# *Supervised and Unsupervised Machine Learning for Large, Noisy STEM Data*
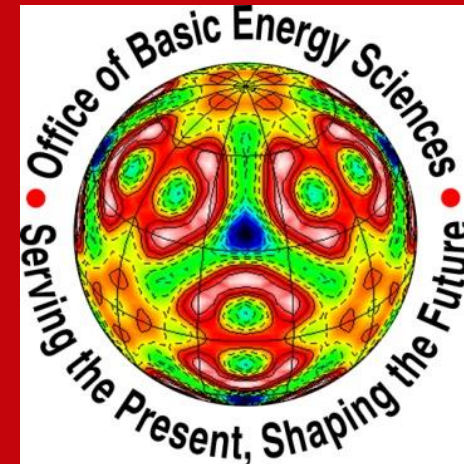
Paul M. Voyles
Materials Science and Engineering

WISCONSIN
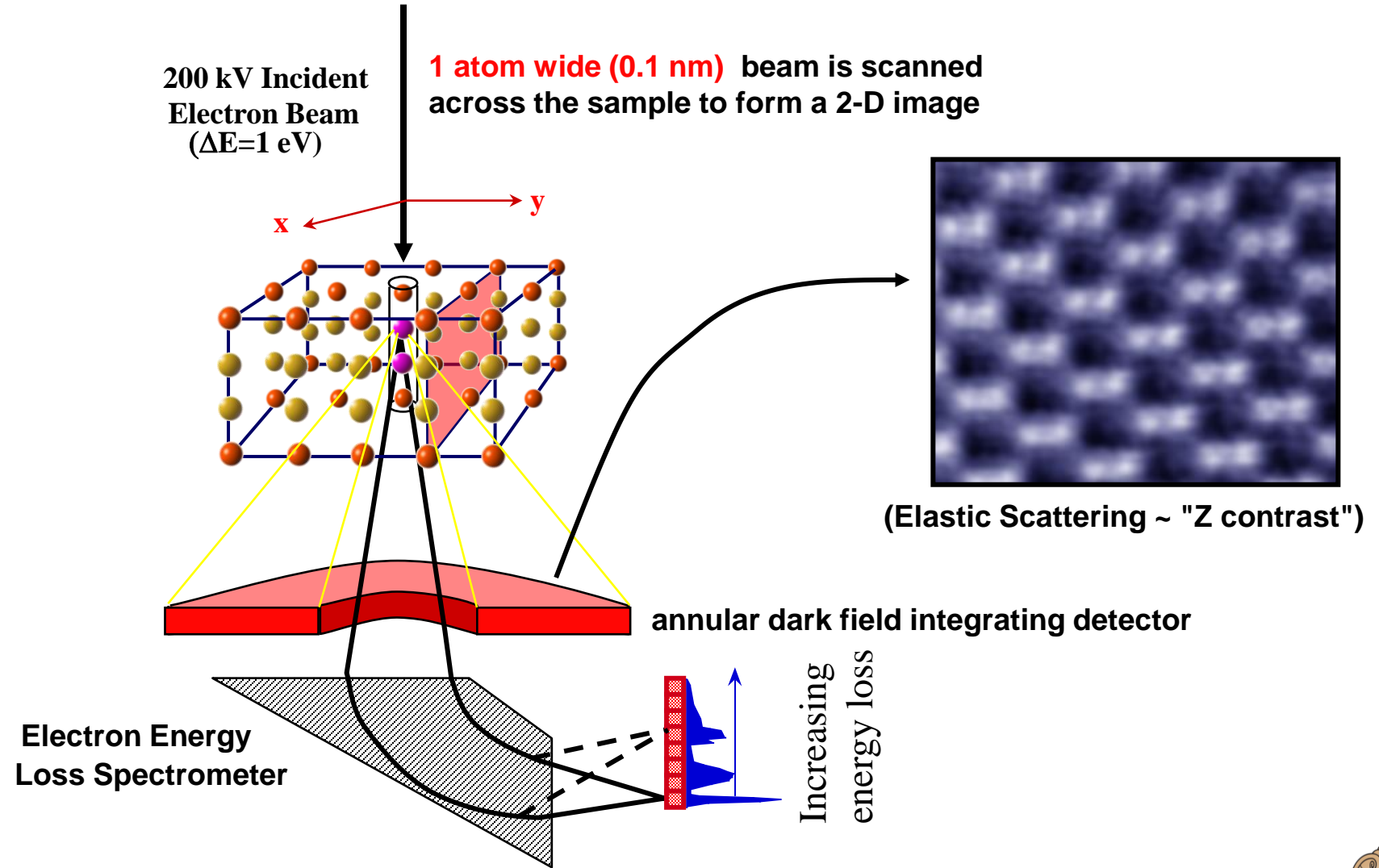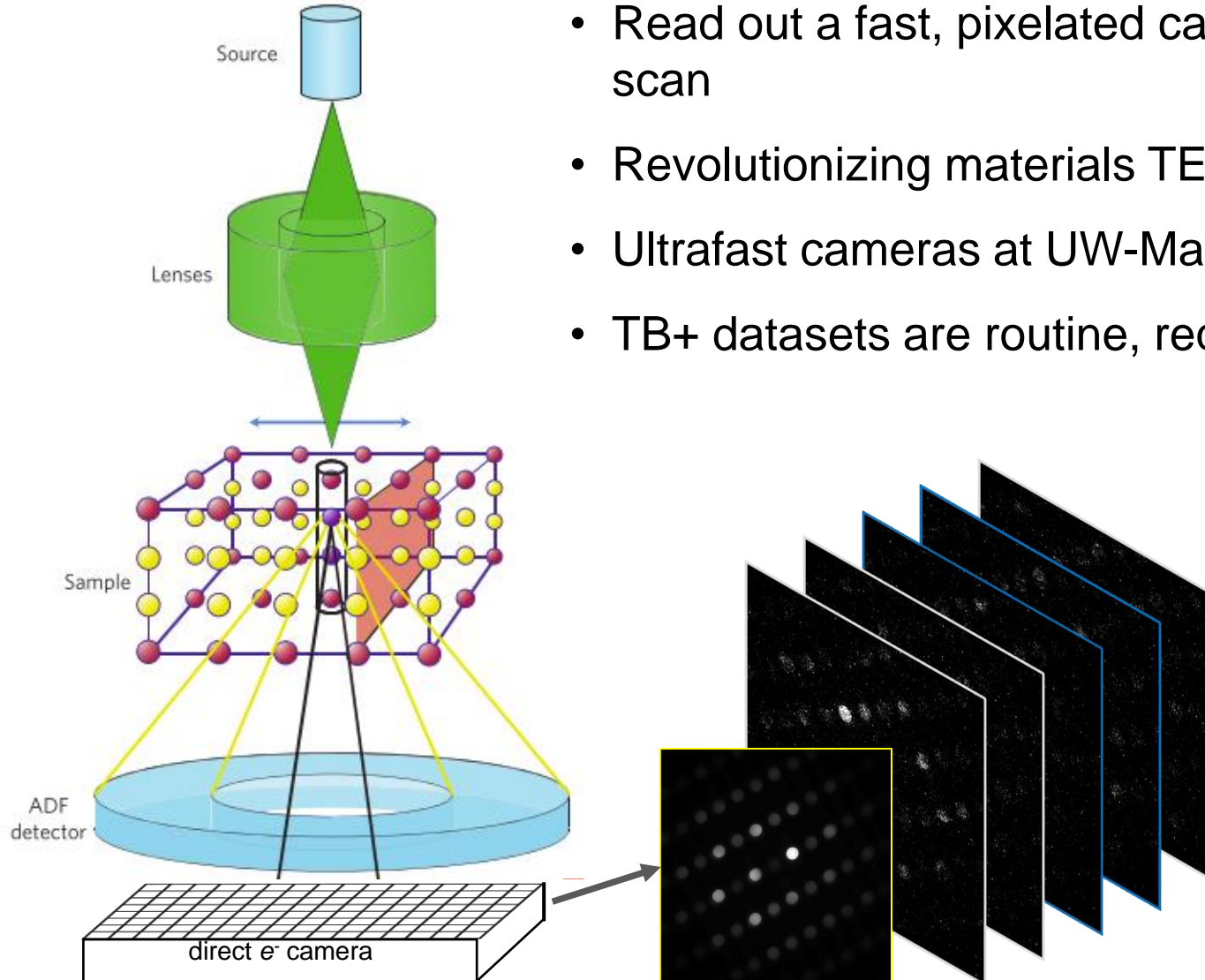UNIVERSITY OF WISCONSIN–MADISON

# *Scanning Transmission Electron Microscopy*

- Modern commercial instrument offers:
  - 1 nA into a 1 Å probe
  - 9000 EEL spectra/sec
  - Atomic-resolution STEM images at 100 nsec / pixel
  - Atomic-resolution TEM images at 3000 fps
- Routine acquisition of O(10 GB) datasets
- Specialize instruments go (much) faster

**200 kV Incident Electron Beam ($\Delta E$=1 eV)**

**1 atom wide (0.1 nm)** beam is scanned across the sample to form a 2-D image

x    y

**(Elastic Scattering ~ "Z contrast")**

**annular dark field integrating detector**

**Electron Energy Loss Spectrometer**

Increasing energy loss

# 4D STEM Data: $I(r_x, r_y, k_x, k_y)$

- Read out a fast, pixelated camera at every position of a STEM probe scan

- Revolutionizing materials TEM & STEM

- Ultrafast cameras at UW-Madison and LBL

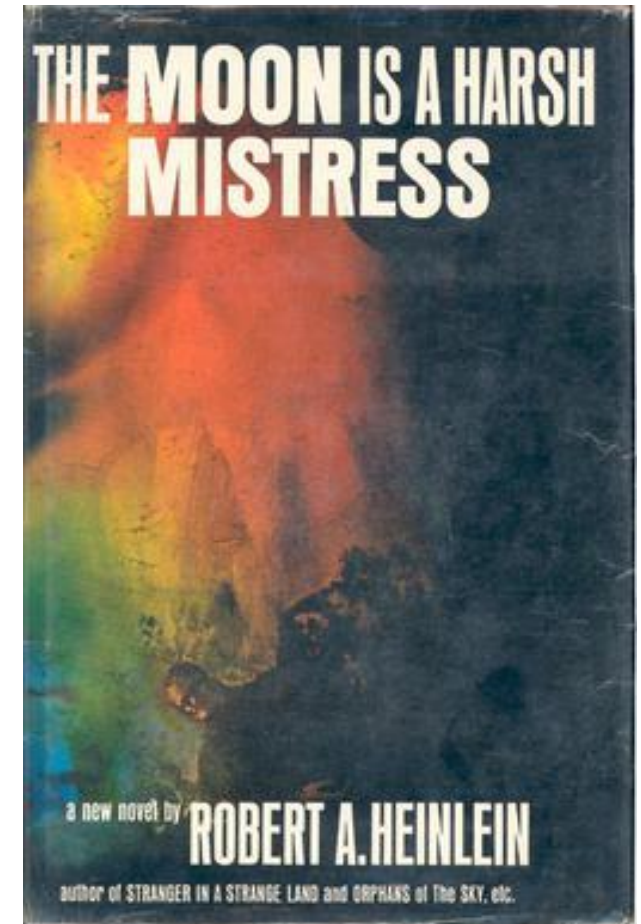- TB+ datasets are routine, requiring large scale computing and AI / ML

| Maximum Readout Speed | |
|---|---|
| Readout size | frames per second |
| 1024×1024 | 1,900 |
| 512×512 | 7,380 |
| 256×256 | 26,000 |
| 256×128 | 49,000 |
| 256×64 | 86,000 |

Source

Lenses

Sample

ADF detector

direct e⁻ camera

# *TANSTAAFL*

## *There **ain't no such thing as a** **f**ree **l**unch*

- ML / AI methods always rely on some form of prior information about

- Unsupervised learning:
  - prior information about the mathematical structure of the data
  - applications in distortion correction, denoising, spectral unmixing, and signal clustering

- Supervised learning:
  - prior information from example, already analyzed data
  - applications in finding features in images, connecting images and spectra to physical quantities of interest

- Examples of both as they apply to STEM



THE MOON IS A HARSH MISTRESS

a new novel by ROBERT A. HEINLEIN

author of STRANGER IN A STRANGE LAND and ORPHANS of The SKY, etc.

# *Distortion Correction in Scanning Images*
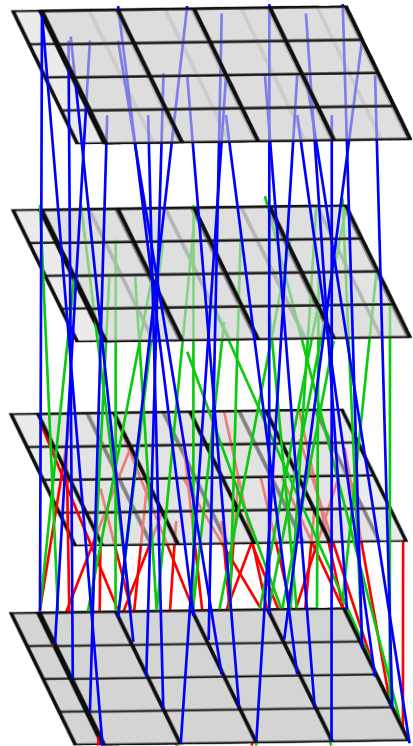


- Imperfect instruments

- Scanning images are subject to distortion arising from instrumental instabilities

- Distortions can be corrected from a series of images if the object is unchanged and the distortions are random

- More prior information:
  - Higher frequency distortions are smaller in magnitude (*e.g.* electronic jitter vs floor vibrations)
  - Lowest frequency distortions are rigid motion of the sample + shear of the image
  - Higher frequency distortions have zero mean over many images

# *Distortion Correction by Non-Rigid Registration*



$f_3$**$f_3$**
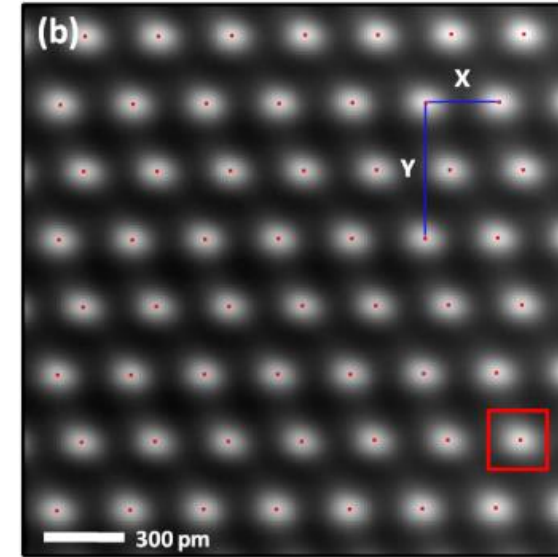
$f_2$**$f_2$**

T3,r

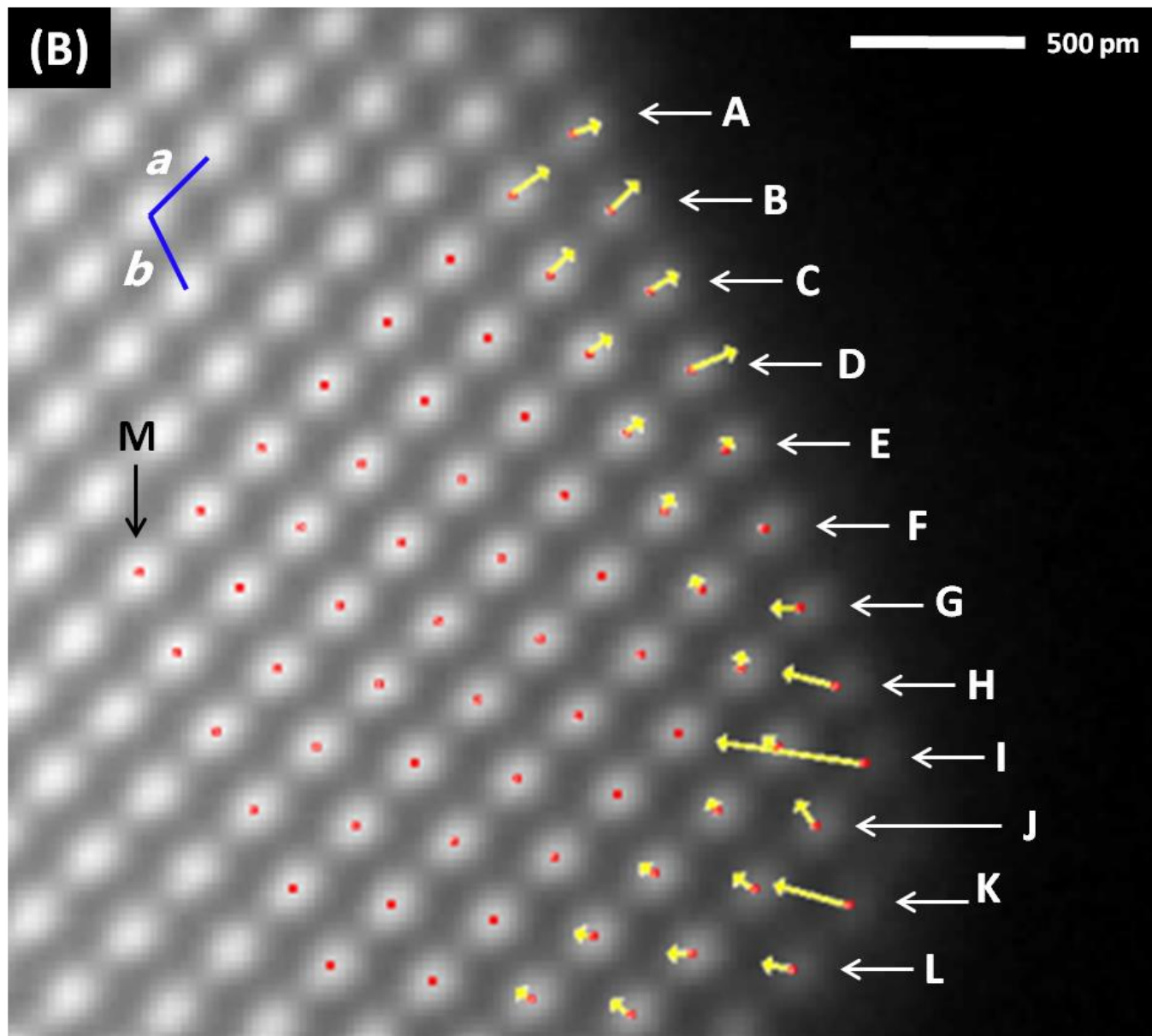$f_1$**$f_1$**
T2,r

T1,r

$f_0$**$f_0$**$f_r$

Raw series    Non-rigid aligned series    average image

- Average image has high SNR and low distortion
- Enables measurement of atom positions with <1 pm precision

