## Journal Name

ARTICLE TYPE

Cite this: DOI: 00.0000/xxxxxxxxx

## Supplementary Information For: Exploring the Transferability of Machine-Learning Models for Analyzing XRD Data of Shocked Microstructures: From Single Crystal to Polycrystals

Daniel Vizoso,<sup>a</sup> Phillip Tsurkan,<sup>b</sup> Ke Ma,<sup>b</sup> Avinash M. Dongare,<sup>b</sup> and Rémi Dingreville<sup>a‡</sup>

### S1 Machine-learning model architectures

### S1.1 Convolutional Autoencoder Architecture and Training

The convolutional autoencoder used for this work is composed of an encoding model and a symmetric decoding model. The encoding model is composed of three convolutional layers, each with a kernel size of 5. The specific encoder architecture is defined as:  $Conv_{4\times121}^{ReLU} \times Conv_{8\times117}^{ReLU} \times Conv_{16\times113}^{ReLU} \times Flat_{1\times1808} \times Linear_{1808\times L_D=20}$ . The nomenclature "Conv" describes a convolutional layer, "Linear" a dense linear layer, and "Flat" a flattening operation. The subscript represents the dimension of a single input after being fed through each layer, with the first number indicating the number of channels and the second number indicating the number of components in each channel. The superscript indicates the activation function that was used after its respective layer, with the ReLU activation function being used after each convolutional operation in the encoder. The architecture of the decoding portion of the autoencoder is a mirrored version of the architecture of the encoder, with the exception that the ReLU activation function is called after the initial inverse linear operation, and is not called after the final transpose convolutional operation. The network architecture has been selected based on prior work<sup>1,2</sup> which demonstrated that this general architecture was appropriate to perform a regression task on similar problems, using similar data formats.

Training of the autoencoder was performed using the mean squared error (MSE) function as the loss function, where the value of the loss was computed by calculating the mean squared difference between the input XRD profiles and the reconstructed XRD profiles produced by the decoder. Optimization of the autoencoder's parameters was performed using the AdamW optimizer as implemented in PyTorch, with a learning rate of  $10^{-3}$  and a weight decay of  $10^{-5}$ , with parameter optimization being performed over a total of 1000 training epochs.

#### S1.2 MLP Architecture and Training

The MLP is composed of six linear layers that gradually expand the number of dimension up until the fourth layer, followed by two linear layers that contract the number of dimensions down to the number of microstructural descriptors that the model is trained to predict. The specific architecture of the MLP is defined as:  $\text{Linear}_{1\times32}^{\text{ReLU}} \times \text{Linear}_{1\times32}^{\text{ReLU}} \times \text{Linear}_{1\times64}^{\text{ReLU}} \times \text{Linear}_{1\times128}^{\text{ReLU}} \times \text{Drop}_{p=0.25} \times \text{Linear}_{1\times32}^{\text{ReLU}} \times \text{Linear}_{1\times64}^{\text{ReLU}}$ . The nomenclature "Drop" refers to a dropout layer, with the subscript indicating the probability p that determines what fraction of the data at that layer is set to zero. The dropout layer is used as a regularization technique to prevent overfitting during training. As in the architecture of the autoencoder, the ReLU activation function is used after each linear layer except for the last layer.

Training of the MLP was also performed using the MSE function as the loss function, with the loss being computed by comparing the known microstructural state descriptor values to the predicted values produced by the MLP. Optimization of the MLP's parameters was also performed using the AdamW optimizer, with a learning rate of  $10^{-4}$  and a weight decay of  $10^{-5}$ . The number of training epochs for the MLP was adjusted based on the composition of the training dataset, as certain shock-loading orientations or combinations of orientations were more prone to overfitting than others. When trained with the complete set of XRD profiles captured from the  $\langle 111 \rangle$  loading direction single-crystal simulation which contained 51 different timesteps, the MLP was trained for 1 000 epochs. For all other single-crystal loading direction shock simulations as well as for a truncated set from the  $\langle 111 \rangle$  loading direction simulation that only included the same timesteps as the other loading directions, the MLP was trained for 3 000 epochs.

‡ email: rdingre@sandia.gov

<sup>&</sup>lt;sup>a</sup> Center for Integrated Nanotechnologies, Sandia National Laboratories, NM, USA

<sup>&</sup>lt;sup>b</sup> Department of Materials Science and Engineering, Institute of Materials Science, University of Connecticut, Storrs, CT, USA

#### S1.3 Training Datasets and Test-Train Splits

As listed in Table S1, a total of 11 different datasets were used for model training, with each dataset being composed of a single shock-loading simulation or multiple shock-loading simulations that used different loading orientations. For each training dataset, 11 different test-train splits were created in order to quantify how the creation of the test-train splits modified model training and prediction performance. The distribution of data between the test and training sets was kept the same when training the autoencoders and the MLPs. All models that included the  $\langle 111 \rangle$  loading orientation as part of a set of multiple loading orientations used the truncated set (dataset 2), so that the number of timesteps included for each loading orientation would be the same. Test-train splits of 70% training and 30% test were used for model training, with a fixed batch size of 32 samples per batch.

