
30 Information Processing Systems 25* (eds. Pereira, F., Burges, C. J. C., Bottou, L. &
31
32
33 624 Weinberger, K. Q.) 1223–1231 (Curran Associates, Inc., 2012).
34
35
36 625
37 626 37. Jia, Y. *et al.* Caffe. in *Proceedings of the ACM International Conference on
38
39 Multimedia - MM '14* 675–678 (ACM Press, 2014). doi:10.1145/2647868.2654889
40
41
42 627
43 628 38. Cheng, Y., Wang, D., Zhou, P. & Zhang, T. A Survey of Model Compression and
44
45
46 629 Acceleration for Deep Neural Networks. (2017).
47
48 630 39. Chen, P.-H. C. *et al.* An augmented reality microscope with real-time artificial
49
50
51 631 intelligence integration for cancer diagnosis. *Nat. Med.* 1–5 (2019).
52
53
54 632 doi:10.1038/s41591-019-0539-7
55
56
57 633 40. Pan, S. J. & Yang, Q. A Survey on Transfer Learning. *IEEE Trans. Knowl. Data
58
59 Eng.* **22**, 1345–1359 (2010).
60
61
62
63
64
65

- 1
2
3
4
5 635 41. Field, K. G. *et al.* Heterogeneous dislocation loop formation near grain boundaries
6
7
8 636 in a neutron-irradiated commercial FeCrAl alloy. *J. Nucl. Mater.* **483**, 54–61
9
10 637 (2017).
11
12
13
14 638 42. Schindelin, J. *et al.* Fiji: An open-source platform for biological-image analysis.
15
16
17 639 *Nat. Methods* **9**, 676–682 (2012).
18
19
20 640 43. Schneider, C. A., Rasband, W. S. & Eliceiri, K. W. NIH Image to ImageJ: 25 years
21
22 of image analysis. *Nat. Methods* **9**, 671–675 (2012).
23
24
25 642 44. Niitani, Y., Ogawa, T., Saito, S. & Saito, M. ChainerCV: a Library for Deep
26
27 Learning in Computer Vision. 1217–1220 (2017). doi:10.1145/3123266.3129395
28
29
30
31 644 45. Module: exposure — skimage v0.18.0 docs. Available at: https://scikit-image.org/docs/stable/api/skimage.exposure.html#skimage.exposure.equalize_ada
32
33
34 645 pthist. (Accessed: 23rd April 2021)
35
36
37 646
38
39
40 647 46. Module: filters — skimage v0.19.0.dev0 docs. Available at: <https://scikit-image.org/docs/dev/api/skimage.filters.html#skimage.filters.gaussian>. (Accessed:
41
42
43 648
44
45 649 23rd April 2021)
46
47
48 650 47. Shorten, C. & Khoshgoftaar, T. M. A survey on Image Data Augmentation for
49
50 Deep Learning. *J. Big Data* **6**, 60 (2019).
51
52
53
54 652 48. Garrido, G. & Joshi, P. *OpenCV 3.x with Python By Example Second Edition Make
55
56 the most of OpenCV and Python to build applications for object recognition and
57
58 augmented reality.* (2018).
59
60
61
62
63
64
65

- 1
2
3
4
5 655 49. Kornilov, A. & Safonov, I. An Overview of Watershed Algorithm
6
7
8 656 Implementations in Open Source Libraries. *J. Imaging* **4**, 123 (2018).
9
10 657 50. OpenCV: Image Segmentation with Watershed Algorithm. Available at:
11
12
13 658 https://docs.opencv.org/master/d3/db4/tutorial_py_watershed.html. (Accessed:
14
15
16 659 24th March 2021)
17
18
19 660 51. opencv/fitellipse.py at master · kipr/opencv. Available at:
20
21
22 661 <https://github.com/kipr/opencv/blob/master/samples/python/fitellipse.py>.
23
24
25 662 (Accessed: 21st April 2021)
26
27
28 663
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65



Click here to access/download

Supplementary material for on-line publication only
SupplementaryInformation_v24.docx

CRediT author statement

Mingren Shen: Data Curation, Investigation, Conceptualization, Validation, Visualization, Methodology, Software, Formal analysis, Writing- Original draft, Writing- Reviewing and Editing.

Guanzhao Li: Data Curation, Investigation, Validation.

Dongxia Wu: Data Curation, Investigation, Validation.

Yuhan Liu: Data Curation, Investigation, Validation.

Hima Bharathi Adusumilli: Investigation, Data Curation.

Jacob Greaves: Investigation, Data Curation.

Wei Hao: Investigation, Data Curation.

Nathaniel J. Krakauer: Investigation, Data Curation.

Leah Krudy: Investigation, Data Curation.

Jacob Perez: Investigation, Data Curation.

Varun Sreenivasan: Investigation, Data Curation.

Bryan Sanchez: Investigation, Data Curation.

Oigimer Torres: Investigation, Data Curation.

Wei Li: Conceptualization, Data Curation.

Kevin G. Field: Supervision, Conceptualization, Methodology, Project administration, Funding acquisition, Writing- Reviewing and Editing, Resources.

Dane Morgan: Supervision, Conceptualization, Methodology, Project administration, Funding acquisition, Writing- Reviewing and Editing, Resources.

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: