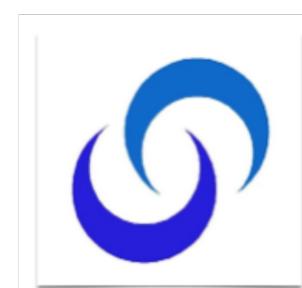
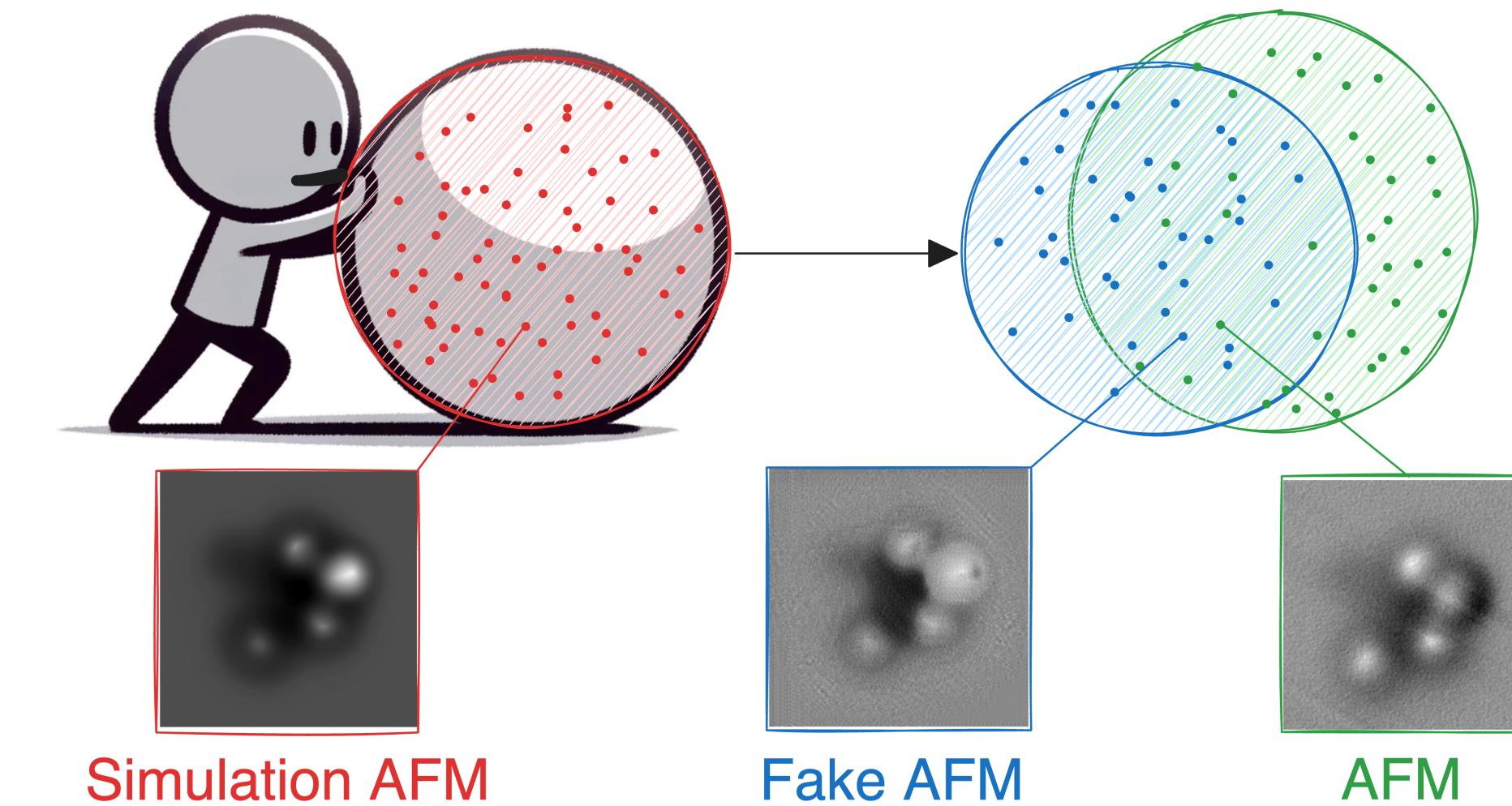


# Enhancing Atomic Force Microscopy Image Analysis: Style Translation and Data Augmentation

Jie Huang ([jie.huang@aalto.fi](mailto:jie.huang@aalto.fi))

17/09/2024

MLM24, Kanazawa, Japan



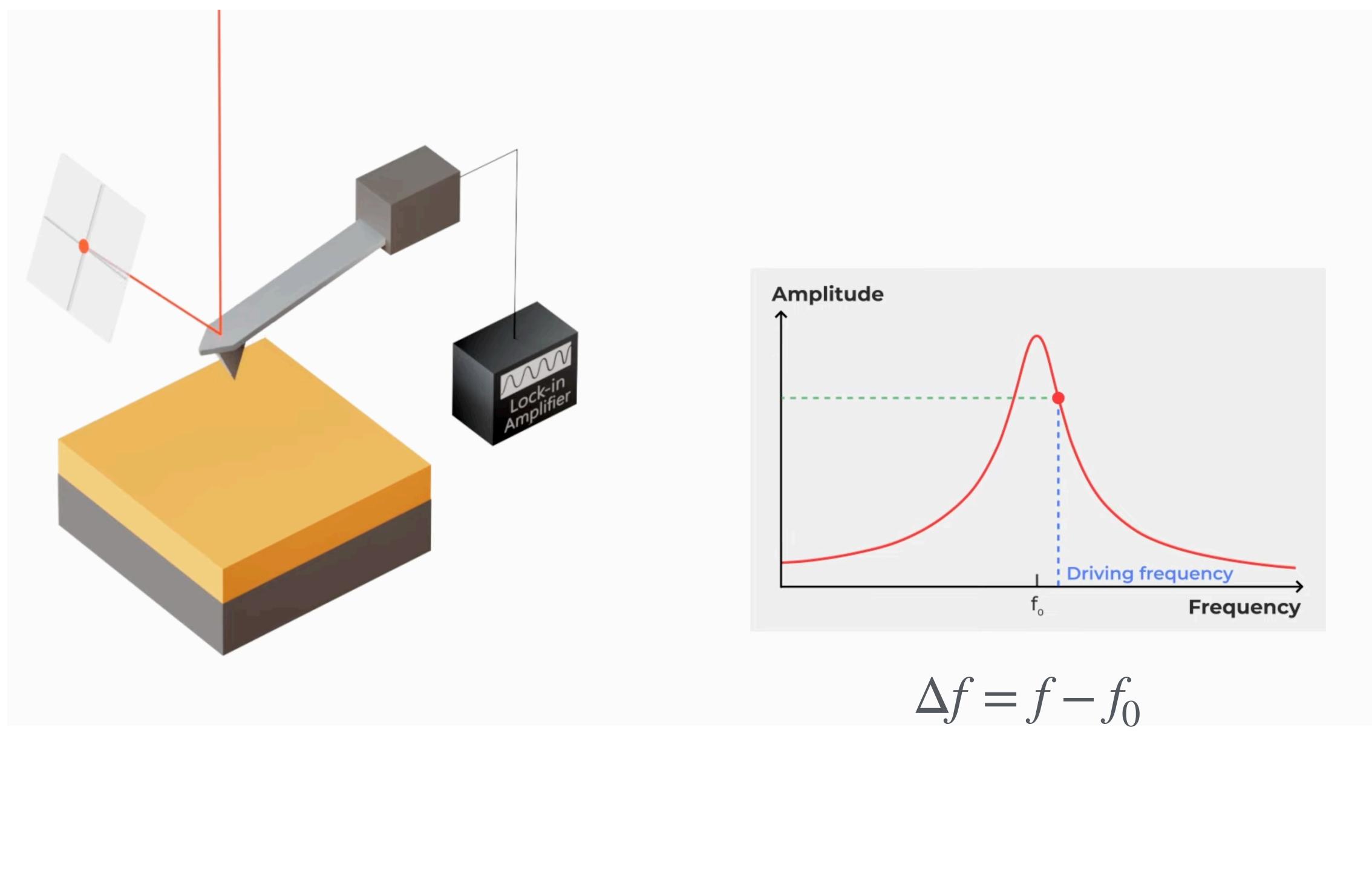
**Surfaces and Interfaces  
at the Nanoscale (SIN)**



**Aalto University  
School of Science**

# Background

## Atomic force microscopy (AFM)



Non-contact mode

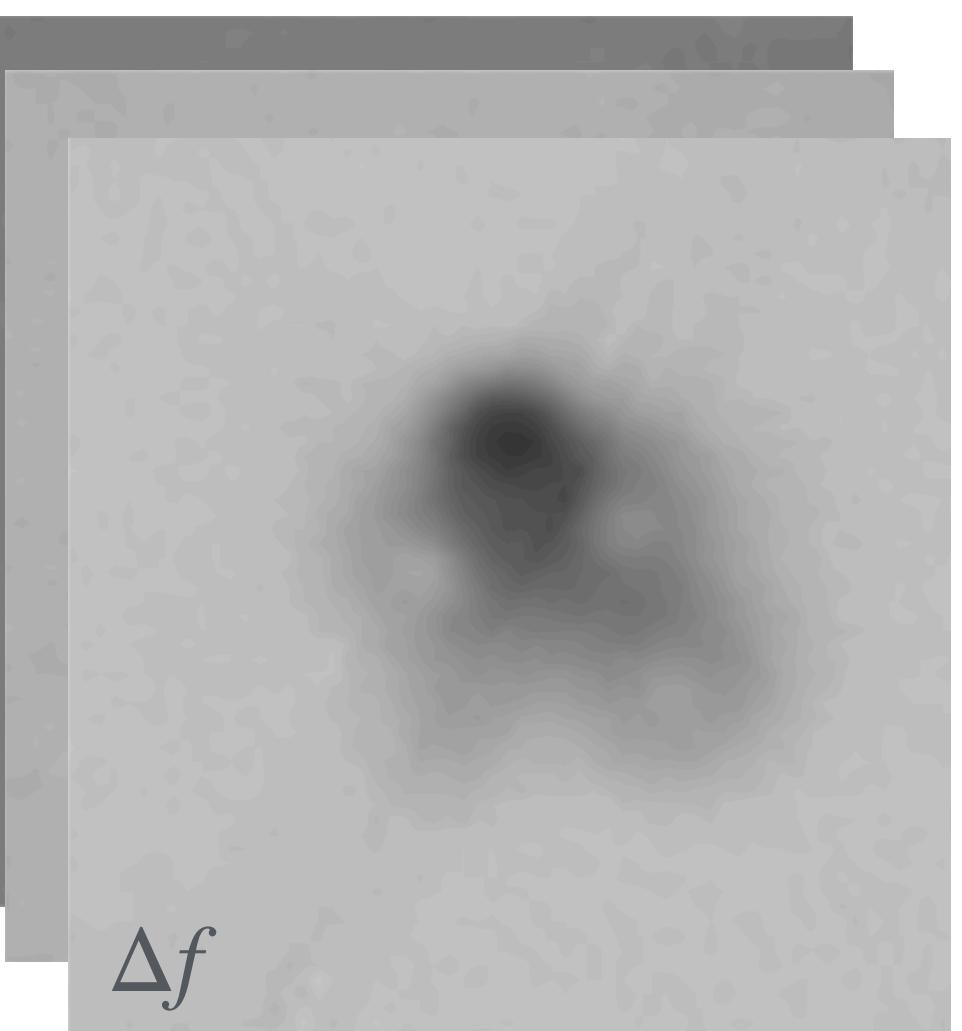
<https://youtu.be/K9FgPWK3Co4?si=hnL9WH7ONWRUkVxJ>

Prepare substrate  
and sample in lab



AFM

AFM images.

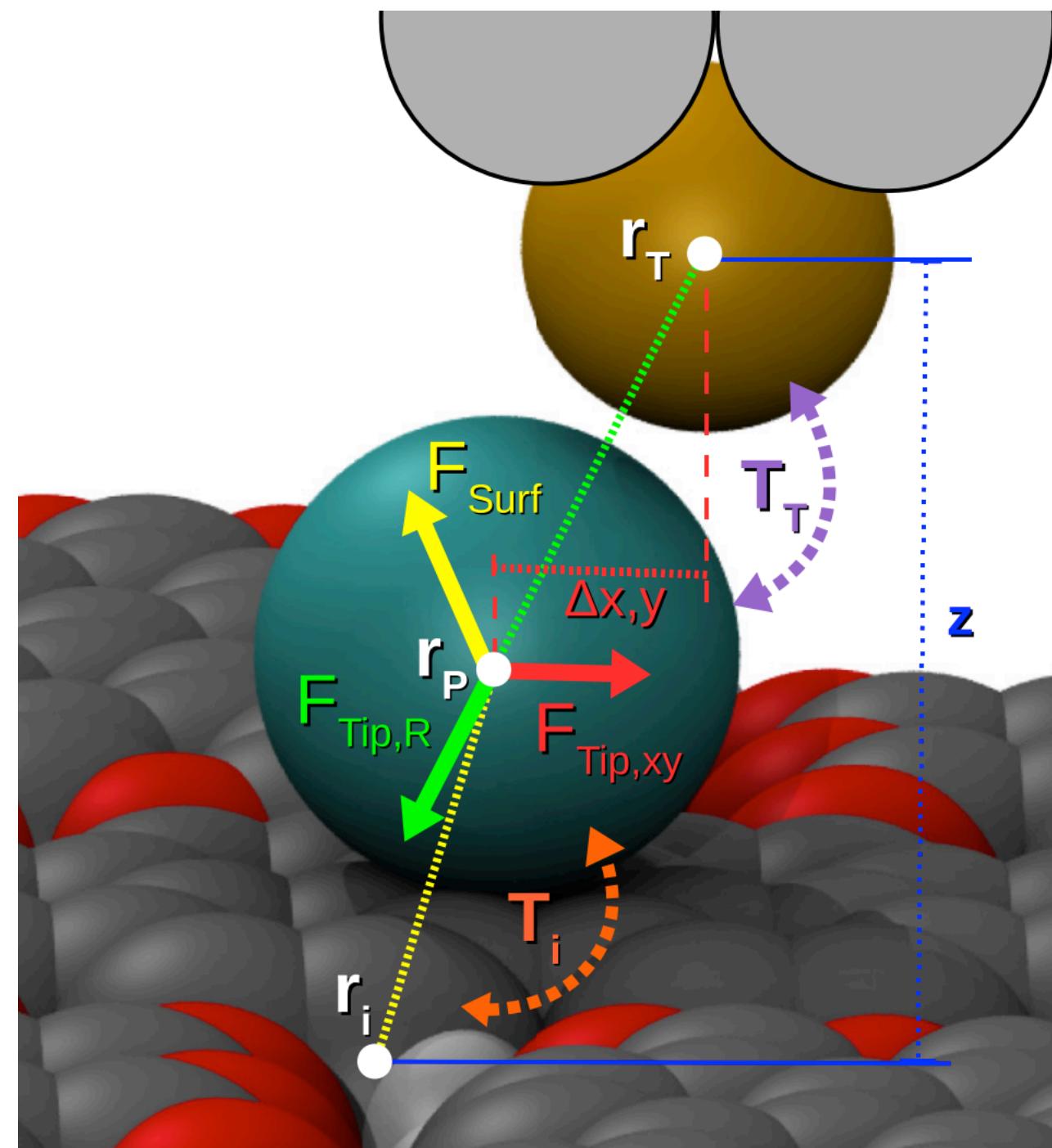


$\Delta f$

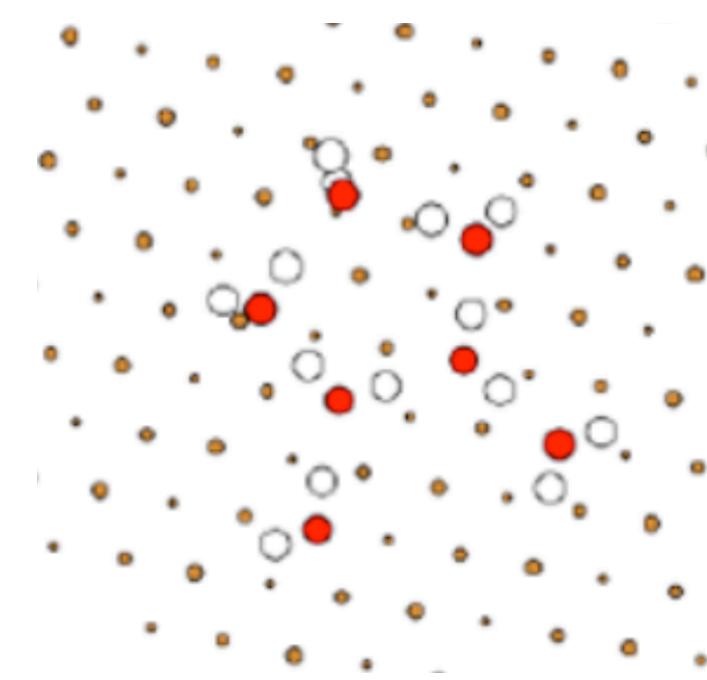
Can we know the exact configuration/structure?  
Where is O and H if it's a sample of water molecules?