Andrew Yankovich   Benjamin Berkels

# Pt on SiO$_2$ Catalyst

- Catalysis happens preferentially at corners and edges of nanoparticles

- Atoms at corners and edges lack some neighboring atoms

- They have shorter bonds that atoms inside the particle

- We can measure those bond lengths more accurately that it is possible to calculate them.



| (C) Atom | Δd (pm) |
|----------|---------|
| A | 9 ± 2 |
| B | 12 ± 4 |
| C | 10 ± 2 |
| D | 15 ± 3 |
| E | 4 ± 2 |
| F | 2 ± 2 |
| G | 8 ± 3 |
| H | 16 ± 5 |
| I | 44 ± 10 |
| J | 9 ± 3 |
| K | 23 ± 7 |
| L | 8 ± 4 |

Nat. Comm. **5**, 5144 (2014)

# Non-local Means Denoising

- Prior information:
  - High-resolution images of crystals contain many repeating features
  - Electron detection experiments are corrupted by Poisson noise

- Result:
  - Non-local means with periodic block matching
  - Similarity measure for Poisson noise
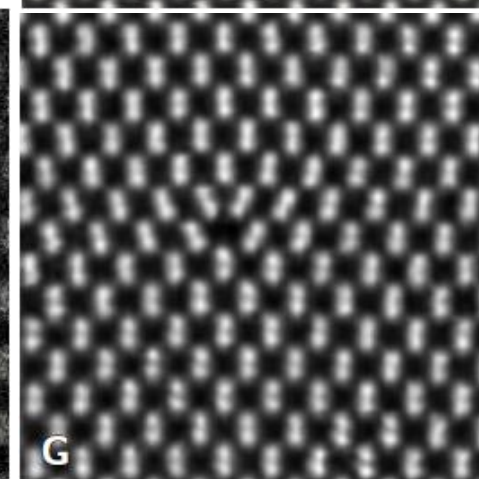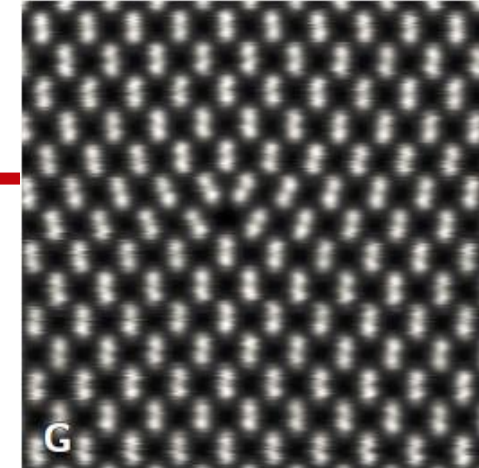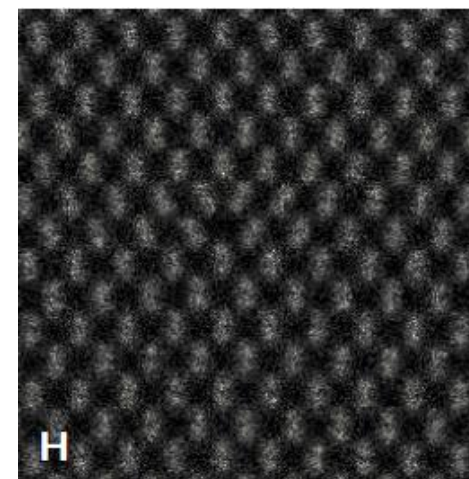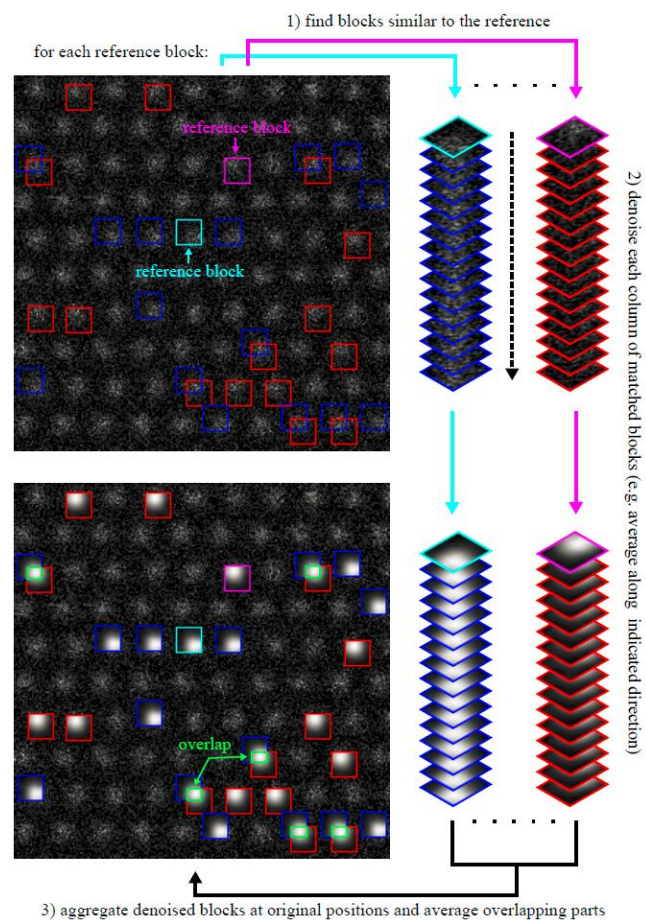  - Better performance than state-of-the-art BM3D.