Table S1 List of Datasets Used For Model Training

Set Number	Dataset Used For Model Training
1	$\langle 111 \rangle$ (complete)
2	$\langle 111 \rangle$ (truncated)
3	$\langle 110 \rangle$
4	$\langle 100 \rangle$
5	$\langle 112 \rangle$
6	$\langle 111 \rangle$ , $\langle 110 \rangle$ , $\langle 100 \rangle$ , $\langle 112 \rangle$
7	$\langle 111 \rangle$ , $\langle 110 \rangle$ , $\langle 100 \rangle$
8	$\langle 112 \rangle$ , $\langle 110 \rangle$ , $\langle 100 \rangle$
9	$\langle 111 \rangle$ , $\langle 110 \rangle$ , $\langle 112 \rangle$
10	$\langle 111 \rangle$ , $\langle 100 \rangle$ , $\langle 112 \rangle$
11	Polycrystalline

### S2 Microstructural state descriptor ranges and histograms for the binned regions of the shock-loading simulations

Table S2 Maximum and minimum pressures, total dislocation densities, phase fractions, and temperatures from the binned regions of the five different shock simulations

Microstructural Descriptor	Single Crystal $\langle 111 \rangle$	Single Crystal $\langle 110 \rangle$	Single Crystal $\langle 100  angle$	Single Crystal $\langle 112 \rangle$	Polycrystalline
Pressure (GPa)	-8.366 to 43.866	-6.864 to 52.750	-8.818 to 49.678	-9.565 to 50.314	-9.491 to 50.713
Total Dislocation Density ( $Å^{-2}$ )	0.0 to 12.304	0.0 to 13.066	0.0 to 3.819	0.0 to 8.665	0.326 to 6.083
Disordered Fraction	0.0 to 0.9495	0.0 to 0.7992	0.0 to 0.9683	0.0 to 0.8022	0.0847 to 0.6893
FCC Fraction	0.0381 to 1.0	0.1609 to 1.0	0.0195 to 1.0	0.1354 to 1.0	0.2310 to 0.9137
HCP Fraction	0.0 to 0.1514	0.0 to 0.1776	0.0 to 0.3510	0.0 to 0.2093	0.0010 to 0.30074
Temperature (K)	298.704 to 1425.779	299.419 to 1415.516	299.340 to 1622.331	299.425 to 1295.620	297.905 to 1369.995

Table S2 provides the maximum and minimum values for the six microstructure descriptors measured within the binned regions of the five different shock simulations performed for this work. As discussed in Section 3.1 of the main text, loading orientation and the initial microstructure had significant impacts on the evolution of the system during the shock-loading process for all of the measured microstructure state descriptors. Histograms of the distributions measured in the binned regions for the six microstructure state descriptors are provided in Figs. S1, S2, S3, S4, S5, and S6.



Fig. S1 Histograms of the pressures measured in the binned regions for set numbers: a) 1, b) 2, c) 3, d) 4, e) 5, f) 11.



Fig. S2 Histograms of the total dislocation densities measured in the binned regions for set numbers: a) 1, b) 2, c) 3, d) 4, e) 5, f) 11.



Fig. S3 Histograms of the disordered phase fractions measured in the binned regions for set numbers: a) 1, b) 2, c) 3, d) 4, e) 5, f) 11.



Fig. S4 Histograms of the FCC phase fractions measured in the binned regions for set numbers: a) 1, b) 2, c) 3, d) 4, e) 5, f) 11.



Fig. S5 Histograms of the HCP phase fractions measured in the binned regions for set numbers: a) 1, b) 2, c) 3, d) 4, e) 5, f) 11.



Fig. S6 Histograms of the temperatures measured in the binned regions for set numbers: a) 1, b) 2, c) 3, d) 4, e) 5, f) 11.

### S3 Validation set regression accuracies for models trained with individual single-crystal shock-loading simulations

Tables S3, S4, S5, S6, and S7 provide the validation-set regression accuracy metrics for models trained with individual single-crystal shock-loading simulations. All metrics shown are MSE values computed using the portion of the corresponding dataset that was set aside for model validation and was not used for model training. MSE values were computed with the microstructure state descriptors being normalized to have values between 0 and 1. The "Average MSE" column shown in the tables provides the average MSE and its standard deviation computed for that microstructure state descriptor across the 11 models that were trained with different test-train splits of the corresponding dataset. The "Best Model MSE" column provides the MSE values for each descriptor from the model that had the lowest average MSE for the six microstructure state descriptors. The "Best Overall MSE" column shows the lowest MSE observed for that microstructure state descriptor across the 11 different models.