# Background

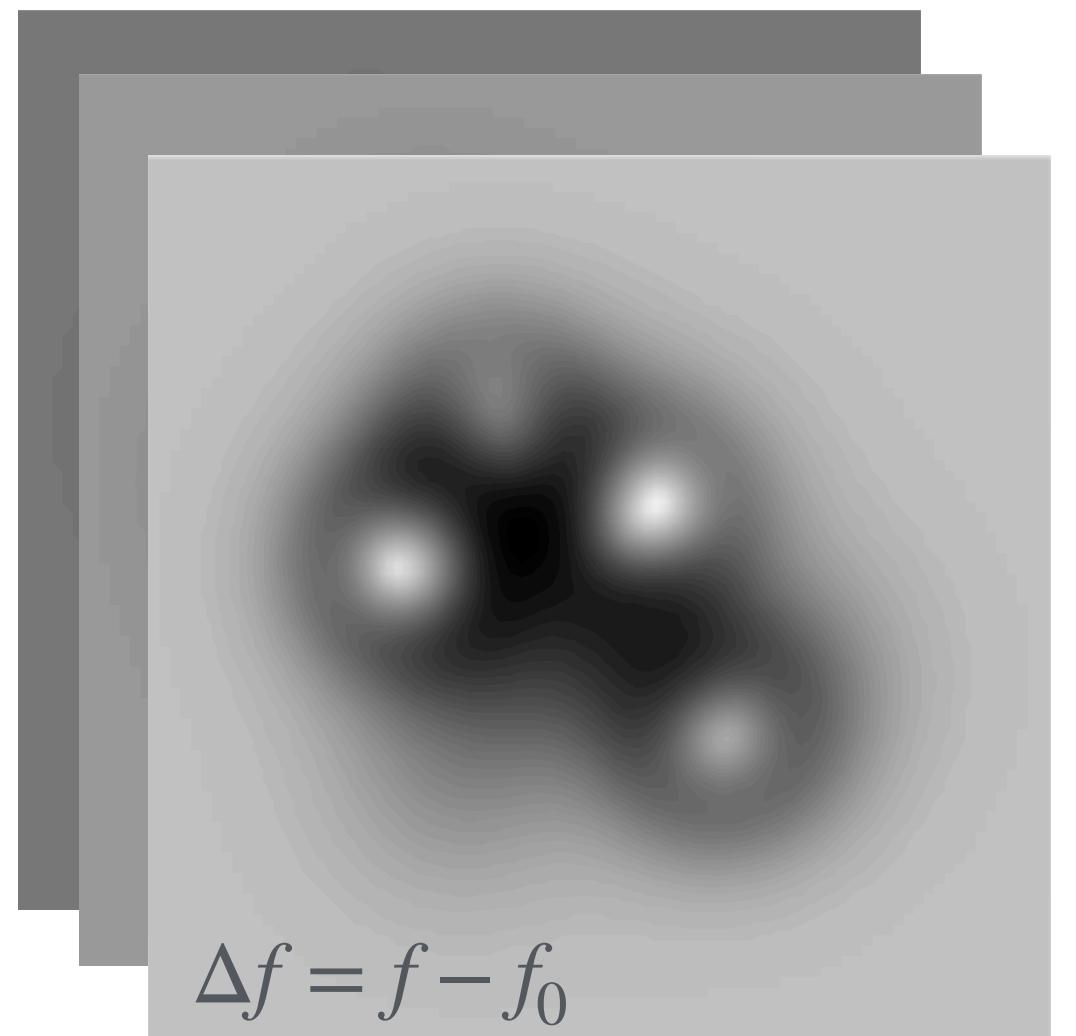
## Probe-Particle AFM (PPAFM)



Prepare substrate  
and sample using  
DFT calculations



PPAFM



Simulation  
AFM images.

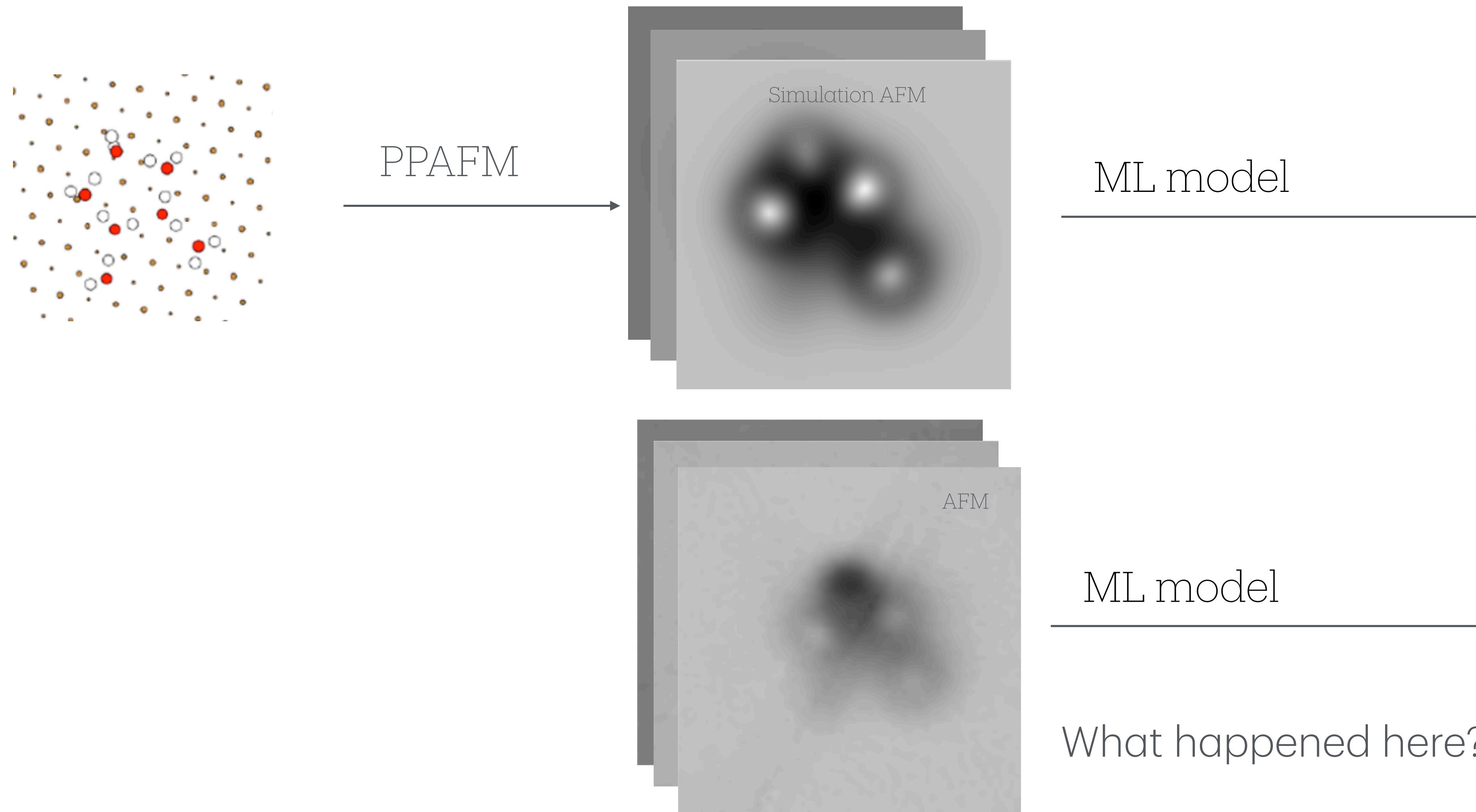
Where is O and H?

<https://github.com/Probe-Particle/ppafm>

Hapala et al., Phys. Rev. Lett., 113, 226101 (2014).

Priante et al., ACS Nano, 2024, DOI: 10.1021/acsnano.3c10958.

# Machine Learning (ML) Applications



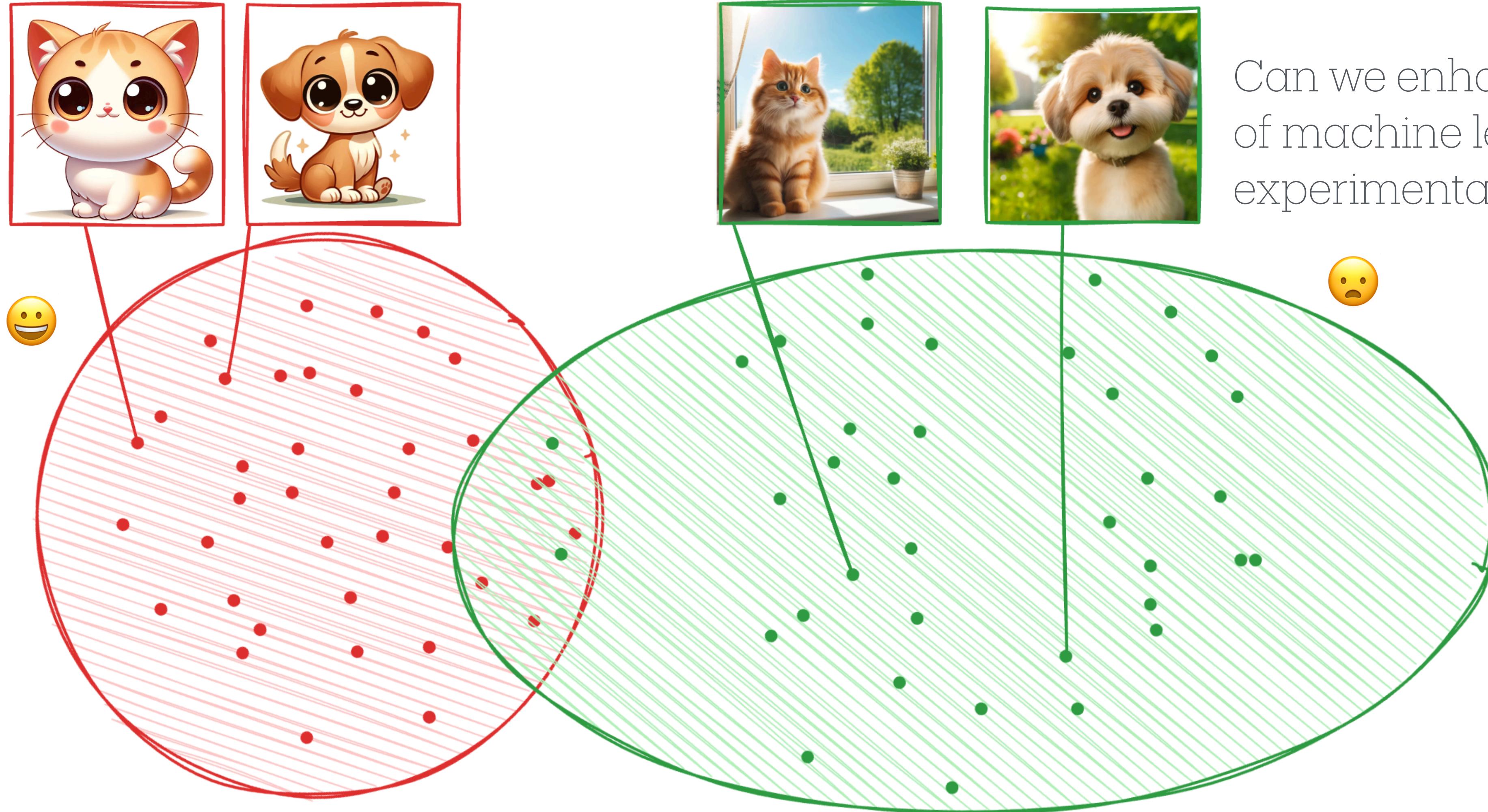
Oinonen et al., MRS Bulletin, 47(9), 895-905 (2022), DOI: 10.1557/s43577-022-00324-3.

Oinonen et al., ACS Nano, 16(1), 89-97 (2021), DOI: 10.1021/acsnano.1c06840.

Kurki et al., \*ACS Nano\*, 18(17), 11130-11138 (2024), DOI: 10.1021/acsnano.3c12654.

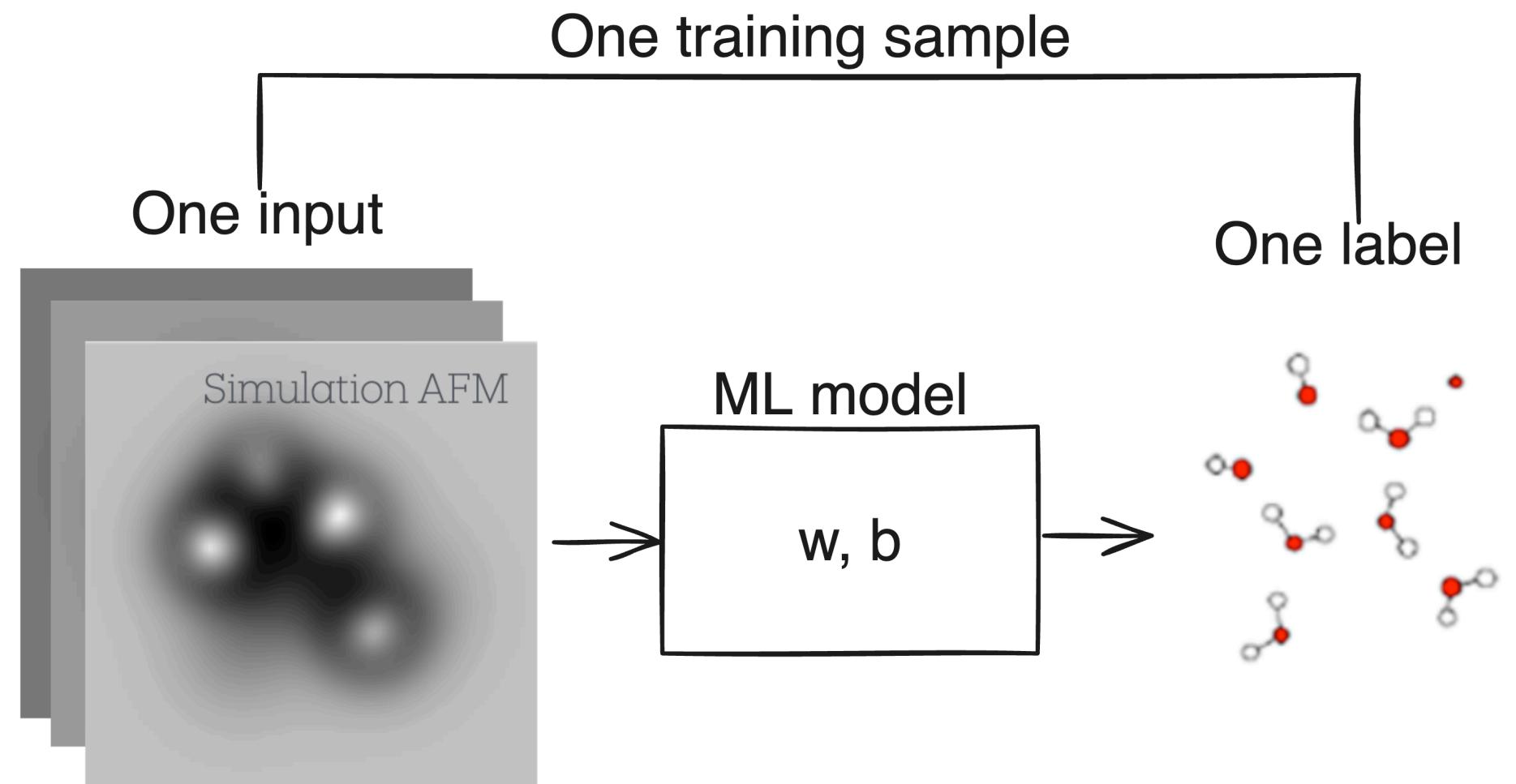
Priante et al., ACS Nano, (2024), DOI: 10.1021/acsnano.3c10958.