N. Mevenkamp, *Adv. Struct. Chem. Imaging* **1,** 3 (2015).
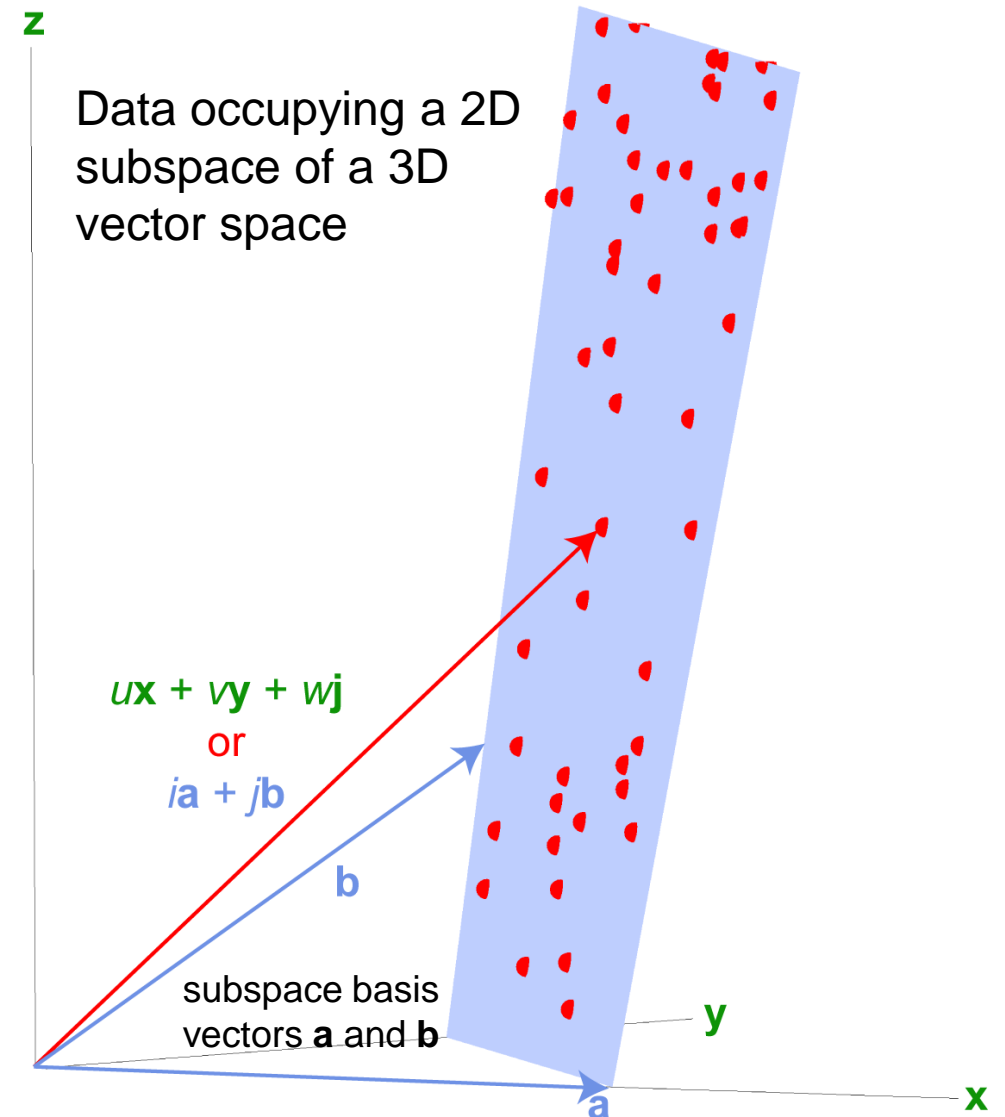
Niklas Mevenkamp
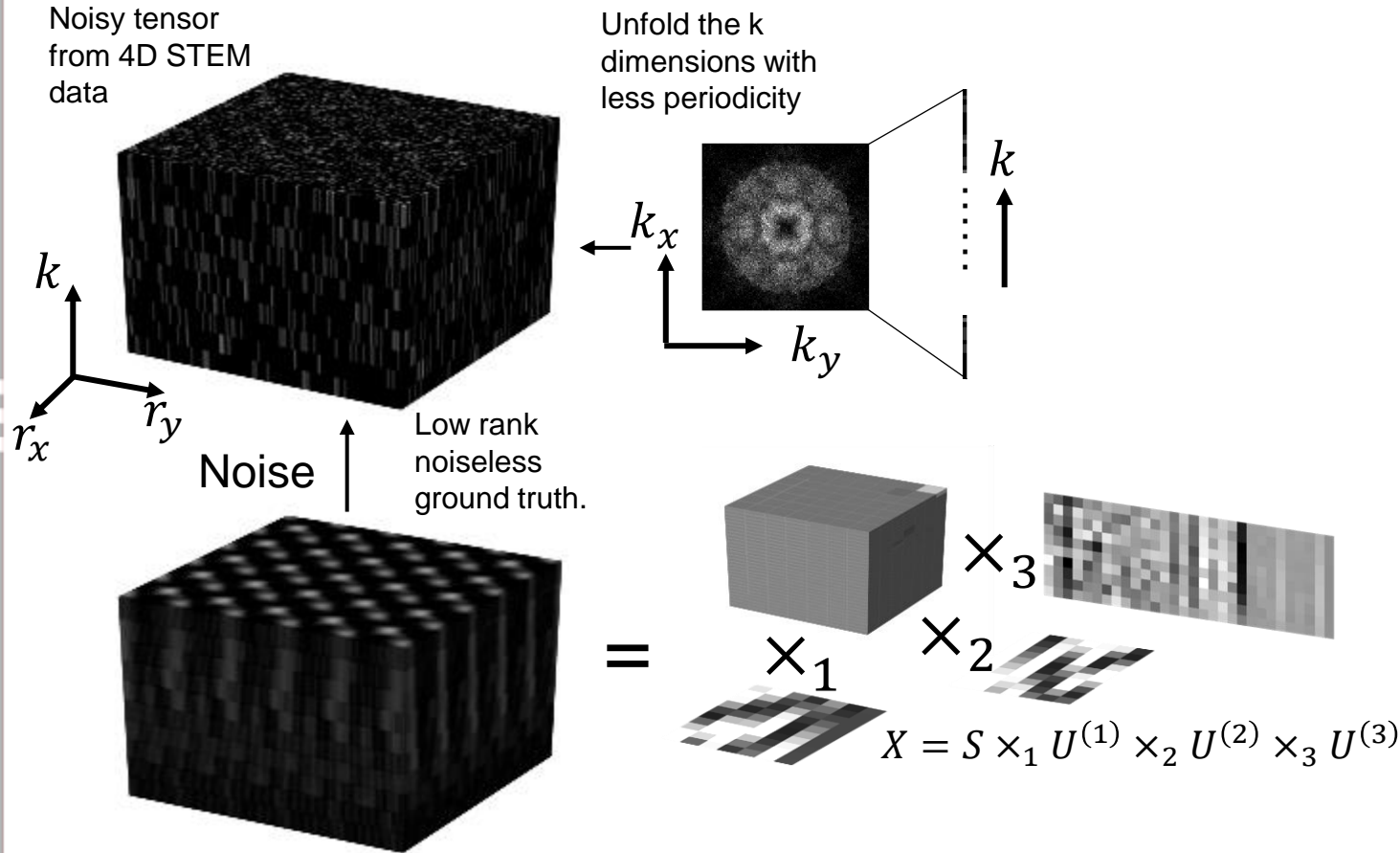
Benjamin Berkels
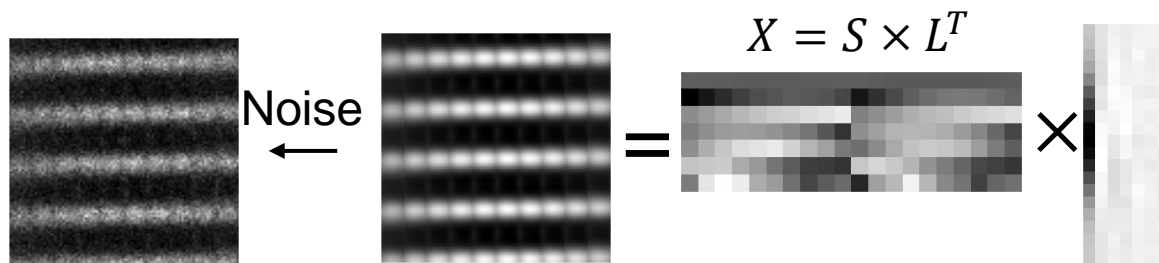
# *Dimensionality Reduction*

- Prior information: data occupy a subspace in the high dimensional vector space of the set of possible measurements

- Applications:
  - denoising by throwing away the signal outside the subspace
  - spectral unmixing to discover prototype signals or for mapping

- Methods for 2D matrices like PCA are widely used, but do not exploit the full structure of higher-dimensional data

M. Bosman, *Ultramicroscopy* **106**, 1024–32 (2006)



Data occupying a 2D subspace of a 3D vector space

$u\mathbf{x} + v\mathbf{y} + w\mathbf{j}$
or
$i\mathbf{a} + j\mathbf{b}$

**b**

subspace basis vectors **a** and **b**

**a**

# *Tensor Singular Value Decomposition*



Noisy tensor from 4D STEM data

Unfold the k dimensions with less periodicity

$k$

$k_x$

$k_y$

$k$

$r_x$ $r_y$

Noise

Low rank noiseless ground truth.

$\times_3$

$=$ $\times_1$ $\times_2$

$X = S \times_1 U^{(1)} \times_2 U^{(2)} \times_3 U^{(3)}$

Noise

$X = S \times L^T$

$=$ $\times$

- Generalization of SVD to work on tensors of any shape
- Can be applied to <span style="color:red">arbitrary high dimensional data</span> while maintaining structure along all dimensions.
- Preserves inherent structure in the data, aiding learning when the data are highly redundant, like atomic-resolution 4D STEM data
- Iterative estimate, not closed-form solution like 2D SVD.
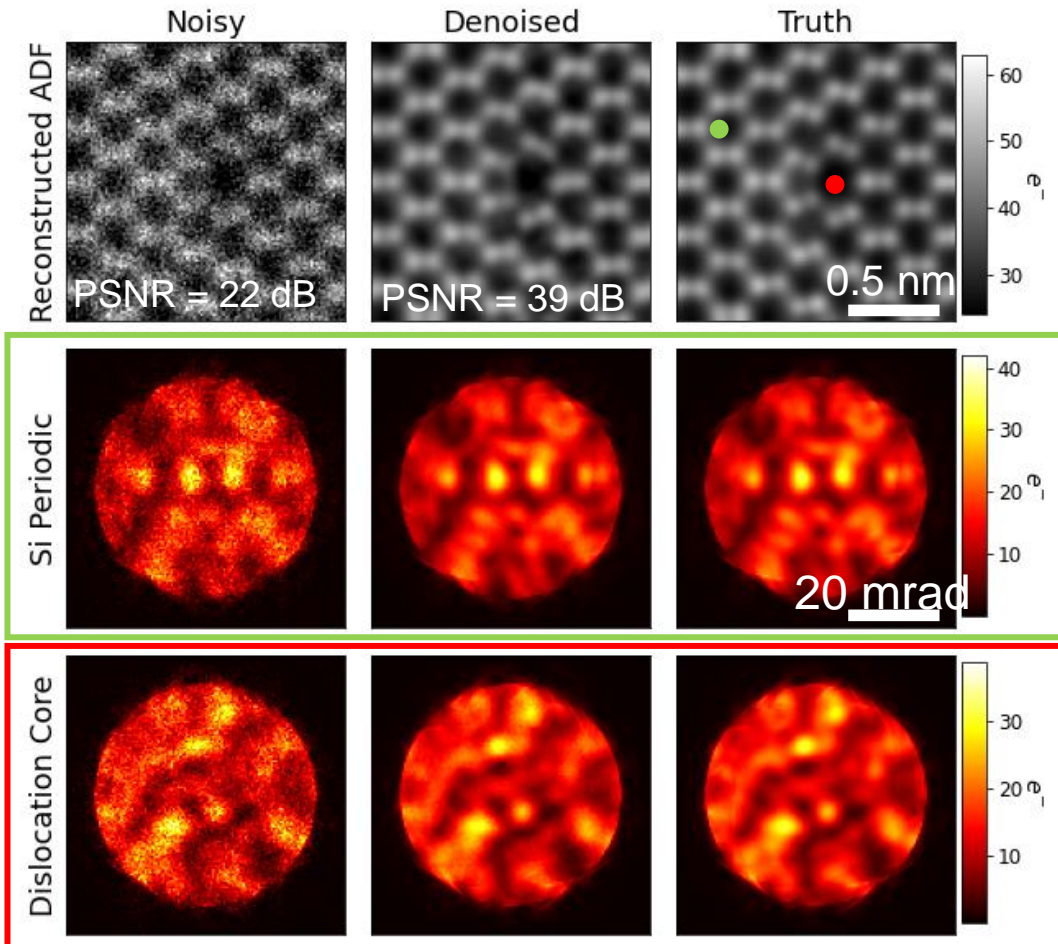- Find low rank ground truth from noisy input data.