Examining the metrics provided in Tables S3, S4, S5, S6, and S7, we observe that models trained on individual single-crystal shockloading simulations achieve similar prediction accuracy for the six microstructure state descriptors in their respective validation sets. Generally, the models are most accurate when predicting the pressure and the phase fractions of the disordered or FCC phases, and tend to have lower accuracy when predicting the total dislocation density and HCP phase fraction. The prediction accuracy for the temperature does vary depending on which shock-loading orientation was used to train the models, with the models trained with the  $\langle 100 \rangle$  and  $\langle 112 \rangle$  loading orientation simulations having temperature prediction accuracies similar to their respective phase fraction prediction accuracies, while the other three shock-loading orientations had temperature prediction accuracies which were noticeably worse. Examining the average and standard deviations for the MSE provided in the tables, it can be noted that the creation of the testtrain split does have a significant impact on the performance of the models that are trained with those datasets. Comparisons between the MSE's from the best models to the best overall MSE's, it can be noted that the lowest average MSE would not have the best MSE for every microstructure state descriptor, and at times the model with the lowest average MSE would have individual MSE values for specific microstructure state descriptors that were above average compared to the other models trained with that same data.

Table S3 Valida	tion set MSE 1	metrics for m	odels trained	with the co	omplete (111	shock-loading	simulation
-----------------	----------------	---------------	---------------	-------------	--------------	---------------	------------

Microstructure Descriptor	Average MSE	Best Model MSE	Best Overall MSE
Pressure	$7.16 \pm 5.25 \times 10^{-4}$	$4.37 \times 10^{-4}$	$3.80 \times 10^{-4}$
Total Dislocation Density	$1.79 \pm 0.633 \times 10^{-3}$	$1.14 \times 10^{-3}$	$1.09  imes 10^{-3}$
Disordered Fraction	$1.09 \pm 0.411 \times 10^{-3}$	$5.10 \times 10^{-3}$	$5.10  imes 10^{-4}$
FCC Fraction	$1.32 \pm 0.586 \times 10^{-3}$	$5.77 imes10^{-4}$	$5.77 imes10^{-4}$
HCP Fraction	$4.11 \pm 1.71 \times 10^{-3}$	$2.53 imes10^{-3}$	$2.18 imes10^{-3}$
Temperature	$1.54 \pm 0.433 \times 10^{-3}$	$1.11  imes 10^{-3}$	$8.41  imes 10^{-4}$

Table S4 Validation set MSE metrics for models trained with the truncated  $\langle 111 \rangle$  shock-loading simulation

Microstructure Descriptor	Average MSE	Best Model MSE	Best Overall MSE
Pressure	$8.17 \pm 5.50  imes 10^{-4}$	$2.60 \times 10^{-4}$	$2.60 \times 10^{-4}$
Total Dislocation Density	$1.97 \pm 0.596 \times 10^{-3}$	$1.03  imes 10^{-3}$	$1.03  imes 10^{-3}$
Disordered Fraction	$5.32 \pm 2.99 \times 10^{-4}$	$3.31  imes 10^{-4}$	$2.60  imes 10^{-4}$
FCC Fraction	$6.62 \pm 5.36 \times 10^{-4}$	$3.42 \times 10^{-4}$	$2.33  imes 10^{-4}$
HCP Fraction	$2.24 \pm 1.18 \times 10^{-3}$	$1.79  imes 10^{-3}$	$1.38  imes 10^{-3}$
Temperature	$1.34 \pm 0.594 \times 10^{-3}$	$1.08  imes 10^{-3}$	$7.08 imes10^{-4}$

Table S5 Validation set MSE metrics for models trained with the  $\langle 110 \rangle$  shock-loading simulation

Microstructural Descriptor	Average MSE	Best Model MSE	Best Overall MSE
Pressure	$4.27 \pm 2.49 \times 10^{-4}$	$2.28  imes 10^{-4}$	$2.28  imes 10^{-4}$
Total Dislocation Density	$1.26 \pm 1.05 \times 10^{-3}$	$5.36 imes10^{-4}$	$5.36 imes10^{-4}$
Disordered Fraction	$8.14 \pm 8.35 \times 10^{-4}$	$2.24  imes 10^{-4}$	$2.24  imes 10^{-4}$
FCC Fraction	$5.20 \pm 3.82 \times 10^{-4}$	$2.66  imes 10^{-4}$	$2.02  imes 10^{-4}$
HCP Fraction	$2.14 \pm 1.33 \times 10^{-3}$	$1.09  imes 10^{-3}$	$1.00  imes 10^{-3}$
Temperature	$1.18\pm 0.771\times 10^{-3}$	$1.26  imes 10^{-3}$	$6.76  imes 10^{-4}$

Table S6 Validation set MSE metrics for models trained with the  $\langle 100 \rangle$  shock-loading simulation