# Distribution shift



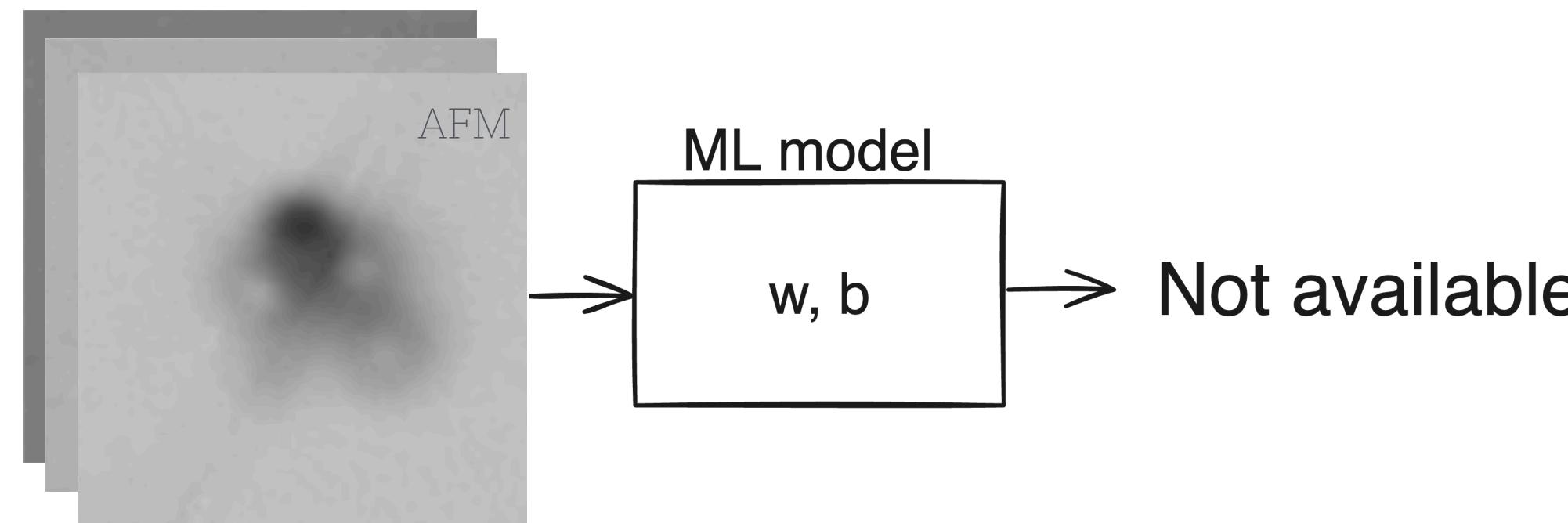
Can we enhance the performance  
of machine learning(ML) models on  
experimental AFM images?

# Challenges in machine learning



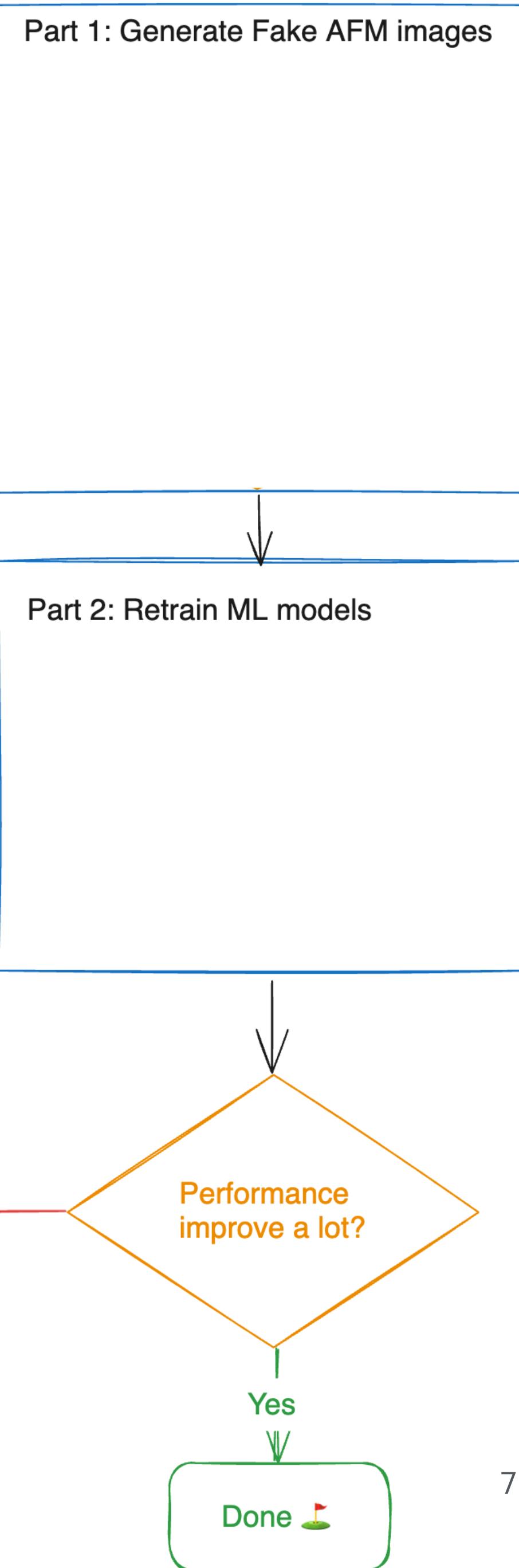
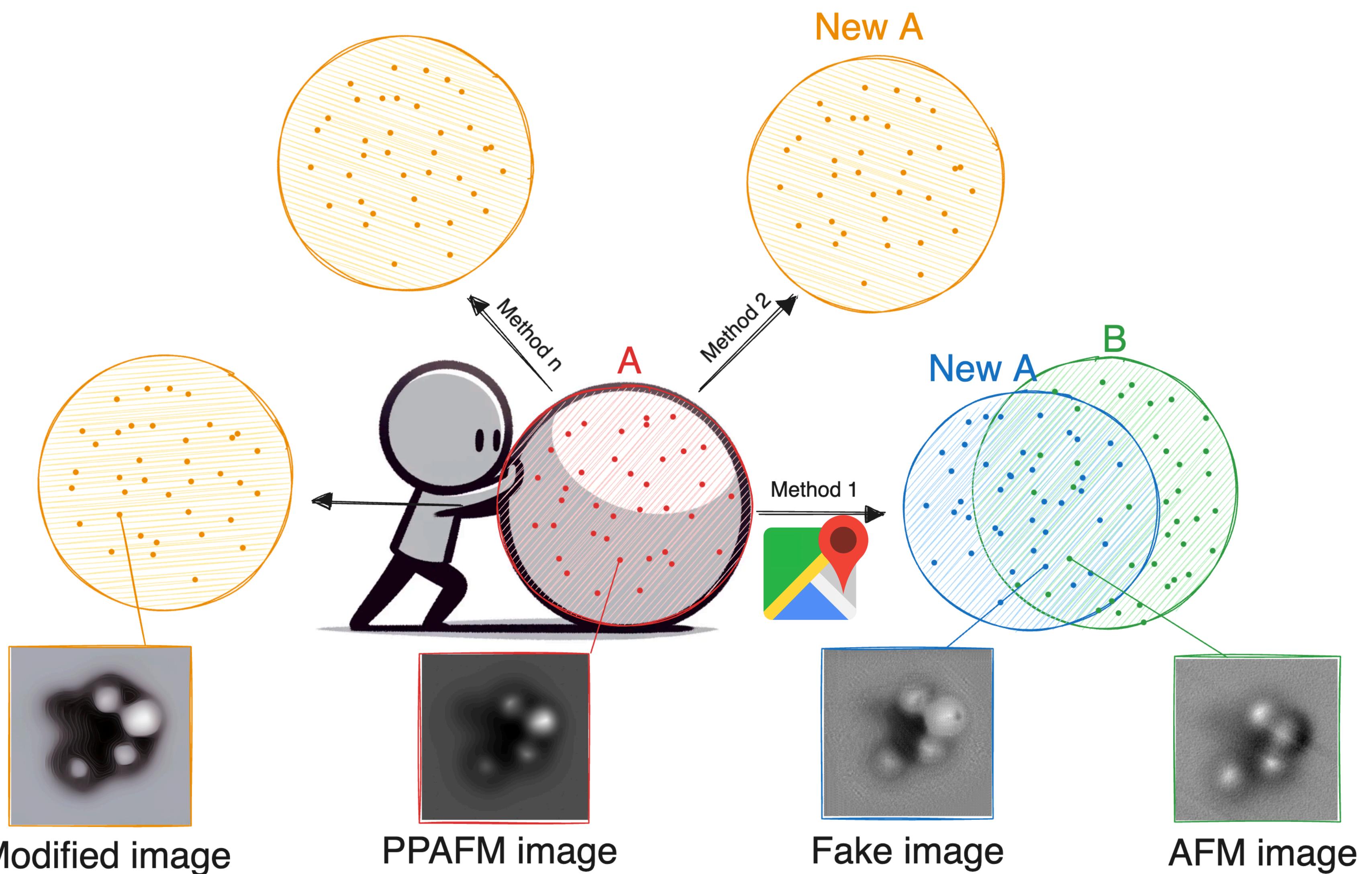
16 000 training samples  
in Water-Au111 dataset.

Cost too much money & time.



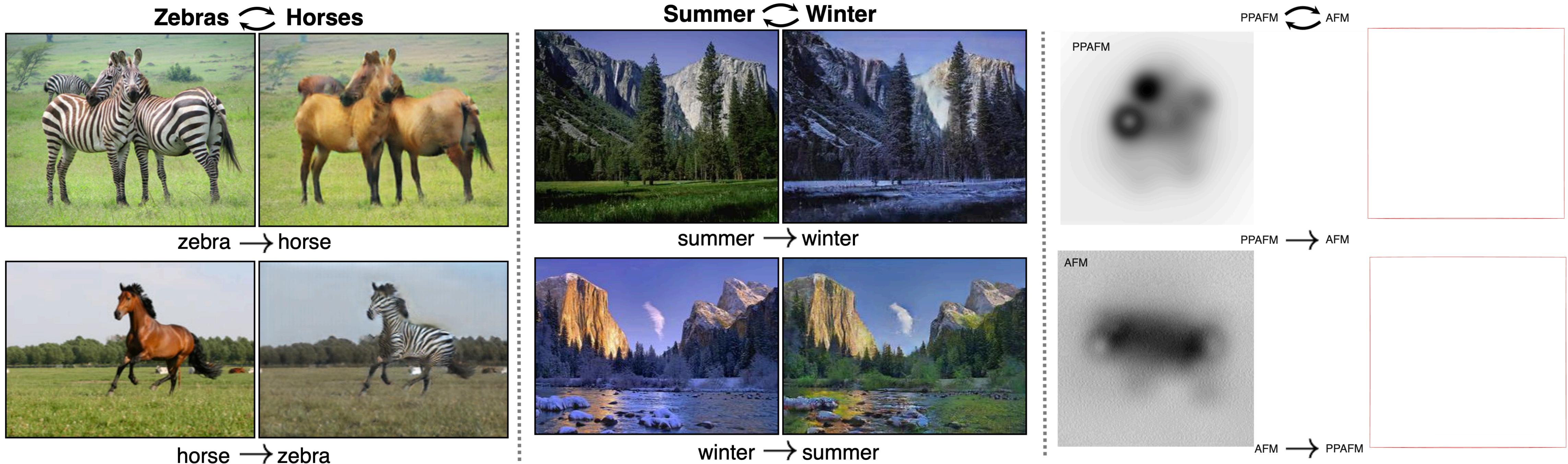
Label not available.

# Hypothesis & Workflow

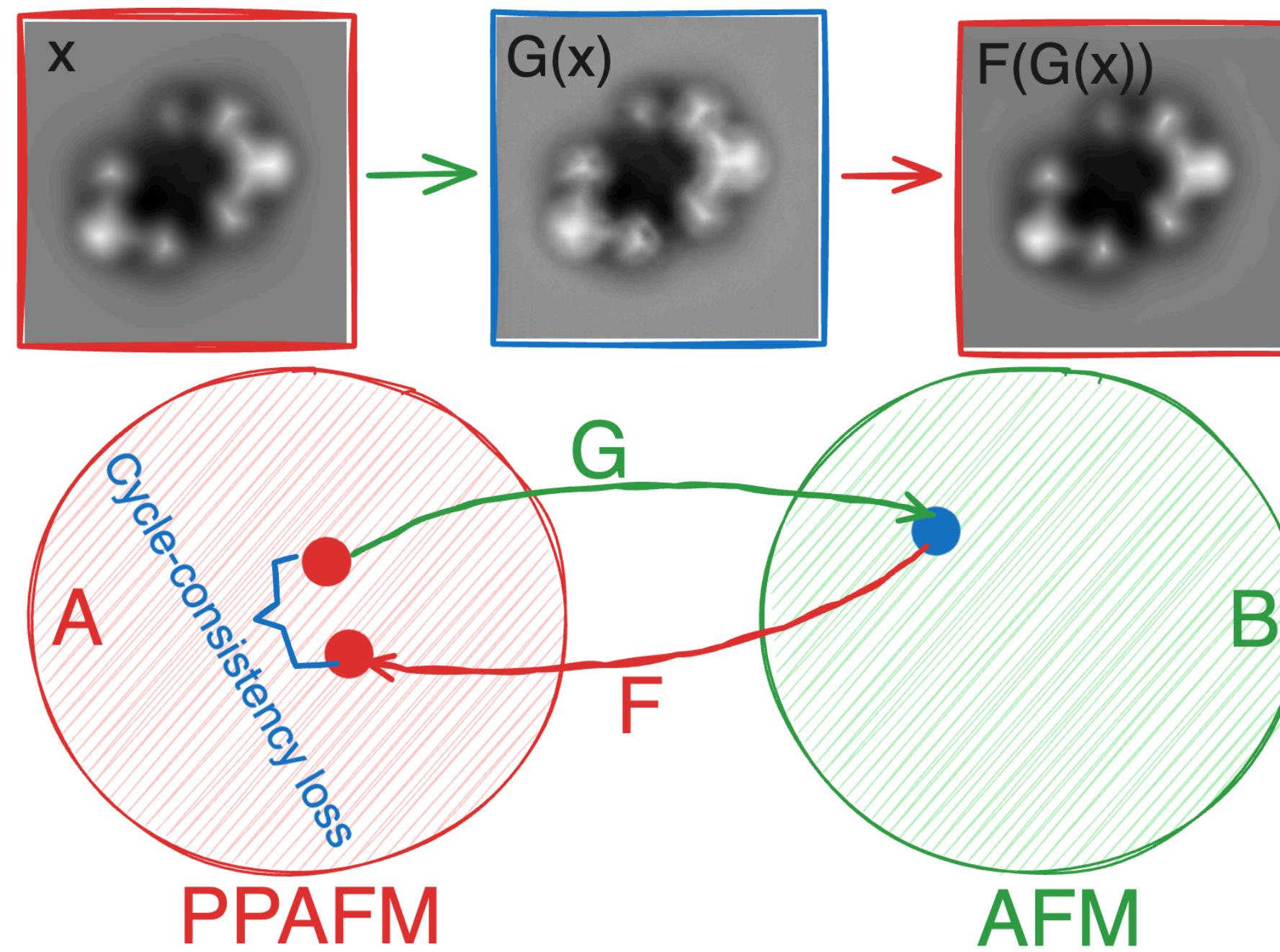
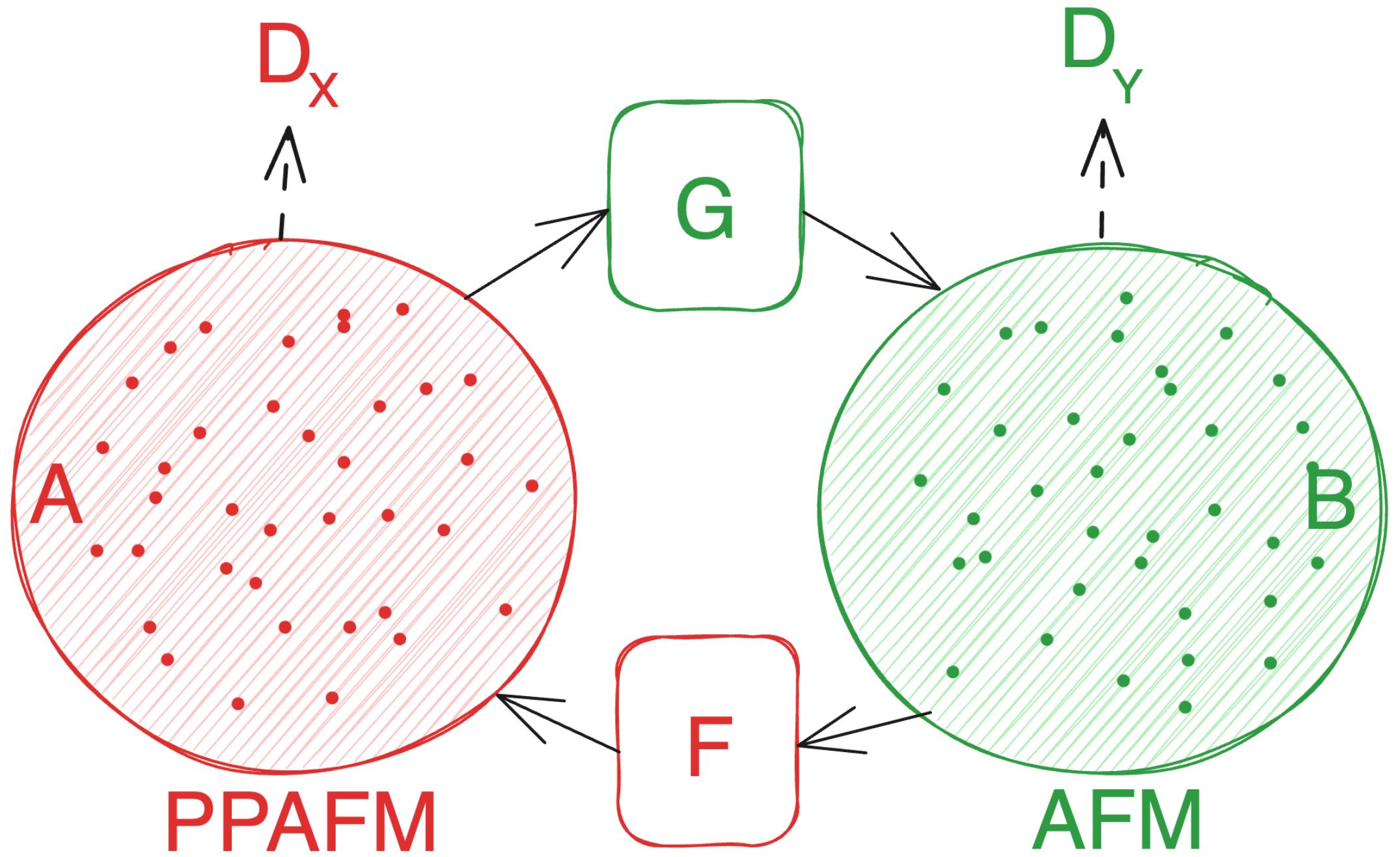