10

# *Tensor SVD Performance: Simulated 4D STEM*

Simulations: Si [110] dislocation core

Experiments: $SrTiO_3$ [100]
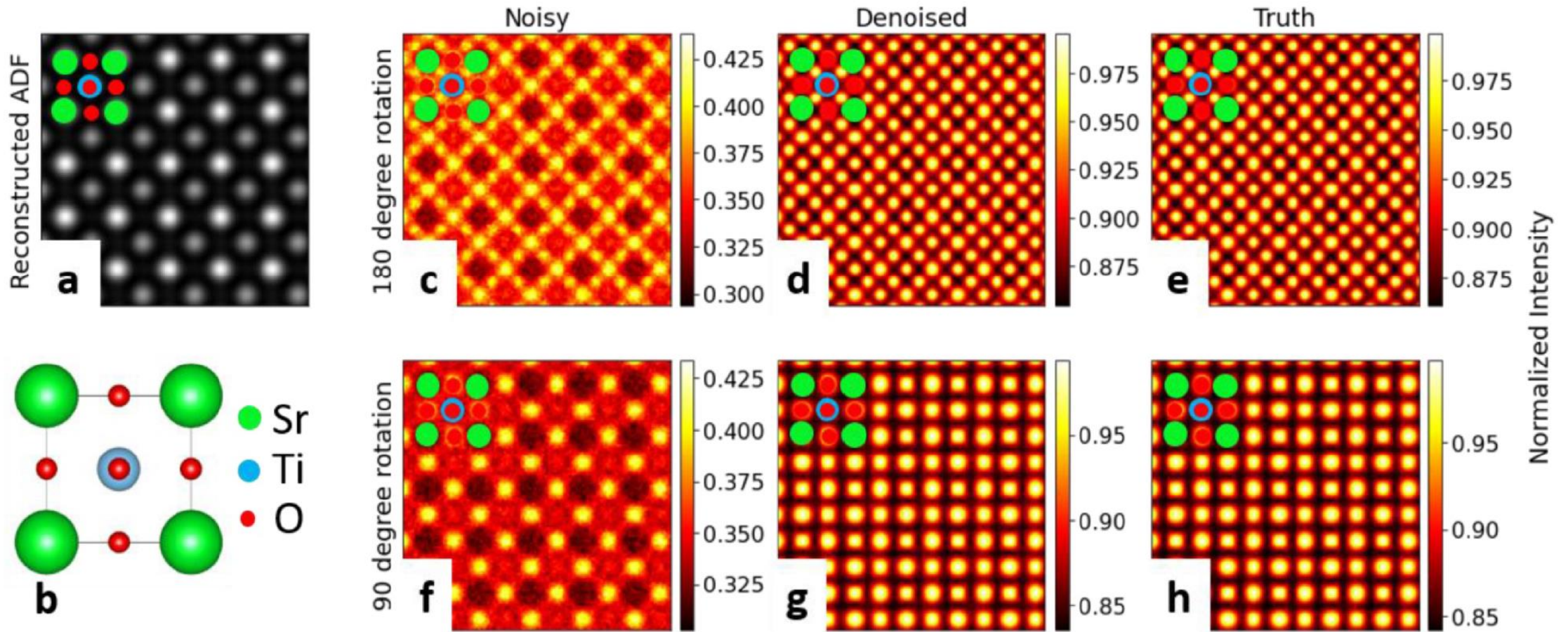


Input size: 1.6 GB    Processing time: 525.6 sec

Input size: 2.8 GB  Processing time:538.9 sec

- Processing time on a desktop with moderate computing power (single Xeon E5-2603 CPU).

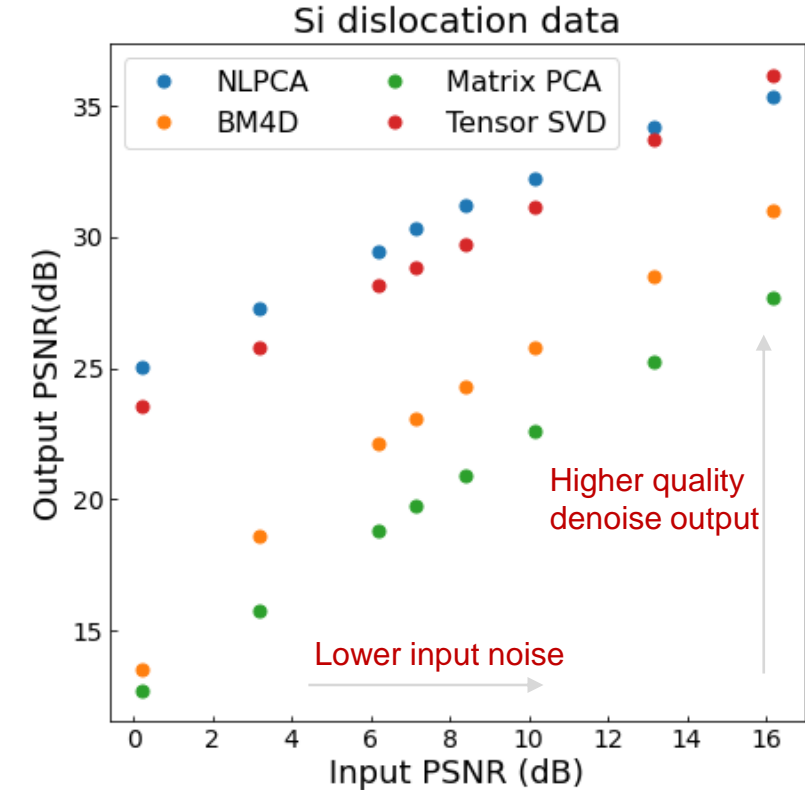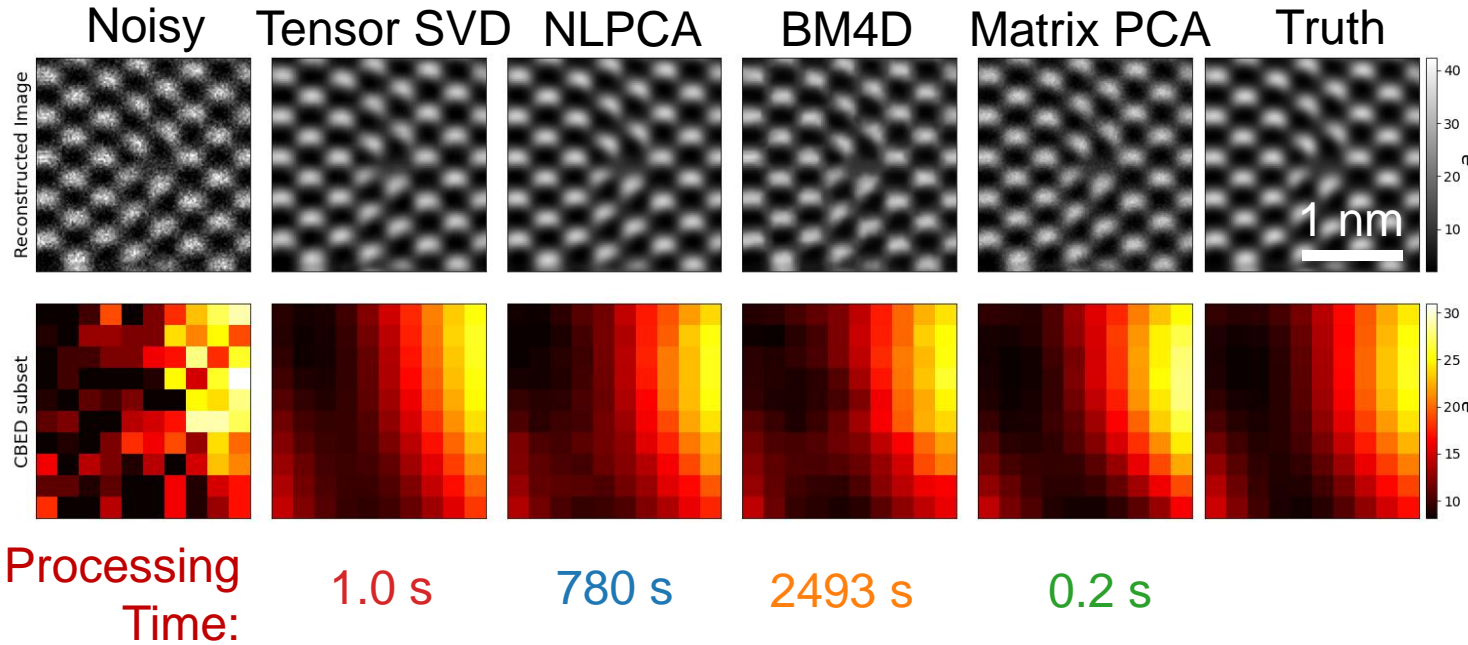# *Tensor SVD Improves Symmetry Information*