Microstructural Descriptor	Average MSE	Best Model MSE	Best Overall MSE
Pressure	$4.51 \pm 1.97 \times 10^{-4}$	$2.94  imes 10^{-4}$	$2.40  imes 10^{-4}$
Total Dislocation Density	$2.80 \pm 0.707 \times 10^{-3}$	$1.78  imes 10^{-3}$	$1.78 imes10^{-3}$
Disordered Fraction	$6.23 \pm 2.54 \times 10^{-4}$	$3.56 \times 10^{-4}$	$3.56 \times 10^{-4}$
FCC Fraction	$6.20 \pm 2.07 \times 10^{-4}$	$3.93  imes 10^{-4}$	$3.93  imes 10^{-4}$
HCP Fraction	$8.08 \pm 2.76 \times 10^{-4}$	$6.73  imes 10^{-4}$	$4.17  imes 10^{-4}$
Temperature	$1.00\pm 0.585\times 10^{-3}$	$6.21  imes 10^{-4}$	$4.33  imes 10^{-4}$

Table S7 Validation set MSE metrics for models trained with the  $\langle 112\rangle$  shock-loading simulation

Microstructural Descriptor	Average MSE	Best Model MSE	Best Overall MSE
Pressure	$3.38 \pm 3.64 \times 10^{-4}$	$1.33 \times 10^{-4}$	$1.11 \times 10^{-4}$
Total Dislocation Density	$1.72\pm 0.706\times 10^{-3}$	$9.53  imes 10^{-4}$	$9.53  imes 10^{-4}$
Disordered Fraction	$6.12\pm 6.08\times 10^{-4}$	$3.84  imes 10^{-4}$	$2.87  imes 10^{-4}$
FCC Fraction	$5.62 \pm 8.06 \times 10^{-4}$	$2.29 imes10^{-4}$	$2.16 imes10^{-4}$
HCP Fraction	$1.21 \pm 0.794 \times 10^{-3}$	$6.48  imes 10^{-4}$	$5.37  imes 10^{-4}$
Temperature	$8.61 \pm 6.63 \times 10^{-4}$	$4.70  imes 10^{-4}$	$3.55  imes 10^{-4}$

## S4 Best overall MSE for models trained with individual single-crystal shock-loading orientations and making predictions for different loading orientations

Tables S8, S9, S10, and S11 summarize the transferability of models trained with individual shock-loading orientation simulations by providing the best overall MSE observed when predicting the microstructure state descriptors for shock-loading orientations that were excluded from model training. As discussed in Section 3.3 of the main text, the ability of models trained with data from one shock-loading orientation to make predictions with data from a different shock-loading orientation depends on several factors. When the expected bounds for the microstructure state descriptors are significantly different between the two datasets, the model will tend to under or over-predict the descriptor depending on whether the training dataset had lower or higher bounds than the targeted dataset. This is most prominently observed for models predicting the descriptors from the  $\langle 100 \rangle$  loading orientation dataset or for the predictions made by the models trained with the  $\langle 100 \rangle$  loading orientation dataset, as the best overall MSEs observed for these tasks tend to be higher than the other MSE's reported in the tables.

Table S8 Best overall MSE for the models trained with the truncated  $\langle 111 \rangle$  shock-loading simulation predicting descriptors for other loading orientation simulations

Microstructure	$\langle 110 \rangle$	$\langle 100 \rangle$	$\langle 112 \rangle$
Descriptor	Best Overall MSE	Best Overall MSE	Best Overall MSE
Pressure	$4.34 \times 10^{-3}$	$1.50 \times 10^{-2}$	$4.42 \times 10^{-3}$
Total Dislocation Density	$6.43 \times 10^{-3}$	$1.01  imes 10^{-2}$	$4.17 \times 10^{-3}$
Disordered Fraction	$5.17 \times 10^{-3}$	$8.57 \times 10^{-3}$	$3.53 \times 10^{-3}$
FCC Fraction	$3.52 \times 10^{-3}$	$1.11  imes 10^{-2}$	$4.25 \times 10^{-3}$
HCP Fraction	$4.36 \times 10^{-3}$	$1.27  imes 10^{-2}$	$5.26  imes 10^{-3}$
Temperature	$2.28  imes 10^{-3}$	$8.50  imes 10^{-3}$	$4.07 \times 10^{-3}$

Table S9 Best overall MSE for the models trained with  $\langle 110 \rangle$  shock-loading simulation predicting descriptors for other loading orientation simulations

Microstructure	$\langle 111 \rangle$ Best Overall MSF	$\langle 100 \rangle$ Best Overall MSF	$\langle 112 \rangle$ Best Overall MSF
Descriptor	Best Overall MDE		
Pressure	$9.81 \times 10^{-5}$	$3.12 \times 10^{-5}$	$2.87 \times 10^{-3}$
Total Dislocation Density	$1.08  imes 10^{-2}$	$3.45  imes 10^{-2}$	$1.84  imes 10^{-2}$
<b>Disordered Fraction</b>	$1.97  imes 10^{-3}$	$2.11  imes 10^{-2}$	$5.32  imes 10^{-3}$
FCC Fraction	$3.53  imes 10^{-3}$	$1.05 imes10^{-2}$	$6.25  imes 10^{-3}$
HCP Fraction	$2.08  imes 10^{-2}$	$1.26  imes 10^{-1}$	$3.61  imes 10^{-2}$
Temperature	$1.55  imes 10^{-2}$	$1.13  imes 10^{-2}$	$8.09 \times 10^{-3}$