# How to generate fake AFM?

CycleGAN



# Style translation between PPAFM & AFM



$$\mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) = \mathbb{E}_{y \sim p_{\text{data}}(y)} [\log D_Y(y)] + \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log (1 - D_Y(G(x)))]$$

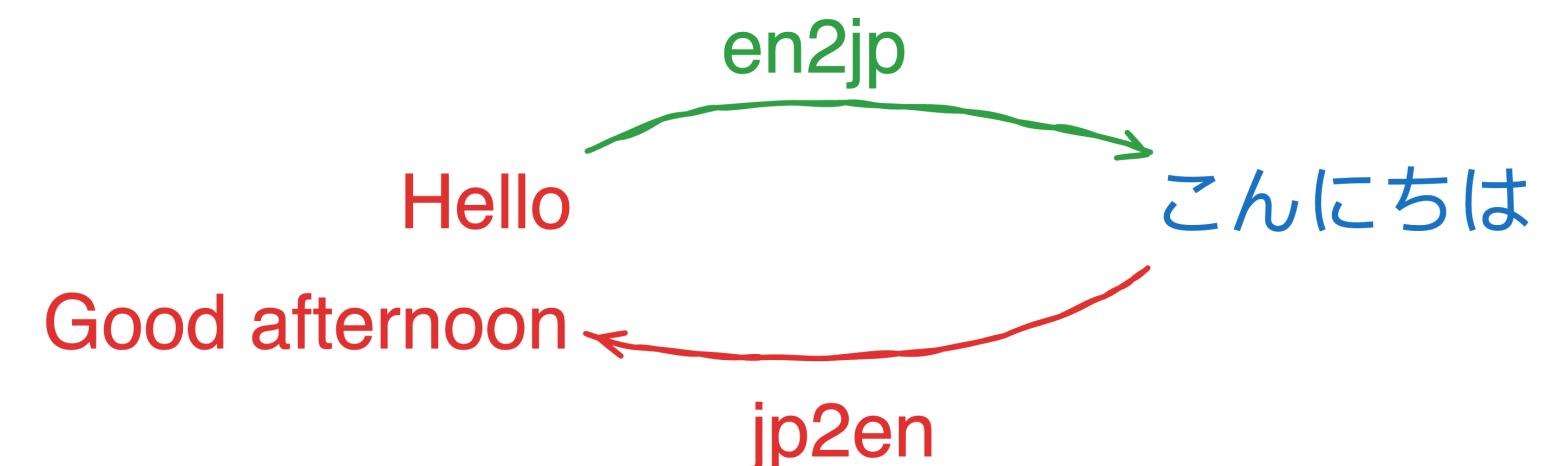
$$G^* = \arg \min_G \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y)$$

$$D_Y^* = \arg \max_D \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y)$$

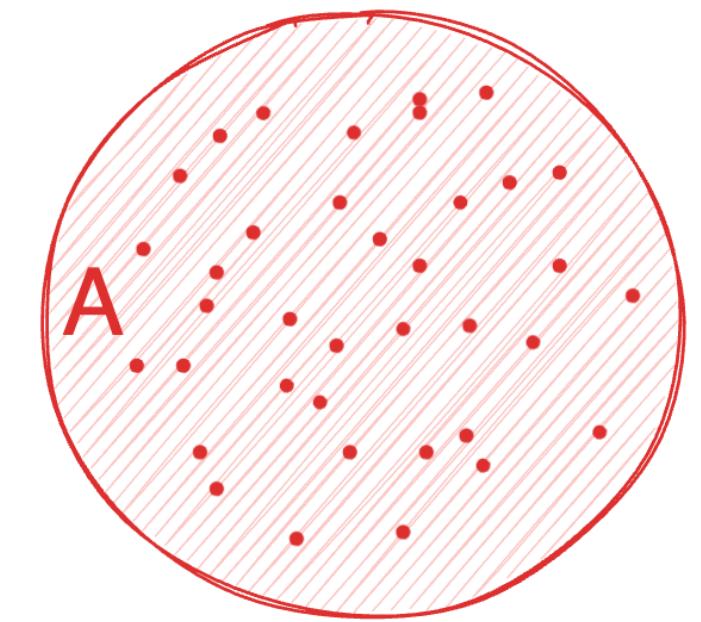
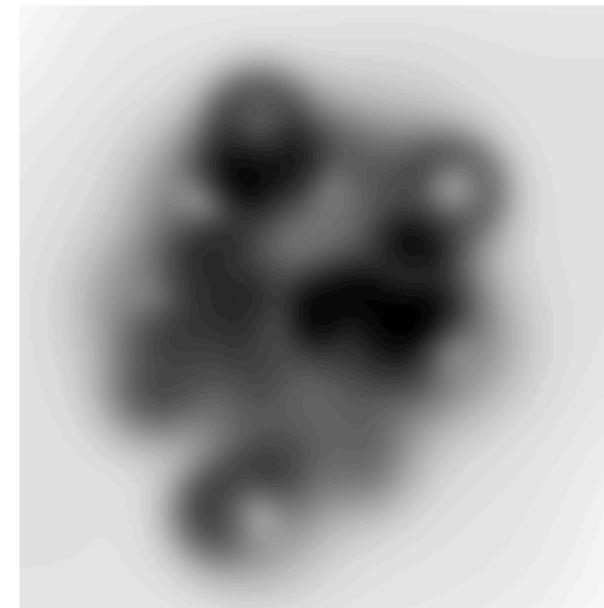
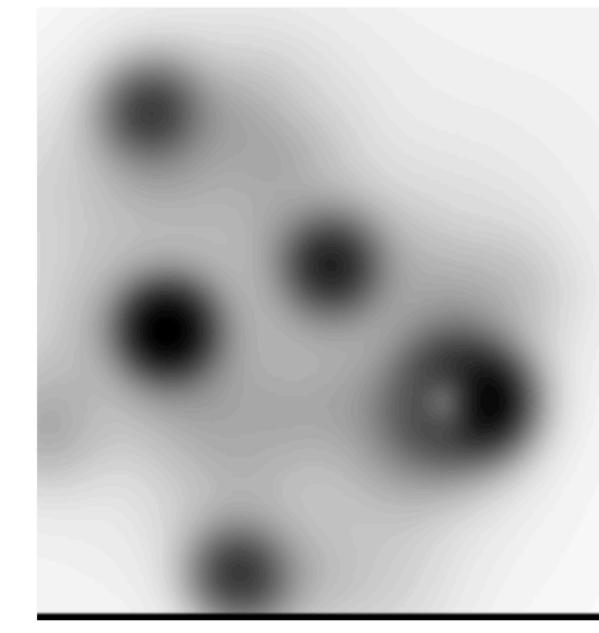
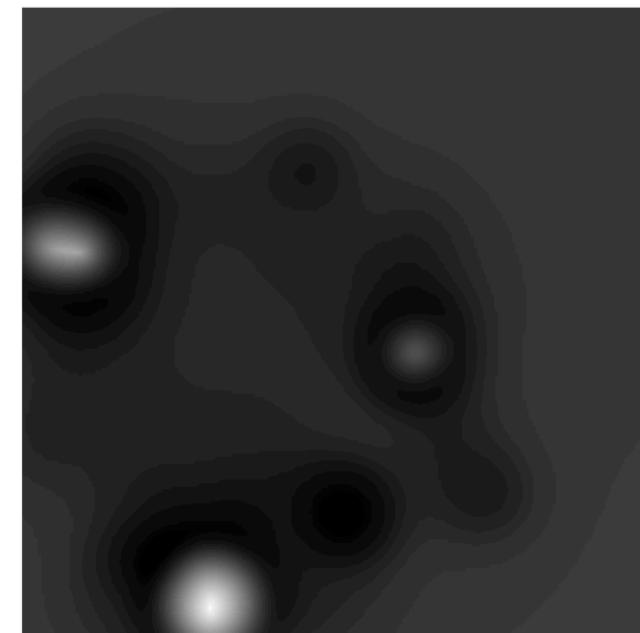
$$G^*, D_Y^* = \arg \min_G \max_D \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y)$$

$$\mathcal{L}(G, F, D_X, D_Y) = \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) + \mathcal{L}_{\text{GAN}}(F, D_X, Y, X) + \lambda \mathcal{L}_{\text{cyc}}(G, F)$$

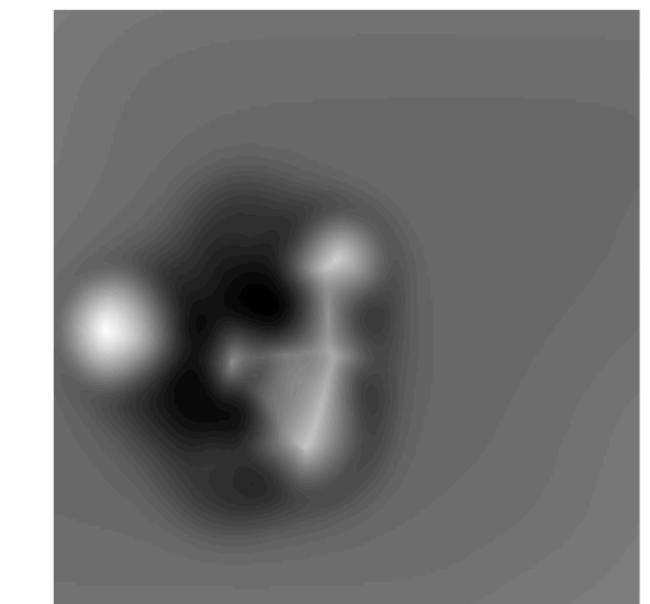
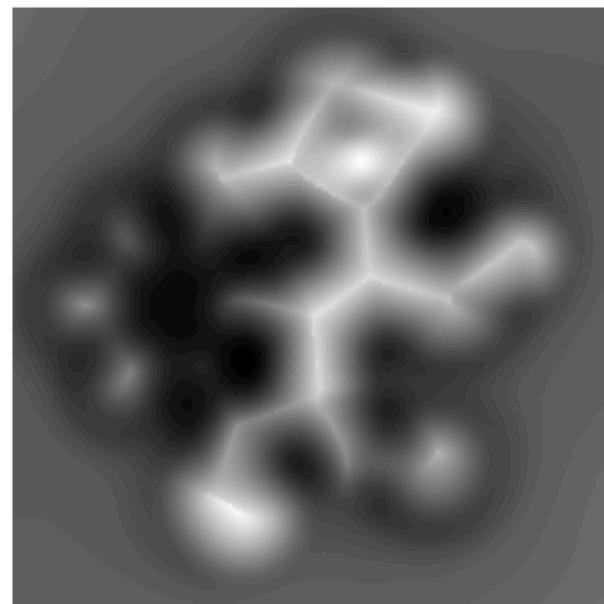
$$G^*, F^* = \arg \min_{G, F} \max_{D_X, D_Y} \mathcal{L}(G, F, D_X, D_Y)$$



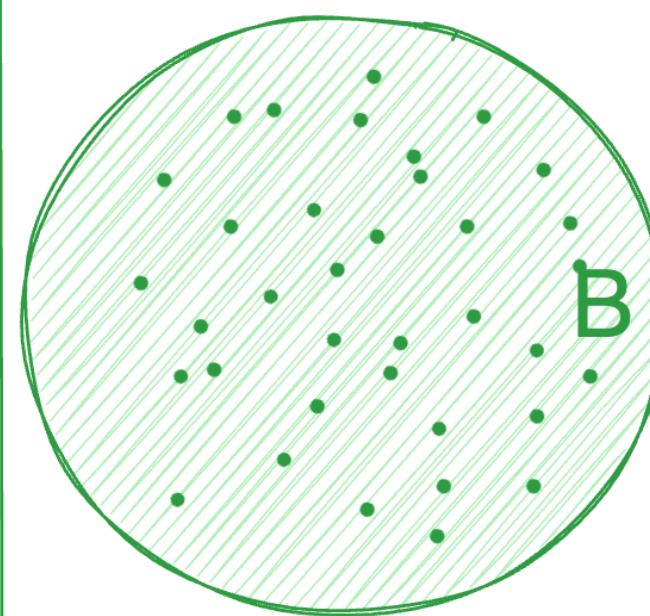
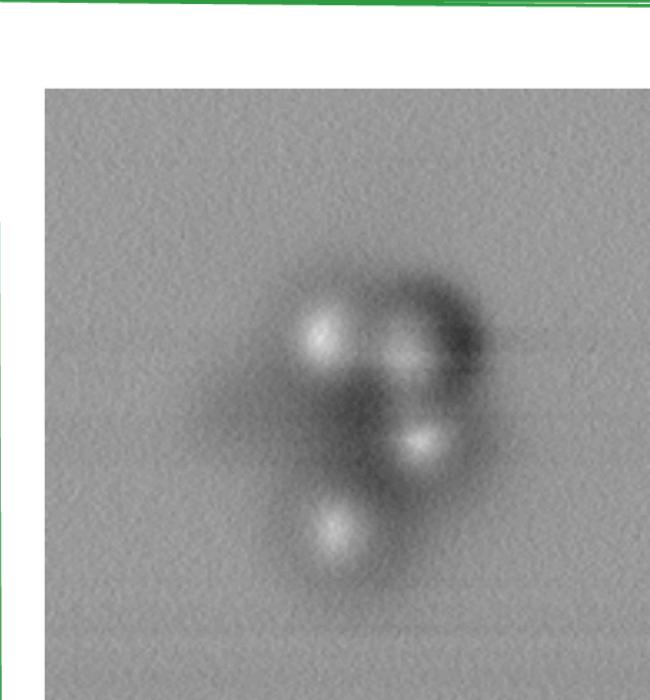
# Training samples