- Symmetry STEM is a new method to extract crystallographic point symmetries from 4D STEM data.

- Noisy 4D STEM data do not report the correct 4-fold symmetry for Sr sites, but denoised data do.

M. Krajnak, J. Etheridge, *Proc. Natl. Acad. Sci.* **117**, 27805–27810 (2020).

# Comparison to Other Denoising Methods

Input size: 10.0 MB



| | Noisy | Tensor SVD | NLPCA | BM4D | Matrix PCA | Truth |
|---|---|---|---|---|---|---|
| Processing Time: | | 1.0 s | 780 s | 2493 s | 0.2 s | |



Si dislocation data
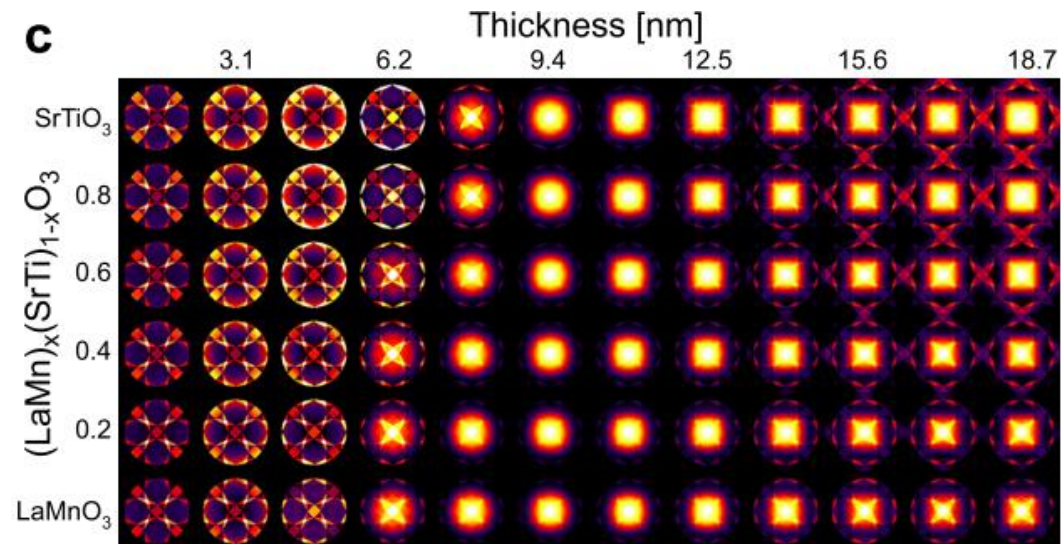
Higher quality denoise output

Lower input noise

- Tensor SVD is tested against non-local principal component analysis (NLPCA), block matching and 4D filtering (BM4D), and matrix PCA.

- Tensor SVD has the best or close to the best denoising performance.

- Tensor SVD is fast and suitable for multi-GB hyperspectral data.
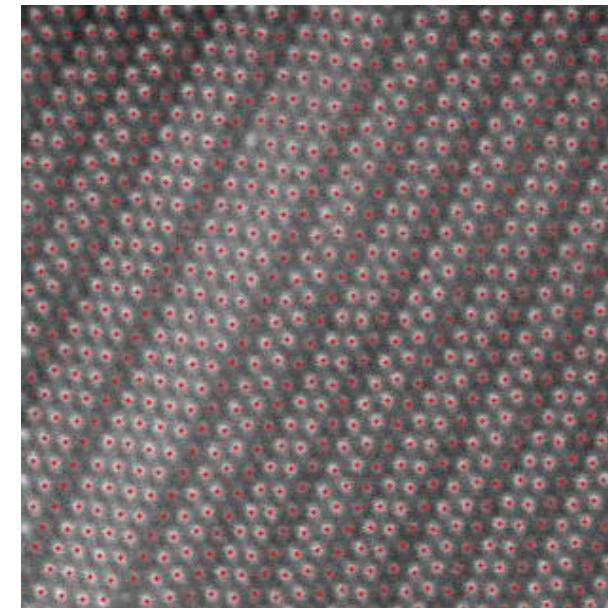
# *Supervised Learning with Neural Networks*

- Prior knowledge is example data, labeled with the result of the analysis

- For STEM, training data can come from simulations

- Limits of the resulting network are not very well determined

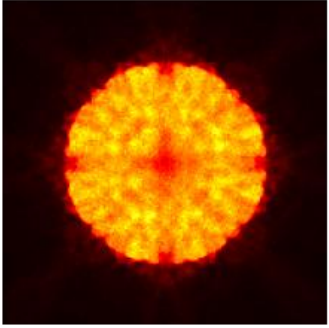Determining sample thickness from 4D STEM data

Finding atomic column locations in HRSTEM images



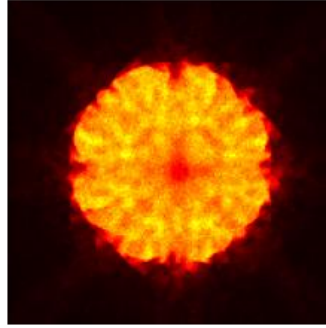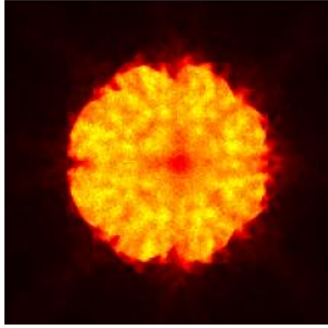CBED patterns changing with sample composition and thickness. C. Ophus, *Appl. Phys. Lett.* **110**, 063102 (2017).

# *Simulated Training Data for CNN*

Ideal PACBED   PACBED with tilt



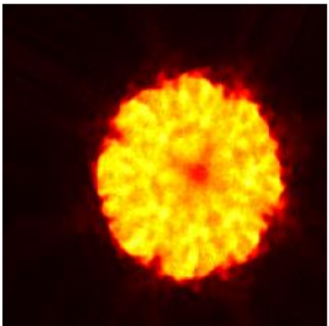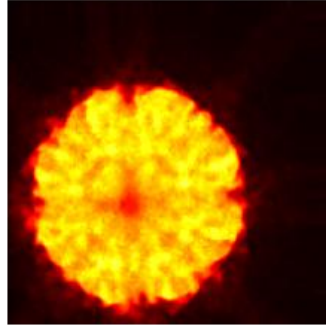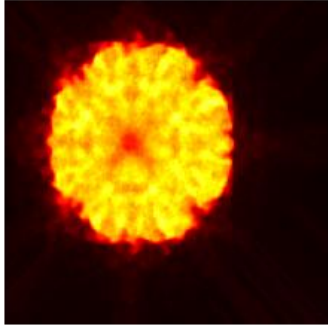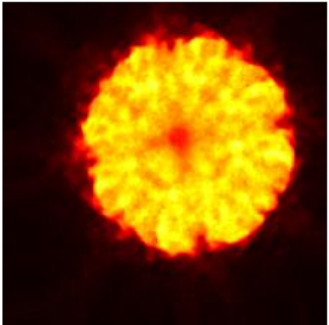After random image augmentation



- **Use multislice simulations to generate automatically labeled training data:**
  - Cover a wide range of possible experiment conditions in simulations, including thickness but also crystal tilt
  - Augment the images after simulations by adding noise, and distortions including shift, zoom, rotation, shear, *etc.*
- Transfer learning:
  - use a vgg-16 network pretrained to recognize features in natural images
  - retrain just the fully connected final layers at first, then tweak the convolutional layers only at the end
- Full training data set is about 750 GB

- RMS deviation for experimental data on the same crystal and orientation is ±1 nm

- Larger thicknesses and more complicated image features work less well.