Table S10 Best overall MSE for the models trained with (100) shock-loading simulation predicting descriptors for other loading orientation simulations

Microstructure Descriptor	$\langle 111 \rangle$ Best Overall MSE	$\langle 110  angle$ Best Overall MSE	$\langle 112 \rangle$ Best Overall MSE
Pressure	$1.13 \times 10^{-2}$	$2.72 \times 10^{-3}$	$2.92 \times 10^{-3}$
Total Dislocation Density	$3.51  imes 10^{-1}$	$5.55  imes 10^{-1}$	$3.91  imes 10^{-1}$
Disordered Fraction	$5.38  imes 10^{-3}$	$1.60  imes 10^{-2}$	$3.57  imes 10^{-3}$
FCC Fraction	$4.51 \times 10^{-3}$	$1.32  imes 10^{-2}$	$6.55  imes 10^{-3}$
HCP Fraction	$2.37  imes 10^{-2}$	$2.32  imes 10^{-2}$	$3.67 \times 10^{-2}$
Temperature	$1.88  imes 10^{-2}$	$2.23  imes 10^{-2}$	$2.02  imes 10^{-2}$

Table S11 Best overall MSE for the models trained with  $\langle 112 \rangle$  shock-loading simulation predicting descriptors for other loading orientation simulations

Microstructure Descriptor	$\langle 111 \rangle$ Best Overall MSE	$\langle 110 \rangle$ Best Overall MSE	$\langle 100 \rangle$ Best Overall MSE
Pressure	$1.28  imes 10^{-2}$	$5.99 \times 10^{-3}$	$5.06 \times 10^{-3}$
Total Dislocation Density	$8.83  imes 10^{-3}$	$3.10 imes10^{-2}$	$5.74  imes 10^{-2}$
Disordered Fraction	$4.11 \times 10^{-3}$	$6.44 \times 10^{-3}$	$1.31 \times 10^{-2}$
FCC Fraction	$2.33  imes 10^{-3}$	$3.26 \times 10^{-3}$	$6.85 \times 10^{-3}$
HCP Fraction	$1.22  imes 10^{-2}$	$1.87 imes10^{-2}$	$3.81  imes 10^{-2}$
Temperature	$1.83  imes 10^{-2}$	$3.12  imes 10^{-2}$	$6.23  imes 10^{-2}$

# S5 Best overall MSE for models trained with three single-crystal shock-loading orientations and making predictions for the excluded single crystal simulation

Table S12 provides the best overall MSE's observed for models that were trained with three single crystal shock-loading orientation simulations when predicting the microstructure state descriptors with data from the single crystal simulation that was excluded from model training. In general, metrics reported in Table S12 are lower than those reported in the tables from Supplementary Information S4, indicating that using multiple shock-loading orientation simulations to train the models improves their transferability to unseen shock-loading orientations.

Table S12 Best overall MSE observed for models trained with three single crystal shock-loading orientations and predicting on the orientation that was excluded from training

Microstructure	$\langle 111 \rangle$	$\langle 110 \rangle$	$\langle 100 \rangle$	$\langle 112 \rangle$
Descriptor	Best Overall MSE	Best Overall MSE	Best Overall MSE	Best Overall MSE
Pressure	$4.37 \times 10^{-3}$	$1.52 \times 10^{-3}$	$3.06 \times 10^{-3}$	$7.88 imes10^{-4}$
Total Dislocation Density	$3.27 \times 10^{-3}$	$1.62 \times 10^{-2}$	$2.71  imes 10^{-2}$	$3.23 \times 10^{-3}$
Disordered Fraction	$1.72 \times 10^{-3}$	$3.36 \times 10^{-3}$	$1.33  imes 10^{-2}$	$2.79 \times 10^{-3}$
FCC Fraction	$1.37 \times 10^{-3}$	$1.74 \times 10^{-3}$	$4.74 \times 10^{-3}$	$2.16  imes 10^{-3}$
HCP Fraction	$4.40 \times 10^{-3}$	$5.46 \times 10^{-3}$	$3.37  imes 10^{-2}$	$3.55 \times 10^{-3}$
Temperature	$1.09  imes 10^{-2}$	$6.38  imes 10^{-3}$	$1.97  imes 10^{-2}$	$8.23  imes 10^{-3}$

### S6 Best overall MSE for models trained with data from single crystal shock-loading simulations and predicting descriptors for the polycrystalline simulation

Tables S13 and S14 summarize the transferability of models trained with data from single crystal shock-loading simulations for predicting microstructure state descriptors from polycrystalline shock simulations. The dataset numbers shown in the tables indicate the dataset used to train the models, with the numbers referring to the datasets defined in Table S1. The metrics shown for dataset 11 in Table S13 are from models trained with the polycrystalline data making predictions with the validation portion of the polycrystalline dataset, which serve as best-possible performances to compare the models trained with single crystal data against.