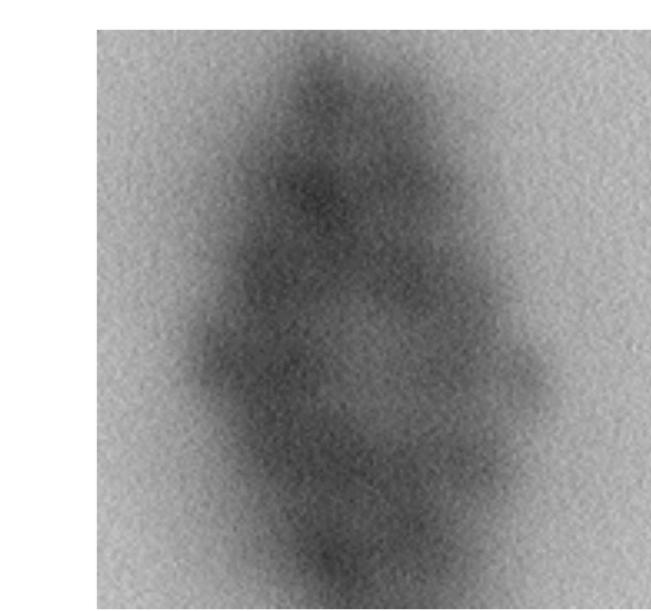
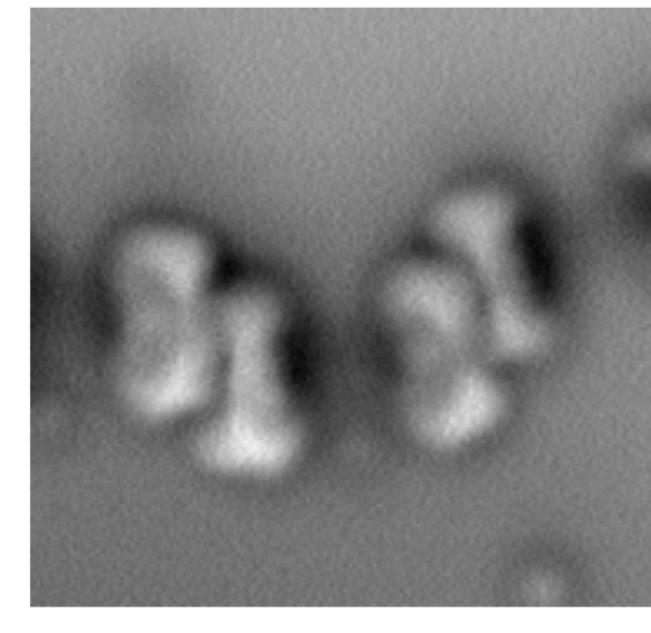
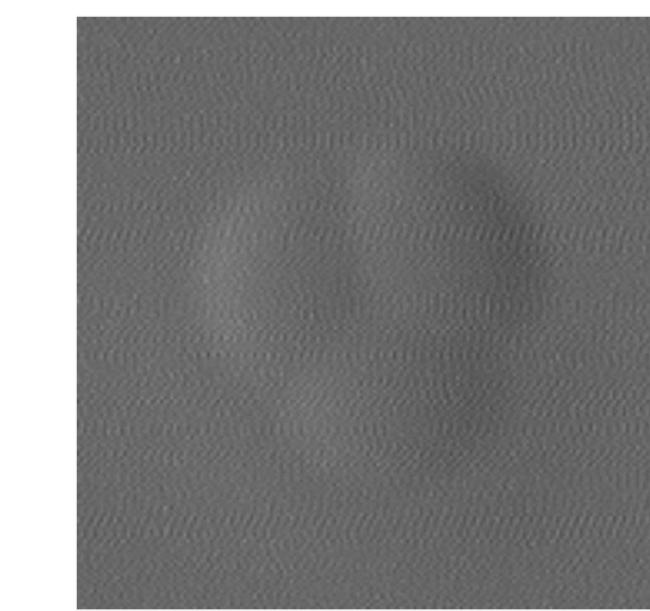
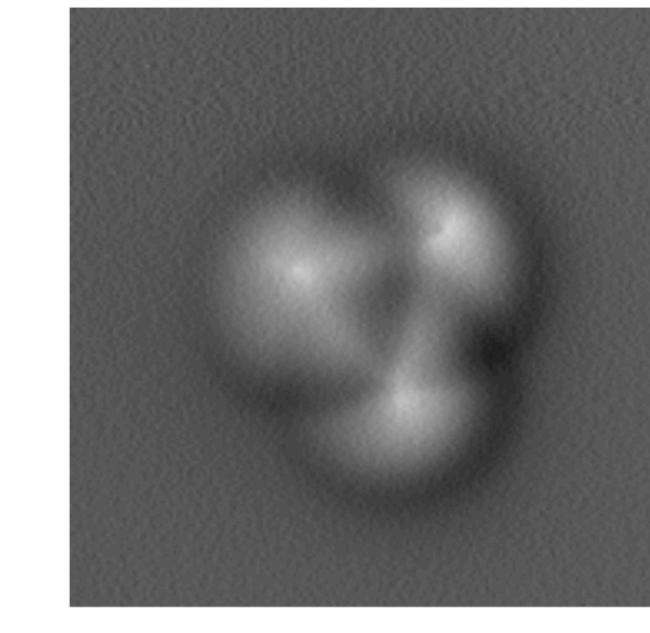
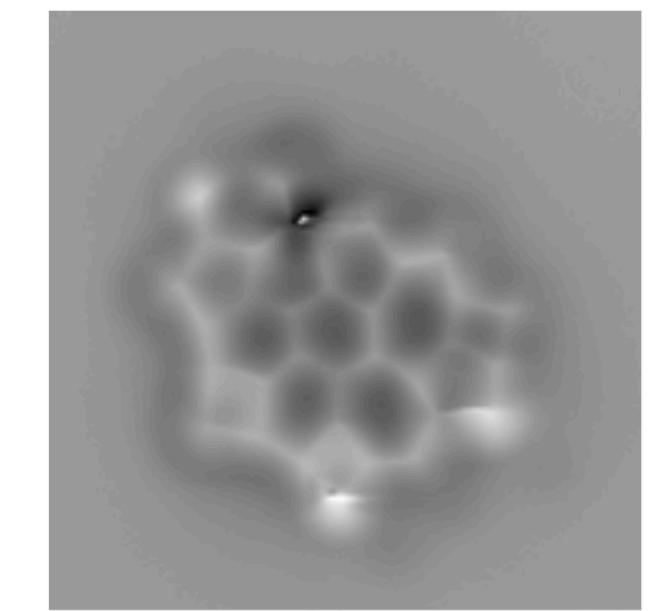
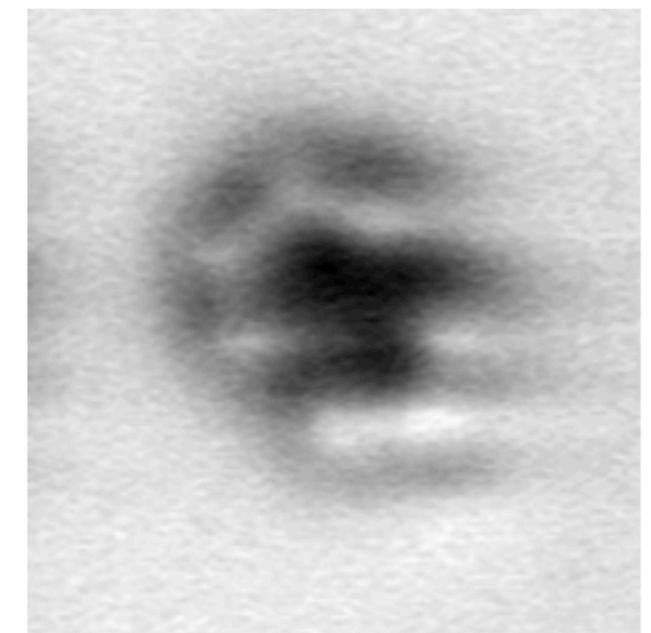
PPAFM Monolayer Water Molecules  
609



Size 192x192, 8-bit grayscale



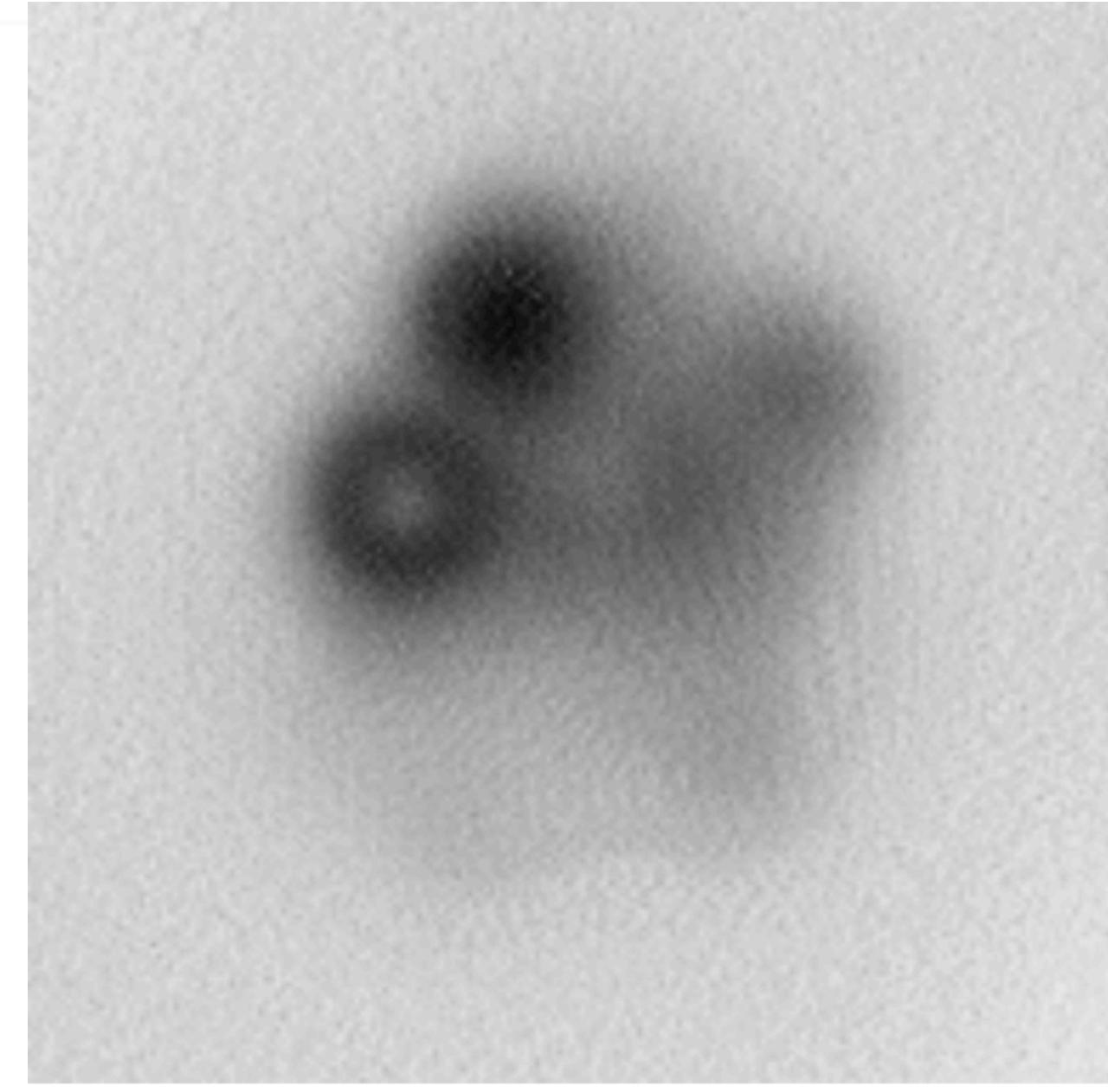
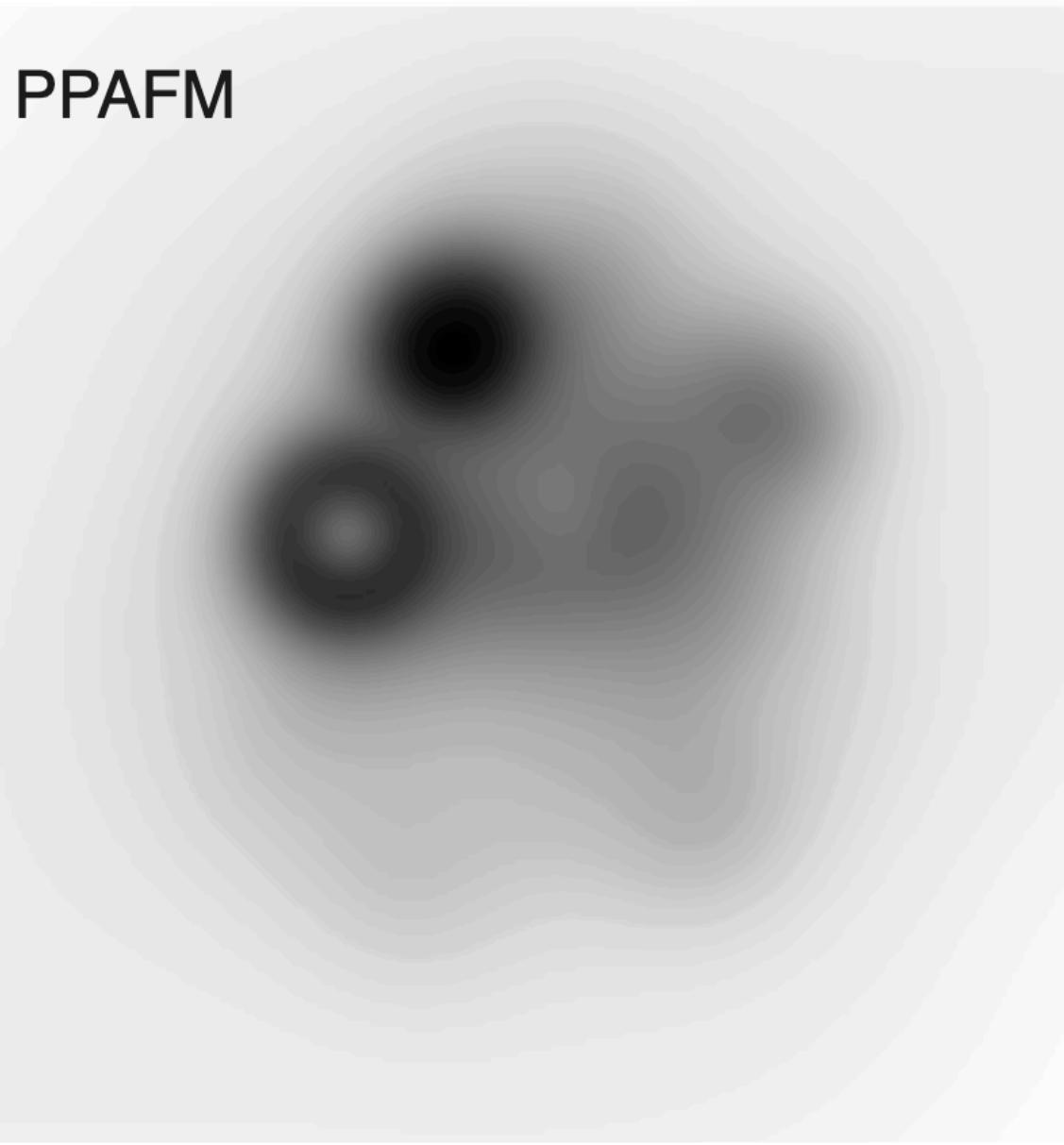
CO tip AFM available  
510