- Network fails for thicknesses outside training data set.

- Network fails for other crystals or even other orientations of the same crystal.

# *Atom Finding: A Common Problem*

Input



**FCN**

Output

Atomic column coordinates

## Lin's AtomSegNet

*Scientific Reports* 11.5386 (2021): 1-15

- Functionalities: **atom segmentation**, noise reduction, background removal, and super-resolution processing.
- Trained on 15 crystal lattices (e.g. $SrTiO_3$, graphene)

## Ziatdinov's AtomNet

*ACS Nano* 11.12 (2017): 12742–12752

- Functionalities: **atom segmentation**, detecting atom species and defects.
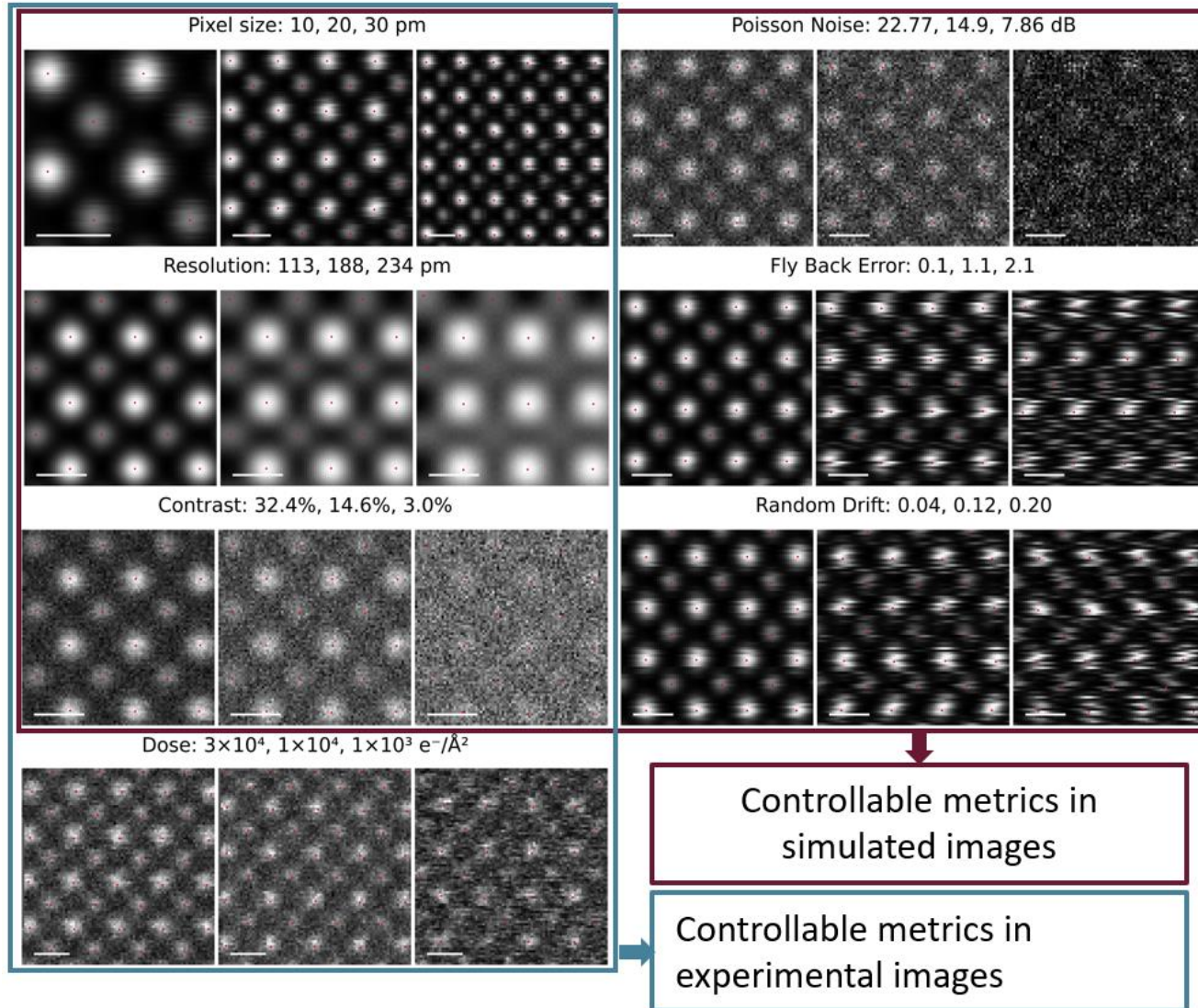- One trained on crystal lattices (e.g. $SrTiO_3$), another trained on hexagonal lattices.

## Ziatdinov's AtomAI

https://github.com/ziatdinovmax/atomai

- Pytorch-based package for training new models for new problems
- We trained a new U-net model on 5 crystal lattices using AtomAI

# *Which Model is "Best"?*



Pixel size: 10, 20, 30 pm

Poisson Noise: 22.77, 14.9, 7.86 dB

Resolution: 113, 188, 234 pm

Fly Back Error: 0.1, 1.1, 2.1

Contrast: 32.4%, 14.6%, 3.0%

Random Drift: 0.04, 0.12, 0.20

Dose: $3\times10^4$, $1\times10^4$, $1\times10^3$ e⁻/Å²

Controllable metrics in simulated images

Controllable metrics in experimental images

- We wanted to use the best network with the least investment of time, but we found no way to evaluate network performance outside their own training and test data.

- Created a benchmark data set with varying image quality:

  - In simulations, vary pixel size, contrast, Poisson noise level, scan distortion

  - In experiments, vary pixel size, spatial resolution, electron dose

  - $WS_2$ and $SrTiO_3$

- ~40 experimental images of various crystal lattices, defects, interfaces

- DOI: 10.18126/e73h-3w6n
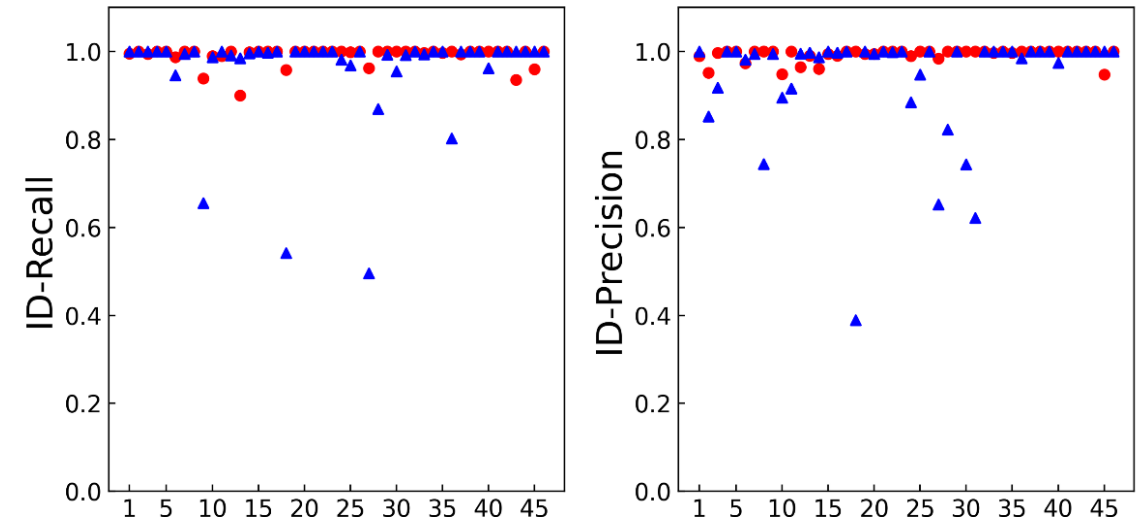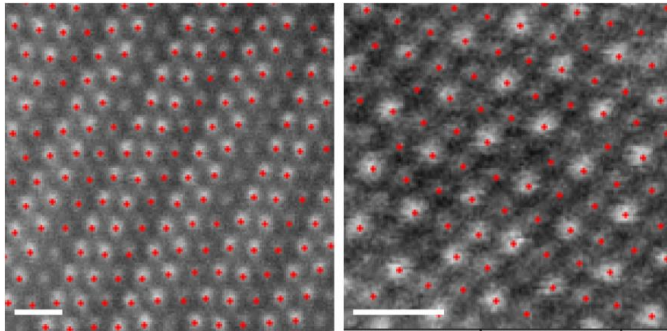
# Model Performance vs Image Quality



- Define acceptable performance as ID-recall > 0.90, ID-precision > 0.95 and $\Delta d < 0.3$ Å