As discussed in Section 3.4 of the main text, some microstructure state descriptors such as temperature and total dislocation density aren't predicted well by any of the of the models trained with single crystal data. Others, such as pressure and phase fractions, can be predicted with accuracies comparable to those observed in the other model transfer cases shown in Tables S8, S9, S10, S11, and S12.

Table S13 Best overall MSE observed for models trained with individual single crystal shock-loading simulations and predicting on the polycrystalline shock simulation

Microstructure Descriptor	Dataset 11 (validation) Best Overall MSE	Dataset 2 Best Overall MSE	Dataset 3 Best Overall MSE	Dataset 4 Best Overall MSE	Dataset 5 Best Overall MSE
Pressure	$3.94 \times 10^{-4}$	$3.52 \times 10^{-3}$	$2.89 \times 10^{-3}$	$1.44 \times 10^{-3}$	$4.44 \times 10^{-3}$
Total Dislocation Density	$3.00 \times 10^{-3}$	$6.47  imes 10^{-2}$	$5.97 imes10^{-2}$	$4.90  imes 10^{-2}$	$5.66 imes10^{-2}$
Disordered Fraction	$1.23  imes 10^{-3}$	$2.70 \times 10^{-3}$	$5.02 \times 10^{-3}$	$1.04 \times 10^{-2}$	$2.78  imes 10^{-3}$
FCC Fraction	$6.78  imes 10^{-4}$	$3.08 \times 10^{-3}$	$3.38  imes 10^{-3}$	$4.24 \times 10^{-3}$	$2.30  imes 10^{-3}$
HCP Fraction	$1.36  imes 10^{-3}$	$4.75  imes 10^{-2}$	$2.61  imes 10^{-2}$	$2.38 imes10^{-2}$	$2.05 imes10^{-2}$
Temperature	$1.71  imes 10^{-3}$	$3.65  imes 10^{-2}$	$1.78  imes 10^{-2}$	$8.99  imes 10^{-3}$	$3.49  imes 10^{-2}$

Table S14 Best overall MSE observed for models trained with multiple single crystal shock-loading simulations and predicting on the polycrystalline shock simulation

Microstructure Descriptor	Dataset 6 Best Overall MSE	Dataset 7 Best Overall MSE	Dataset 8 Best Overall MSE	Dataset 9 Best Overall MSE	Dataset 10 Best Overall MSE
Pressure	$1.31 \times 10^{-3}$	$1.83 \times 10^{-3}$	$1.90 \times 10^{-3}$	$2.24 \times 10^{-3}$	$1.62 \times 10^{-3}$
Total Dislocation Density	$2.90  imes 10^{-2}$	$3.21 \times 10^{-2}$	$3.62 \times 10^{-2}$	$4.44 imes10^{-2}$	$2.17 imes10^{-2}$
Disordered Fraction	$2.48  imes 10^{-3}$	$2.10  imes 10^{-3}$	$2.75  imes 10^{-3}$	$2.82  imes 10^{-3}$	$2.75  imes 10^{-3}$
FCC Fraction	$1.36 \times 10^{-3}$	$1.08 \times 10^{-3}$	$1.56 \times 10^{-3}$	$2.06  imes 10^{-3}$	$1.00 \times 10^{-3}$
HCP Fraction	$4.12 \times 10^{-3}$	$4.40 \times 10^{-3}$	$4.24 \times 10^{-3}$	$1.40  imes 10^{-2}$	$9.28  imes 10^{-3}$
Temperature	$1.10 \times 10^{-2}$	$1.03 \times 10^{-2}$	$1.12\times10^{-2}$	$1.31 \times 10^{-2}$	$1.04 \times 10^{-2}$

### S7 Comparisons with a Simpler Machine-Learning Approach

As a point of comparison to the machine-learning model architectures described in Supplementary Information 1, we present here the model training results for a simpler machine-learning workflow. In this simpler approach, the normalized XRD profiles are passed directly into a decision tree regression model to predict the microstructural state descriptors. For this series of tests, scikit-learn's decision tree model GradientBoostingRegressor<sup>3,4</sup> was used for the regression of the microstructural descriptors. As this implementation only allows for the regression of a single regression target, individual single-output regression models were trained for each of the six microstructural descriptors.

For this comparison, two training datasets were used: Dataset 2 (the truncated  $\langle 111 \rangle$  dataset) and Dataset 7 (the combination of the truncated  $\langle 111 \rangle$ ,  $\langle 110 \rangle$ , and  $\langle 100 \rangle$  datasets). For each of the training datasets, ten different test-train splits of 70% training and 30% test were used to train ten different iterations of the direct regression workflow. For the models trained with Dataset 2, the trained models were used to predict the microstructural state descriptors from Datasets 3, 4, 5, and 11. The models trained with Dataset 7 were used to predict the microstructural state descriptors from Datasets 5 and 11.