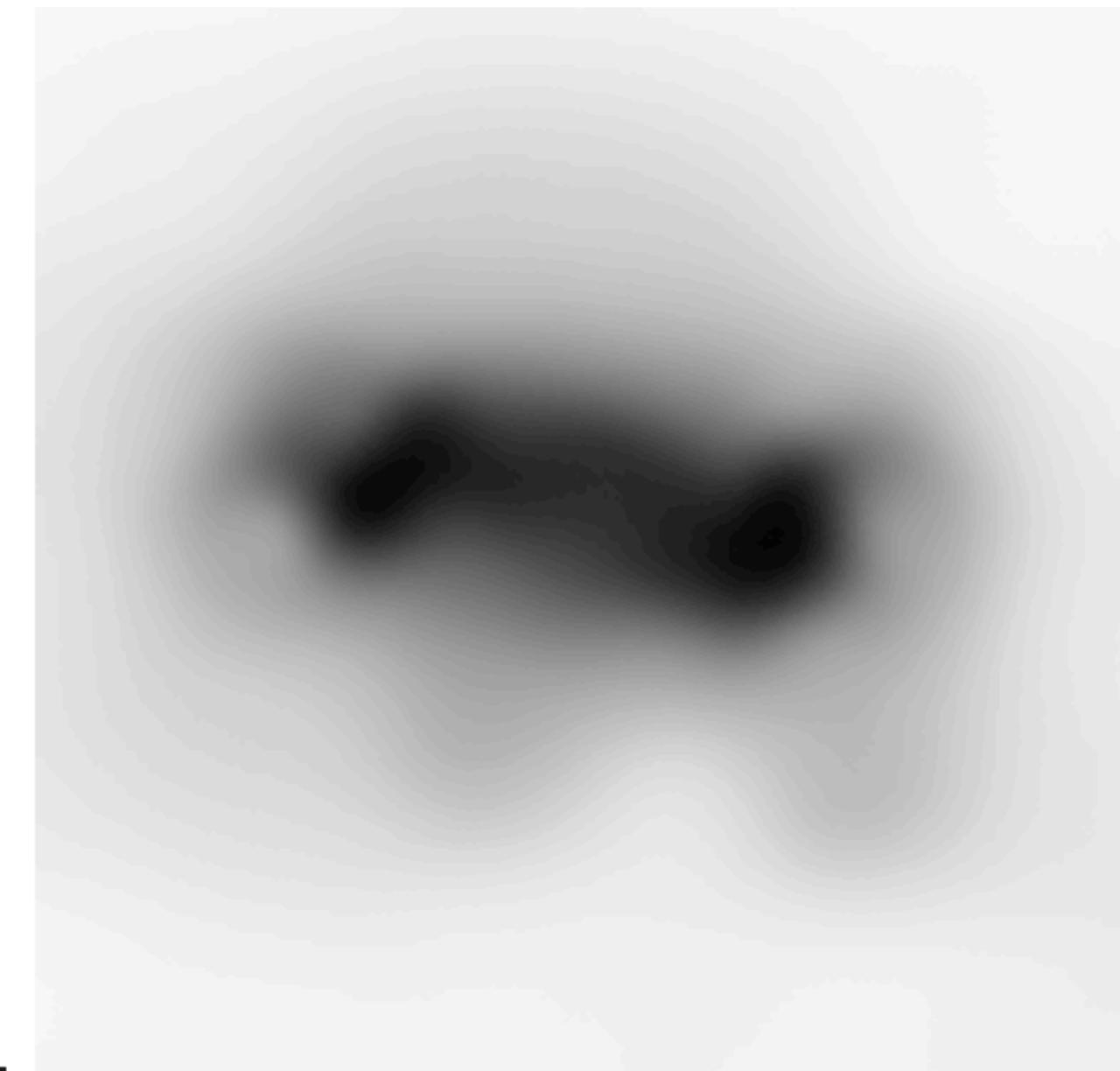
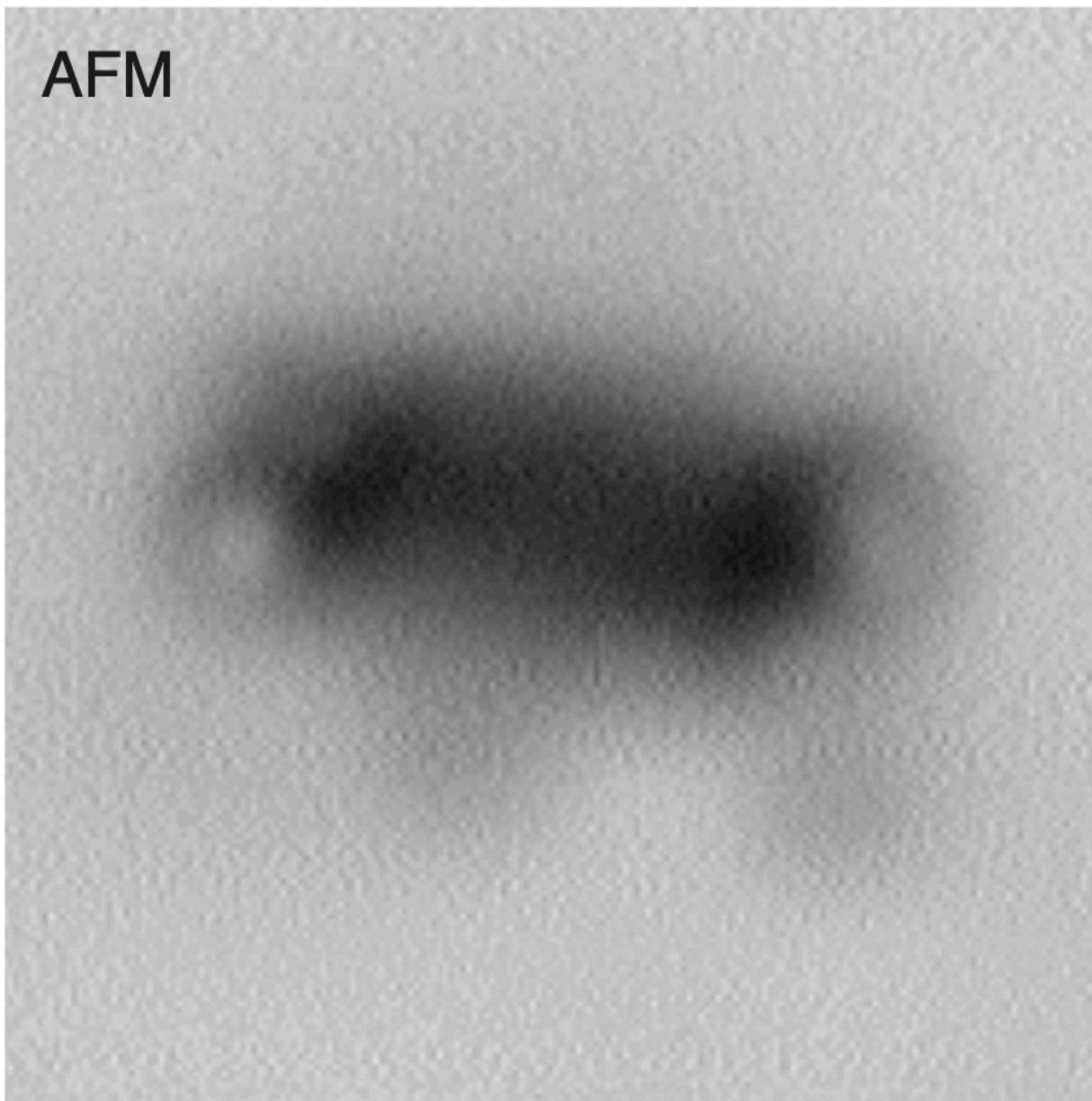
Size 192x192 (most), 8-bit grayscale