- Larger blue polygons Ziatdinov's model is more forgiving of poor image quality

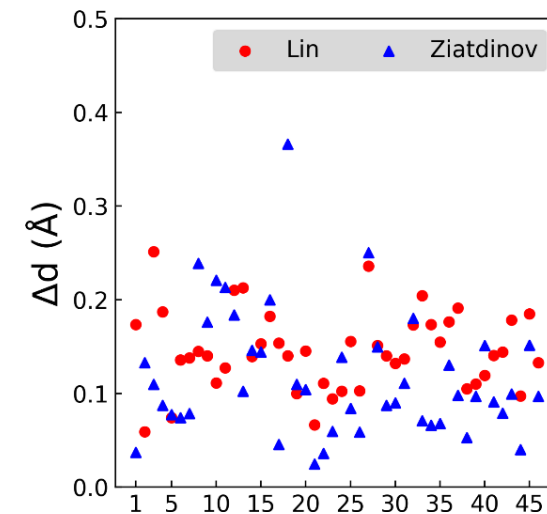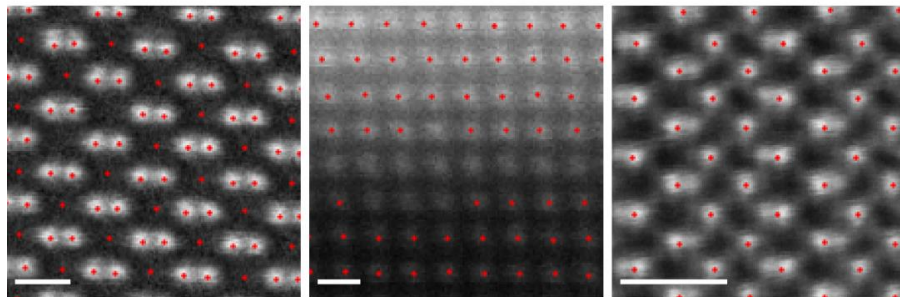- Potential trade-off between model overall performance and general applicability

# *Applicability Outside Training Crystals*

- Models applicable for most crystal lattice, defects, interface, etc.

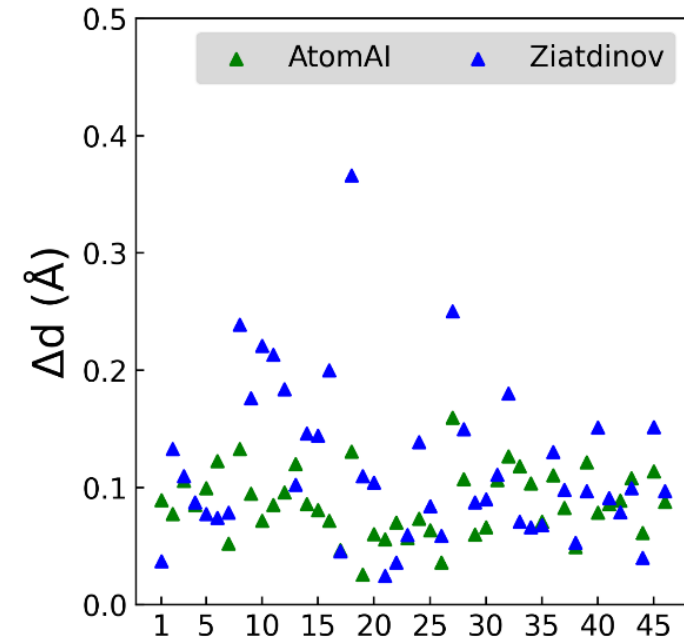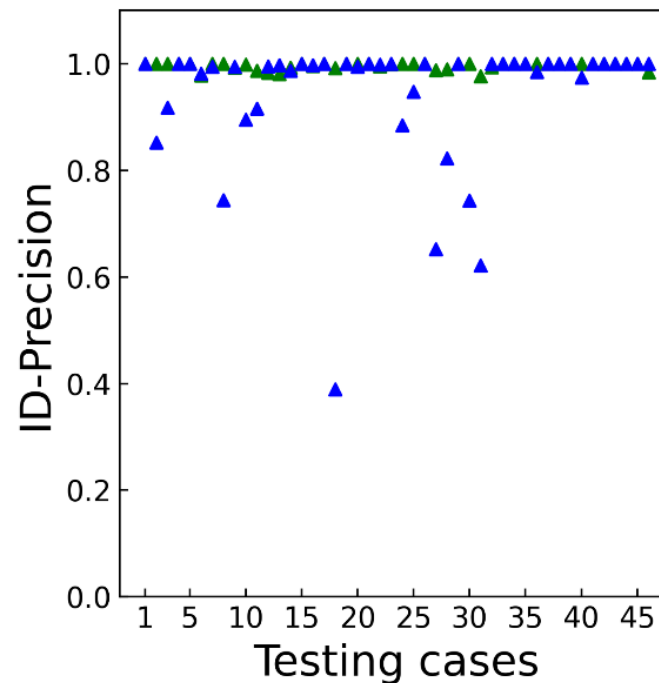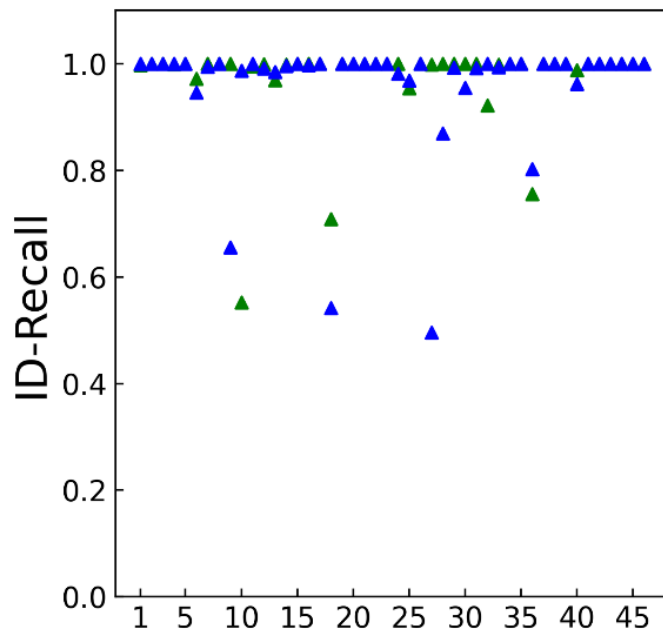- Poor cases for Lin's model due to low SNR



- Poor cases for Ziatdinov's model including FPs in background, TNs in areas of varying contrast and overlapping atoms.

# *Toward a More General Network*

- Used the AtomAI framework to train a network on simulated images from 5 crystal lattices, plus augmentation

- More general than Ziatdinov model while maintaining robustness against image quality.

# *Making Materials ML FAIR*

- **F**indable

- **A**ccessible

- **I**nteroperable

- **R**eusable

- At least in my corner of materials science, ML models are not FAIR

  - Easy find (Github) but hard to run

  - Prior knowledge / training data is often unspecified or unavailable



M. D. Wilkinson, The FAIR Guiding Principles for scientific data management and stewardship. *Sci. Data* **3**, 160018 (2016). DOI: 10.1038/sdata.2016.18

https://www.force11.org/group/fairgroup/fairprinciples

# *FAIR Data and Models*

- Tools for distributing data like figshare and Materials Data Facility are well developed
  - Non-rigid registration: 10.6084/m9.figshare.12592466.v1
  - Non-local denoising: 10.6084/m9.figshare.12592457.v1
  - tensor SVD: 10.18126/vh9q-i1l6
  - 4D STEM CNN: 10.18126/4nm2-0g70
  - Atom finding test data: 10.18126/e73h-3w6n
- Need to be more widely used

- Tools for software exist and are widely used
  - NRR and tensor SVD have python modules compatible with HyperSpy
  - Non-local denoising and the 4D STEM CNN are available on Github
- How often does research-grade software off Github actually work to solve a problem?
- How often can you test the software on the data used to develop it?

**FOUNDRY**
DATA, MODELS, SCIENCE

- Containerized ML models permanently associated with data sets
- Radically reduced barriers to reuse, meta-studies, benchmarking, and more
- Atom finder dataset available now
  - DOI: 10.18126/e73h-3w6n
  - Standard dataset description interface
  - Queriable format (hdf5)
  - Highly accessible metadata

Consumers

Science!

```
From foundry import Foundry
f = Foundry()

X,y = f.load("dataset1", v="1.0")
y_pred = f.run("model1", v="1.0", X)
```

- Models run locally or on distributed endpoints
- Capabilities to pull datasets to desired location or move compute to desired location

Dataset

Function

API layer

API layer

Data Publishers

Model Publishers

Data Provider

Models / Functions

```
f.data.publish("./"
"dataset1", v="1.1")
```

```
f.model.publish("./"
"model1", v="1.1")
```

MATERIALS DATA FACILITY

NIST CHiMaD

DLHub
Data and Learning Hub for Science
https://www.dlhub.org

U.S. DEPARTMENT OF ENERGY

Argonne
NATIONAL LABORATORY

**Dane Morgan,** Paul Voyles, Michael Ferris, Marcus Schwarting, **Ben Blaiszik**

24

# *Summary*

- STEM data are growing in rapidly in size and complexity

- ML / AI methods are essential and developing quickly

- Example applications:
    - distortion correction
    - non-local denoising
    - low-dimensional representations for tensor data
    - determining sample characteristics directly from 4D STEM data
    - atom finding in high-resolution images

- Data and models are all available from the bibliography at tem.msae.wisc.edu

- ***TANSTAAFL and make it FAIR***