#### S7.1 Direct Regression from XRD Profiles

For the direct regression workflow, normalized XRD profiles were passed directly to the gradient boosting decision tree regression model for the prediction of the microstructural state descriptors. Table S15 provides validation set regression accuracy metrics for the models trained with Dataset 2. Comparing the metrics provided in Table S15 with those provided in Table S4, we can observe that the simpler direct regression workflow does have worse regression accuracy metrics compared to the sequential autoencoder-MLP workflow described in Supplementary Information 1, with substantial differences in most MSE metrics, and slightly worse performance in the prediction of the disordered phase fraction by the most accurate direct regression model.

Table S15 Validation set MSE metrics for models trained with the truncated  $\langle 111\rangle$  shock-loading simulation

Microstructure Descriptor	Average MSE	Best Model MSE	Best Overall MSE
Pressure	$1.26 \pm 0.984 \times 10^{-3}$	$6.88 \times 10^{-4}$	$5.62 \times 10^{-4}$
Total Dislocation Density	$2.91 \pm 1.07 \times 10^{-3}$	$1.14 \times 10^{-3}$	$1.14  imes 10^{-3}$
<b>Disordered Fraction</b>	$7.92 \pm 2.22 \times 10^{-4}$	$5.76  imes 10^{-4}$	$5.37 imes10^{-4}$
FCC Fraction	$7.23 \pm 5.18 \times 10^{-4}$	$4.52  imes 10^{-4}$	$4.00  imes 10^{-4}$
HCP Fraction	$6.72 \pm 2.37 \times 10^{-3}$	$4.16  imes 10^{-3}$	$3.82  imes 10^{-3}$
Temperature	$2.42 \pm 1.38  imes 10^{-3}$	$1.74 imes10^{-3}$	$1.16 \times 10^{-3}$

Table S16 provides the lowest observed MSE across the ten different test-train splits for the direct regression models that were trained with Dataset 2 when predicting the microstructural state descriptors for Datasets 3, 4, 5, and 11. These metrics can be directly

compared to the accuracies provided in Tables S8 and S13. Comparing the ability of the direct regression models to make predictions on different loading orientation datasets compared to the performance of the more complex machine-learning workflow described in Supplementary Information 1, we can observe that for most of the microstructural descriptors, the more complex machine-learning workflow significantly outperforms the direct regression approach. However, for Datasets 3, 4, and 5, the direct regression approach does perform similarly or outperform the autoencoder-MLP workflow in predicting the disordered phase fraction and the FCC phase fraction. Also, for Dataset 11, the direct regression approach does have a lower MSE for the prediction of the total dislocation density, although both the direct regression approach and the autoencoder-MLP approach perform poorly for this descriptor.

Microstructure	Dataset 3	Dataset 4	Dataset 5	Dataset 11
Descriptor	Best Overall MSE	Best Overall MSE	Best Overall MSE	Best Overall MSE
Pressure	$1.17 \times 10^{-2}$	$4.10 \times 10^{-2}$	$1.21 \times 10^{-2}$	$1.47 \times 10^{-2}$
Total Dislocation Density	$2.31  imes 10^{-2}$	$2.54 imes10^{-2}$	$1.03 imes10^{-2}$	$4.32 \times 10^{-2}$
Disordered Fraction	$3.63 \times 10^{-3}$	$5.53  imes 10^{-3}$	$2.80  imes 10^{-3}$	$3.71  imes 10^{-3}$
FCC Fraction	$4.22 \times 10^{-3}$	$6.18 \times 10^{-3}$	$2.67 \times 10^{-3}$	$4.31 \times 10^{-3}$
HCP Fraction	$5.91  imes 10^{-2}$	$2.38  imes 10^{-1}$	$4.32 \times 10^{-2}$	$8.13  imes 10^{-2}$
Temperature	$2.86  imes 10^{-2}$	$3.51  imes 10^{-2}$	$1.90  imes 10^{-2}$	$6.15  imes 10^{-2}$

Table S16 Best overall MSE observed for models trained with Dataset 2 and predicting on Datasets 3, 4, 5, and 11

Table S17 provides validation set regression accuracy metrics for direct regression models trained with Dataset 7, as well as prediction accuracy metrics when these models predict microstructural state descriptors for Datasets 5 and 11. Comparing the metrics provided in Table S17 with those provided in Tables S12 and S14, we can observe that the direct regression workflow has the same performance trade-offs when trained with data from multiple shock-loading orientations when predicting on different single crystal or polycrystalline shock-loading simulations. Predictions of the disordered phase fraction and FCC phase fraction are better for the direct regression workflow compared to the autoencoder-MLP workflow for Dataset 5, while predictions for the other microstructural state descriptors are worse. Predictions of the total dislocation density of Dataset 11 are slightly better but still poor with the direct regression workflow compared to the autoencoder-MLP workflow, while the regression accuracies of the other microstructural state descriptors are better for the autoencoder-MLP workflow.