# Training Results

PPAFM  AFM



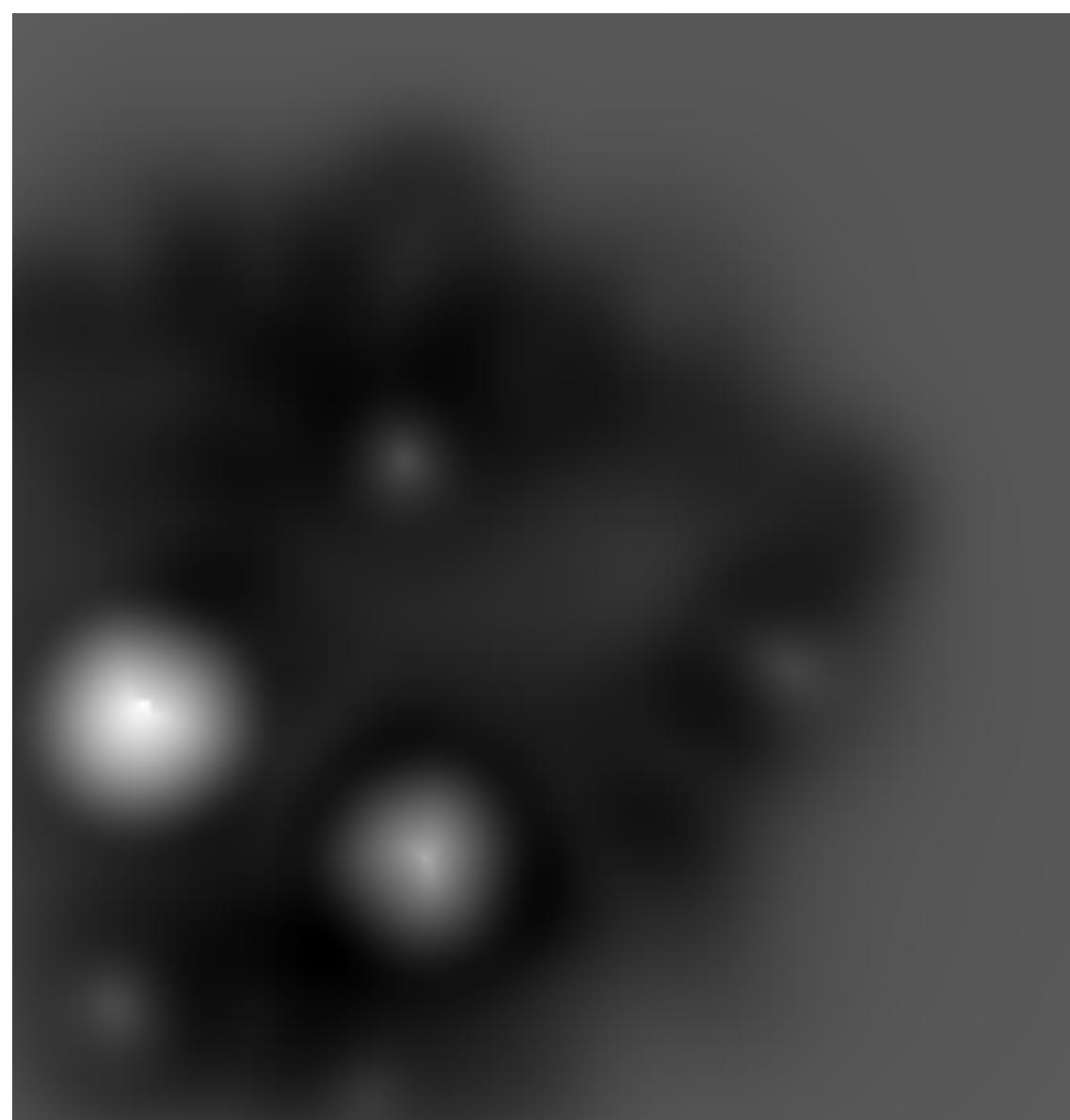
PPAFM → AFM



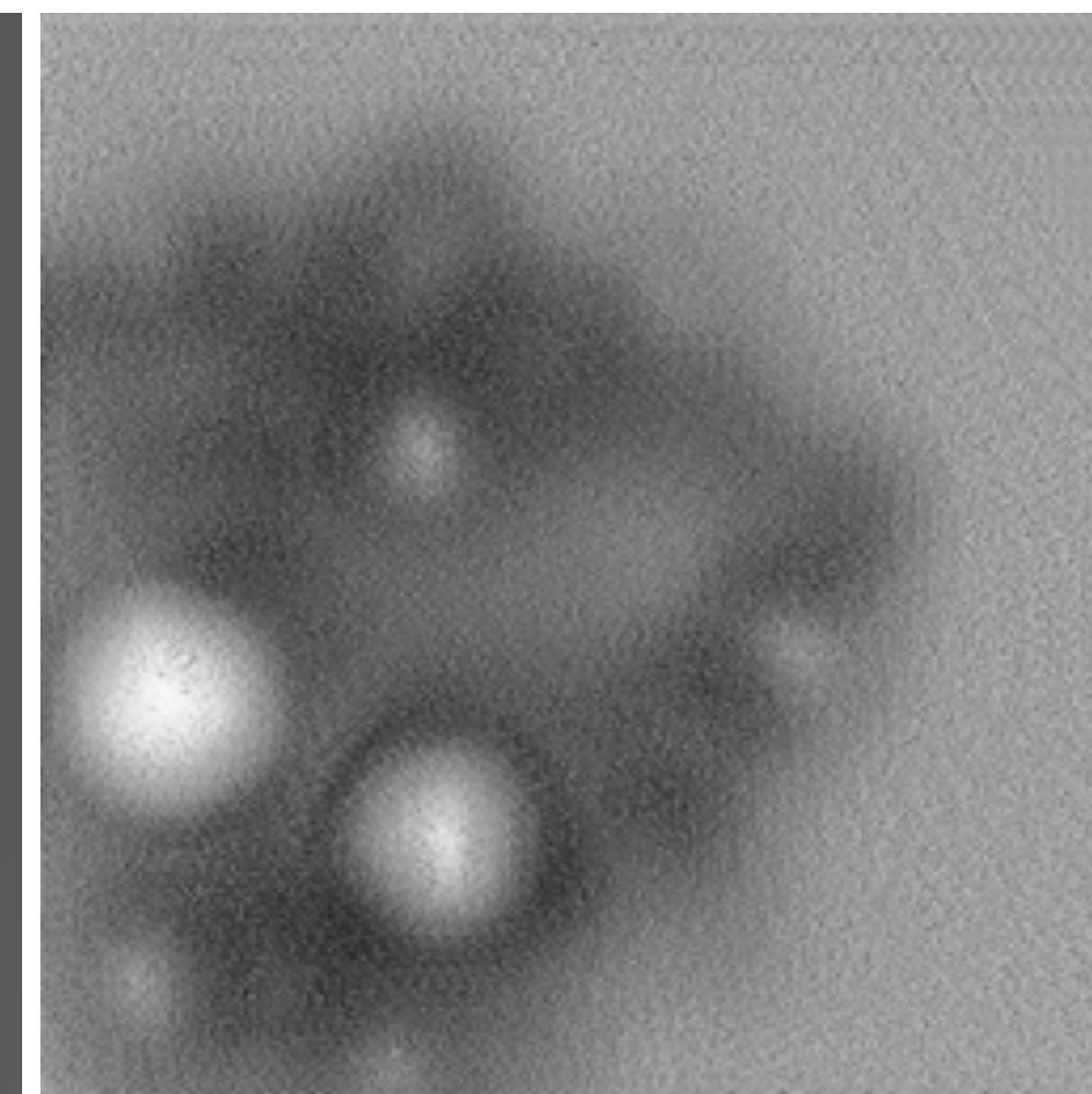
AFM → PPAFM

# More Translation Examples: A2B

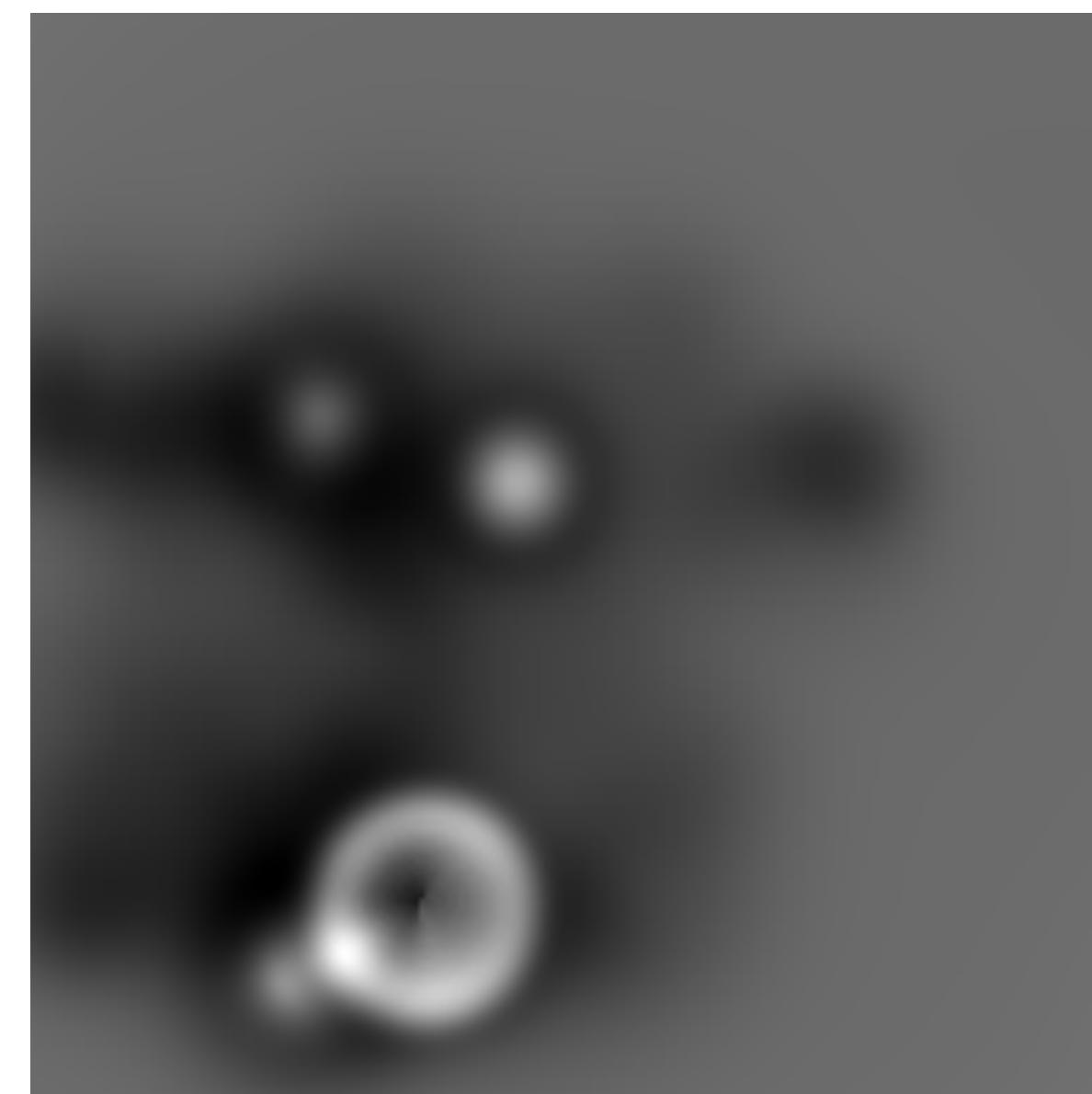
Simulation



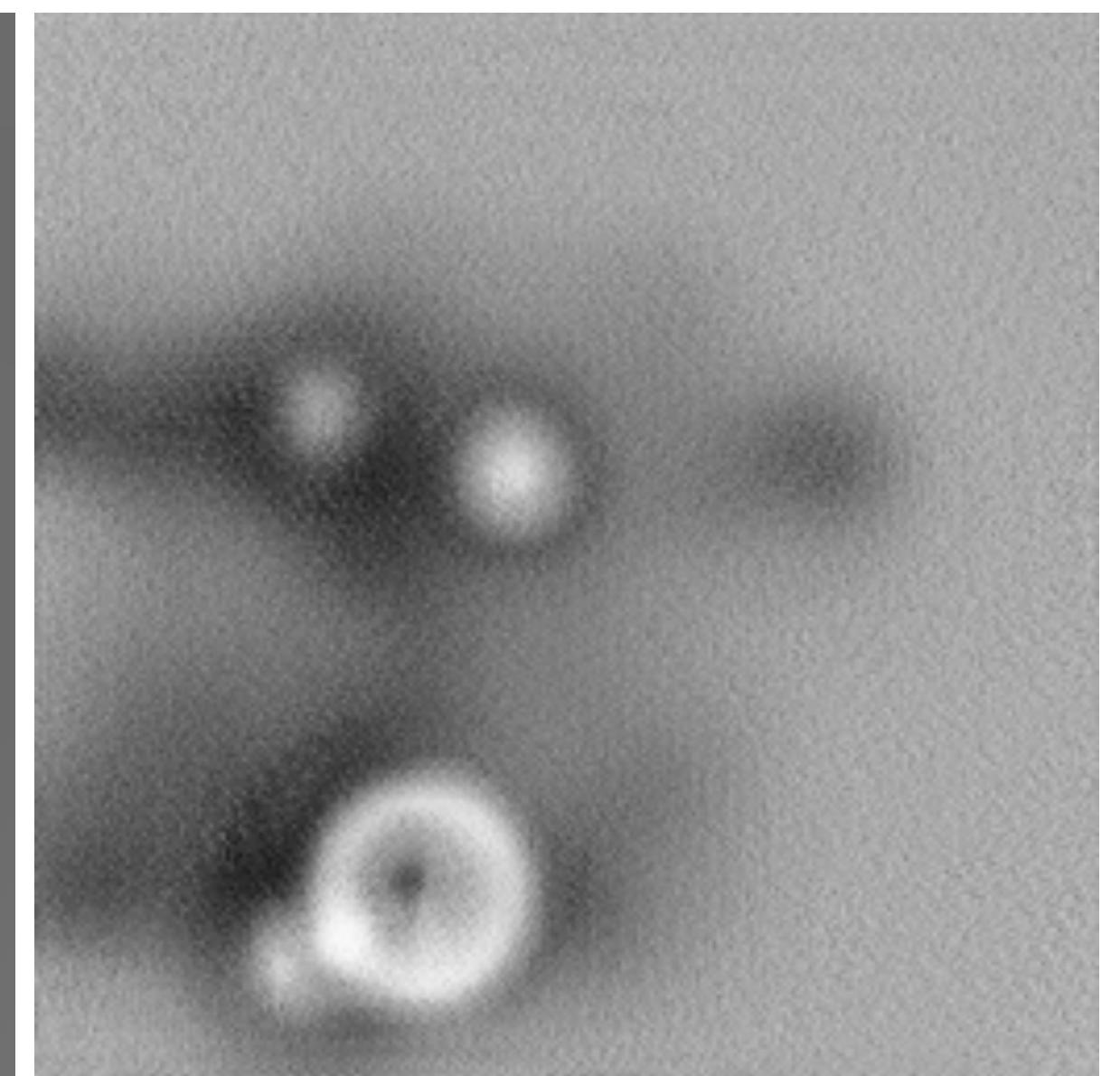
Experimental-style



Simulation

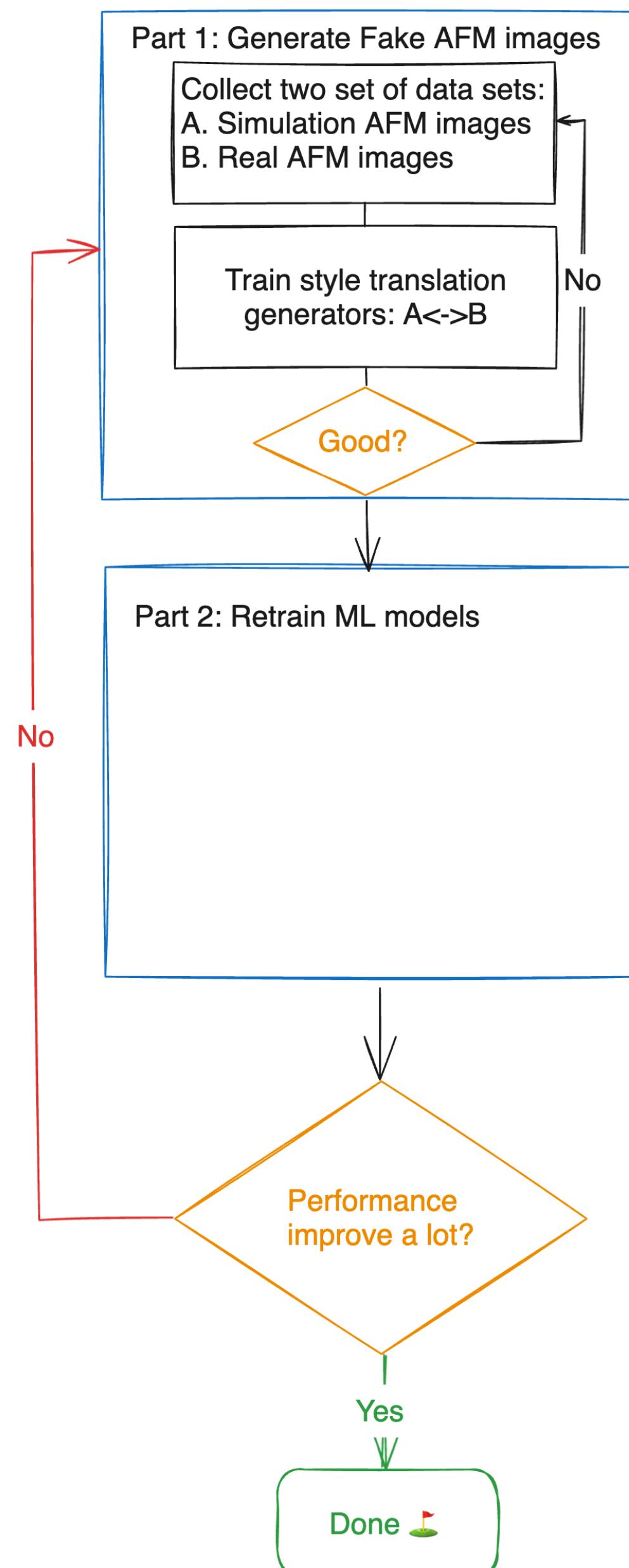
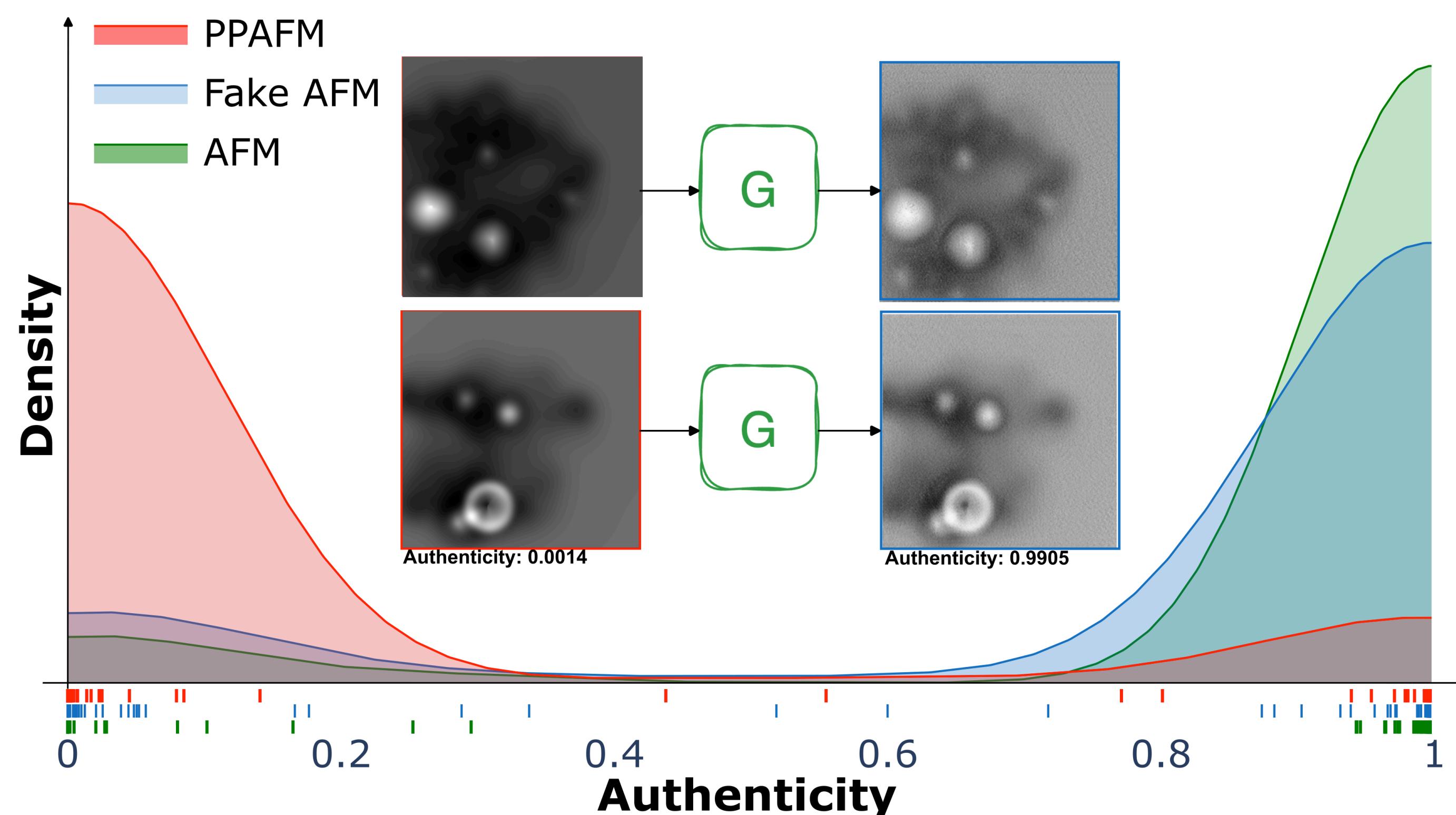
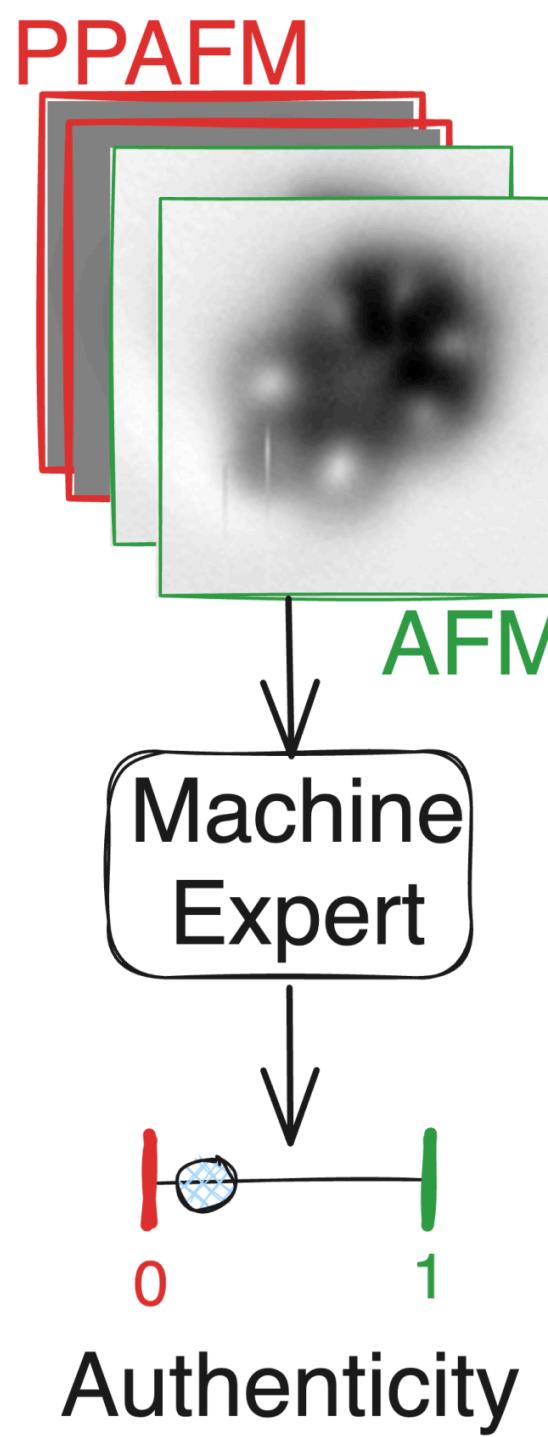


Experimental-style



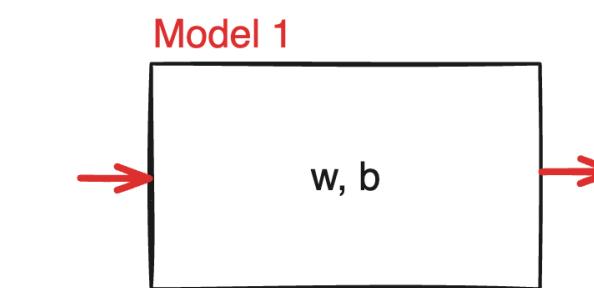
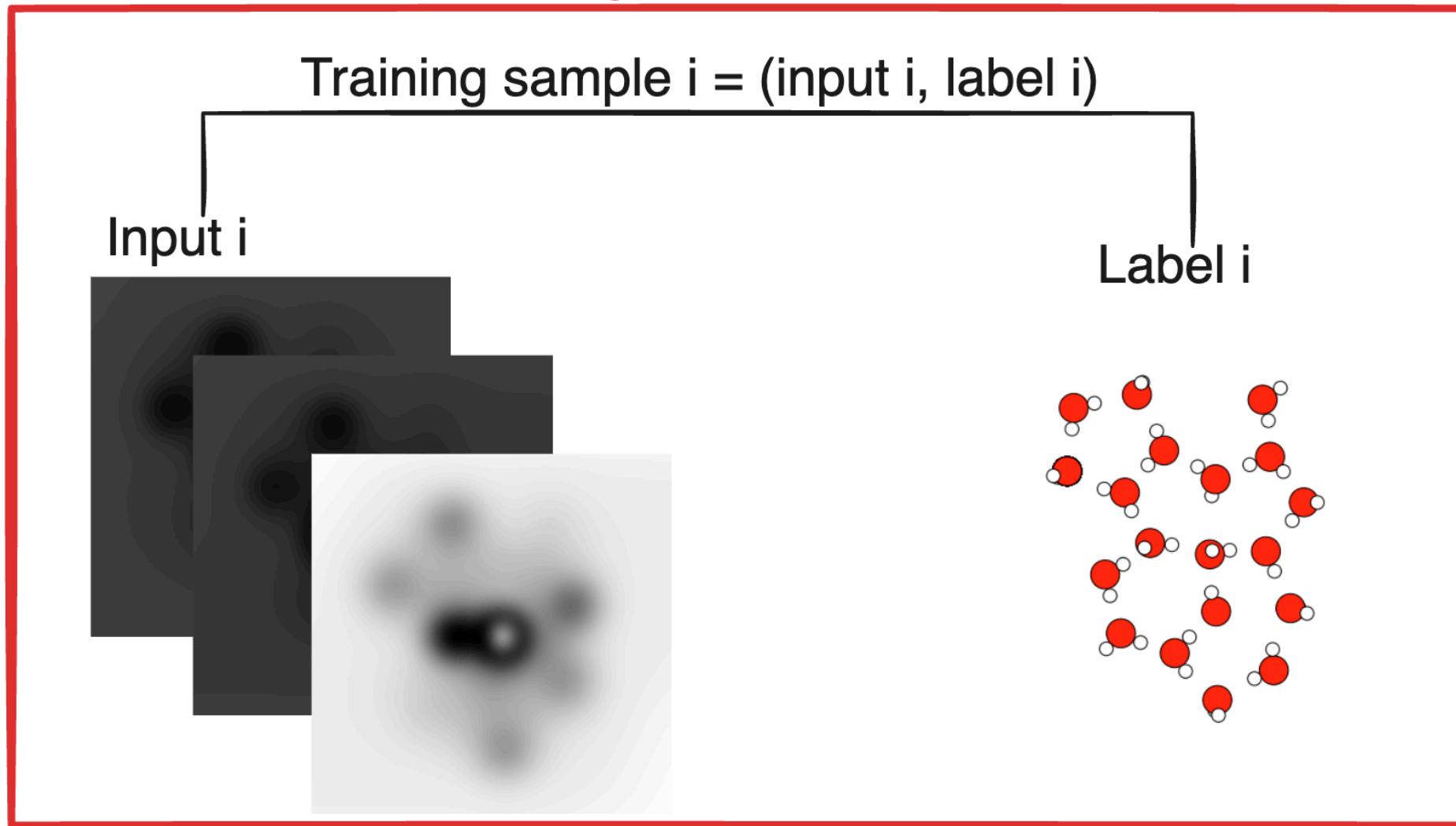
# Dataset Evaluation

From the perspective of a well trained machine expert.