Table S17 MSE metrics for models trained with Dataset 7

Microstructure Descriptor	Dataset 7 (validation) Average MSE	Dataset 7 (validation) Best Overall MSE	Dataset 5 Best Overall MSE	Dataset 11 Best Overall MSE
Pressure	$9.85 \pm 7.85 \times 10^{-4}$	$4.90 \times 10^{-4}$	$2.16 \times 10^{-3}$	$6.18 \times 10^{-3}$
Total Dislocation Density	$2.62 \pm 0.821 \times 10^{-3}$	$1.30 \times 10^{-3}$	$5.08  imes 10^{-3}$	$2.62  imes 10^{-2}$
Disordered Fraction	$6.39 \pm 1.07 \times 10^{-4}$	$4.69  imes 10^{-4}$	$1.18  imes 10^{-3}$	$2.97  imes 10^{-3}$
FCC Fraction	$6.54 \pm 0.912 \times 10^{-4}$	$4.86 \times 10^{-4}$	$1.19 \times 10^{-3}$	$2.37 \times 10^{-3}$
HCP Fraction	$1.48 \pm 0.317 \times 10^{-3}$	$9.21  imes 10^{-4}$	$7.96  imes 10^{-3}$	$4.59 \times 10^{-3}$
Temperature	$2.35 \pm 0.541 \times 10^{-3}$	$1.68 \times 10^{-3}$	$9.38  imes 10^{-3}$	$2.28  imes 10^{-2}$

In summary, the direct regression workflow is capable of learning the relationships that exist between XRD profiles collected during shock-loading simulations and sets of microstructural state descriptors. Direct regression models achieve relatively good (although worse than the autoencoder-MLP workflow) regression accuracies for validation sets. With respect to the transferability of direct regression models to data from shock-loading orientations that were excluded from the training data, there is a consistent trend where the direct regression approach performs better than the autoencoder-MLP workflow when predicting the disordered phase fractions and FCC phase fractions of unseen single crystal shock-loading simulations, while performing worse (in some cases, such as predicting the pressure, substantially worse) on all of the other tested microstructural state descriptors. When predicting the microstructural state descriptors from the polycrystalline shock simulation, the direct regression approach performs worse than the autoencoder-MLP workflow (however, both perform poorly for this particular descriptor), while performing worse for all of the other descriptors. In general, the direct regression approach performs worse than the autoencoder-MLP workflow defined in Supplementary Information 1, and the cases where the direct regression approach performs better than the autoencoder-MLP workflow typically come at the cost of a substantial decrease in performance in other areas that compromise its utility.

### S8 PCA Analysis of XRD Data from Molecular Dynamics Simulations of the Shock Loading of Cu

Principal Component Analysis (PCA)<sup>5</sup> can be used to estimate how many latent dimensions will be needed to accurately encode a dataset by examining the amount of variance that each additional latent dimension captures. Figure S7 provides the amount of variance captured by each latent dimension up to a latent dimension size of 20, as well as the cumulative variance that is captured as the number of latent dimensions is increased. For this analysis, the PCA implementation from sckikit-learn was used with all of the xrd profiles from Dataset 6.



Fig. S7 Explained variance and cumulative variance as a function of the intrinsic dimensionality,  $L_d$ . Note that  $L_d = 20$  captures over 99.5% of the variance in our dataset.

Per the analysis of the variance captured as a function of latent dimension size, a latent dimension of size 20 captures over 99.5% of the variance present in the XRD data collected from molecular dynamics simulations of the shock loading of single crystal Cu along different loading orientations. With this analysis as well as a sensitivity study that was performed by examining the reconstruction performance of autoencoders with the XRD data as a function of encoded dimension size, a latent dimension of size 20 was selected for our machine-learning workflows.

### S9 Shock Compression versus Spall Failure

Figure S8 compares the distributions of the stress, temperature, and dislocation density along the loading axis of the simulation cell for the four different single crystal configurations during (a) shock compression and (b) spall failure. These comparisons visualize how varying the loading orientation changes the evolution of the microstructural descriptors at equivalent stages of the shock simulation.



Fig. S8 Z stress, temperature, and dislocation density versus position along the shock-loading direction for different shock-loading orientation simulations during (a) shock compression and (b) spall failure.

## Supplementary References

- 1 D. Vizoso, G. Subhash, K. Rajan and R. Dingreville, Chemistry of Materials, 2013, 35, 1186–1200.
- 2 D. Vizoso and R. Dingreville, Journal of Applied Physics, 2025, 137, XX.
- 3 J. H. Friedman, The Annals of Statistics, 2001, 29, 1189–1232.
- 4 J. H. Friedman, Computational Statistics & Data Analysis, 2002, 38, 367–378.
- 5 A. Mead, Journal of the Royal Statistical Society, 1992, 41, 27–39.