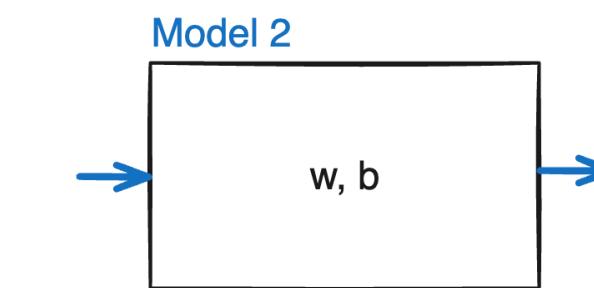
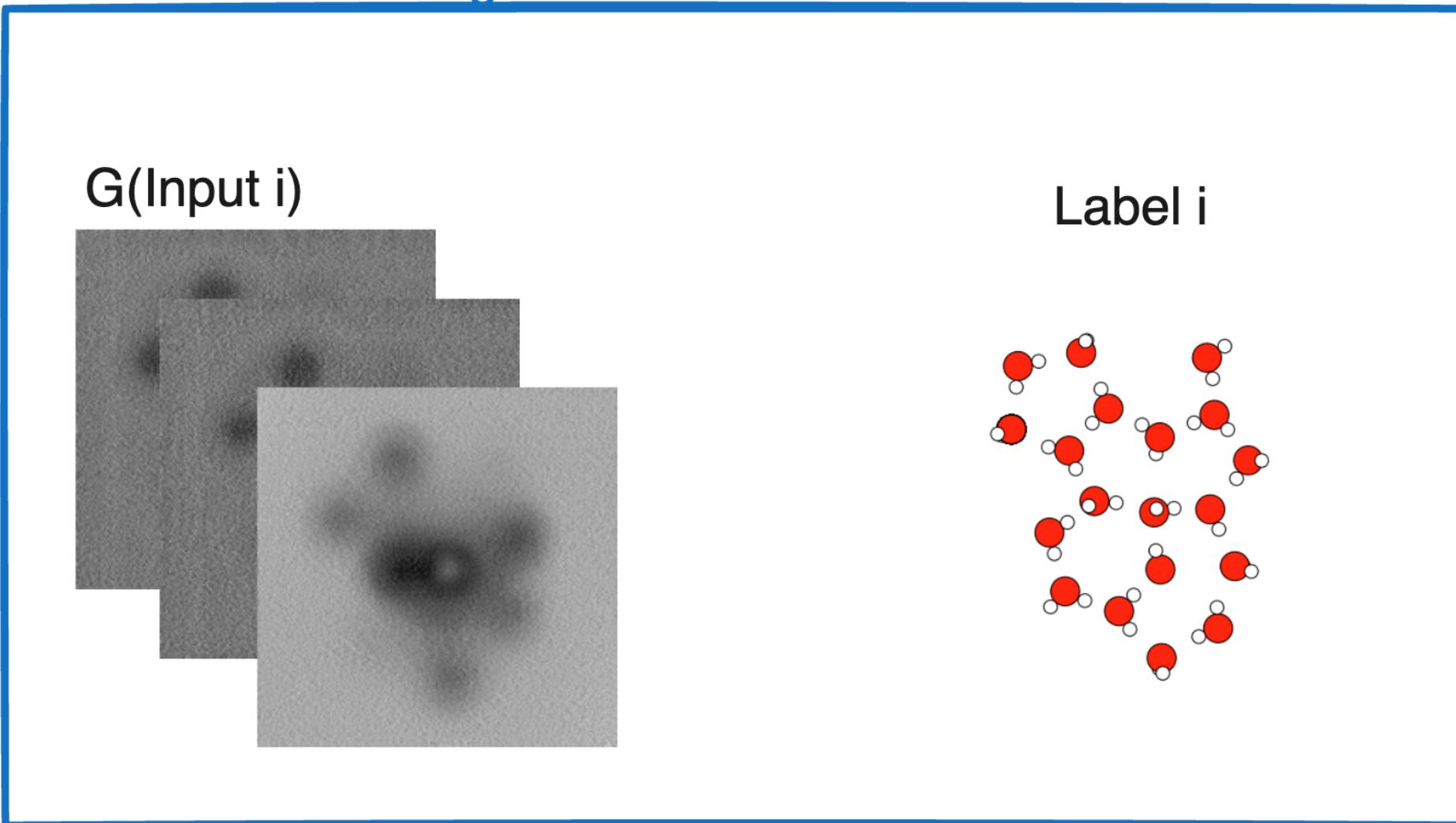


# Retrain structure discovery model with different datasets

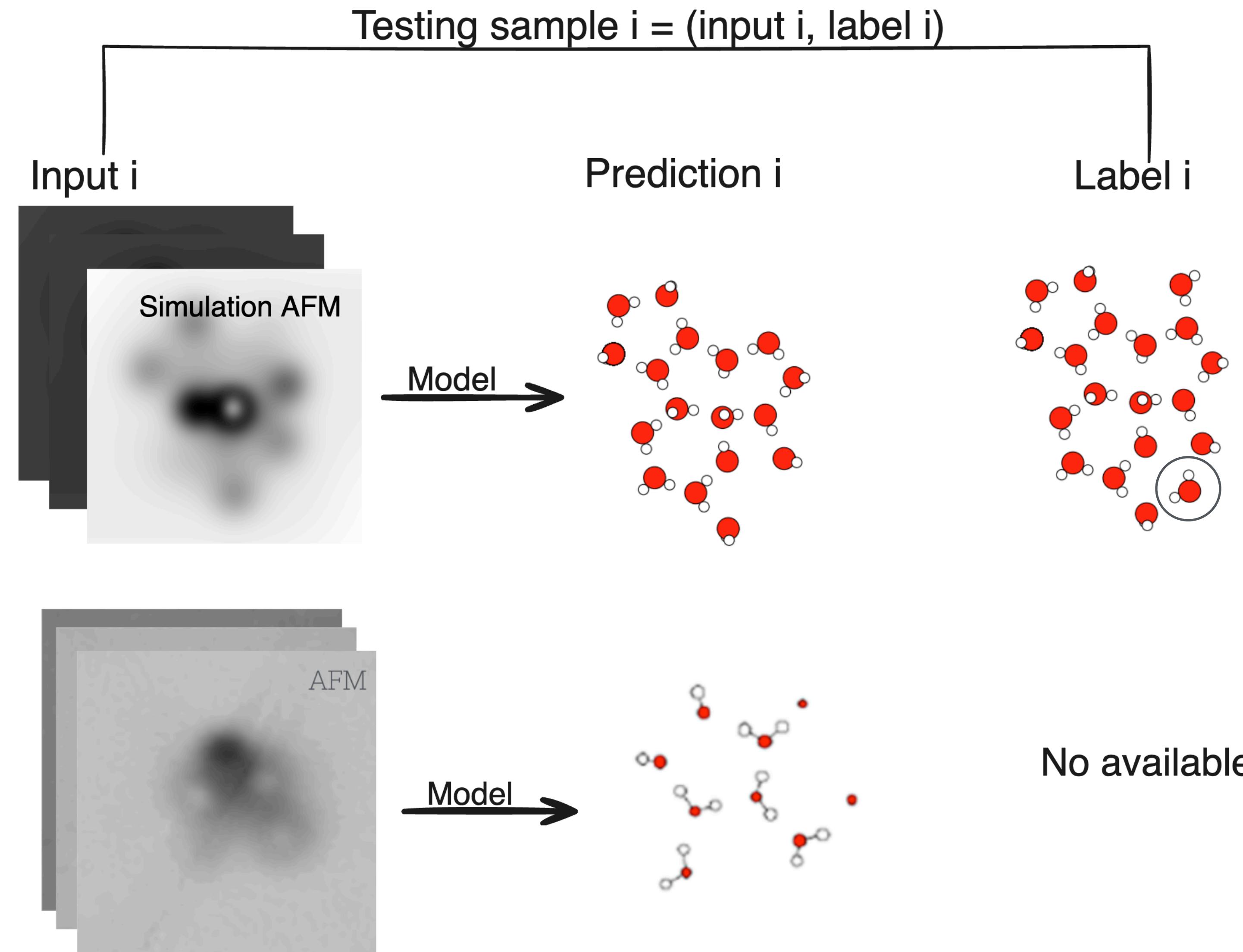
## 1. Simulation AFM training set



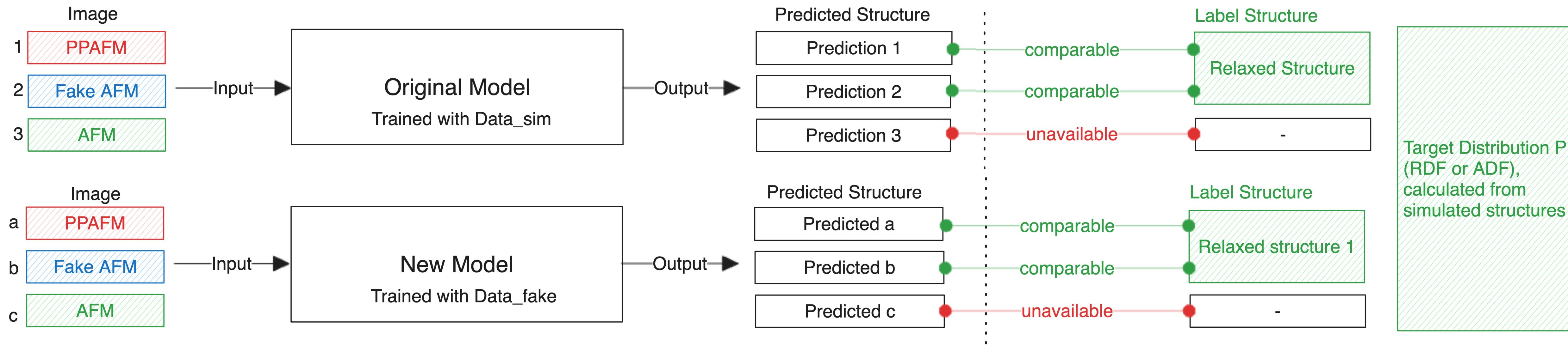
## 2. Fake AFM training set



# The challenge of performance evaluation



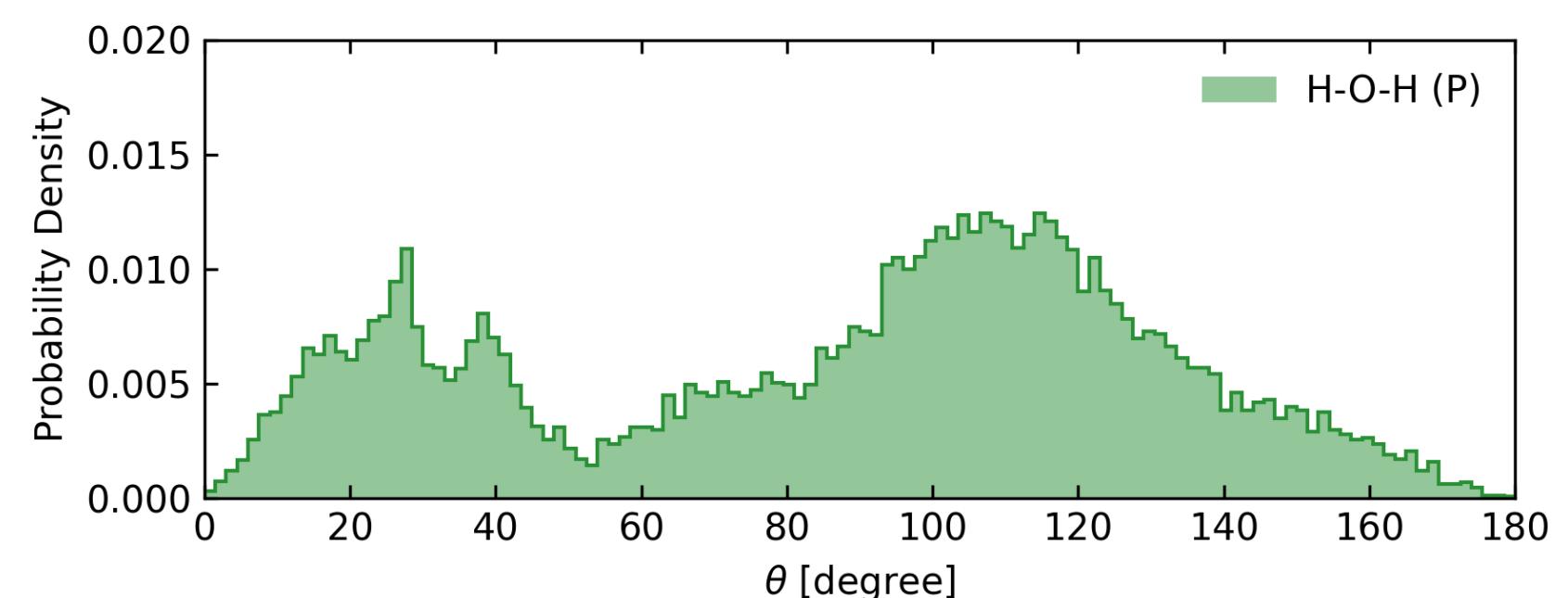
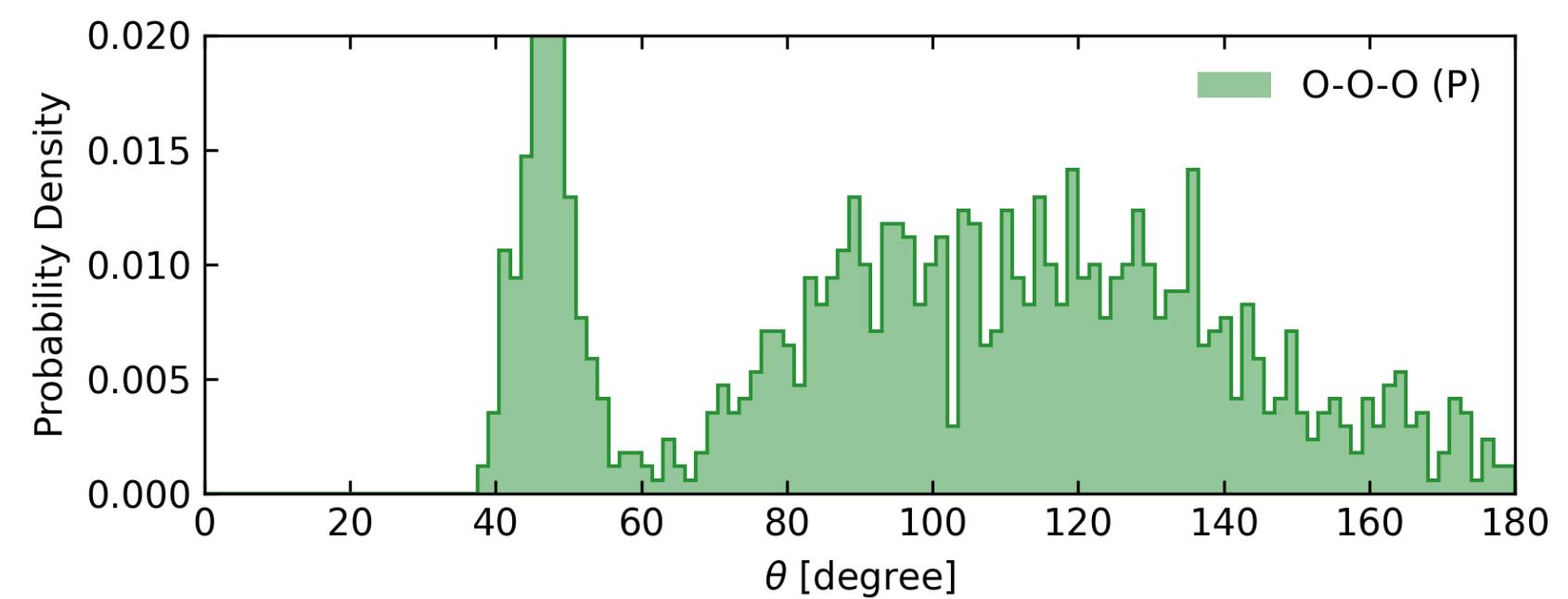
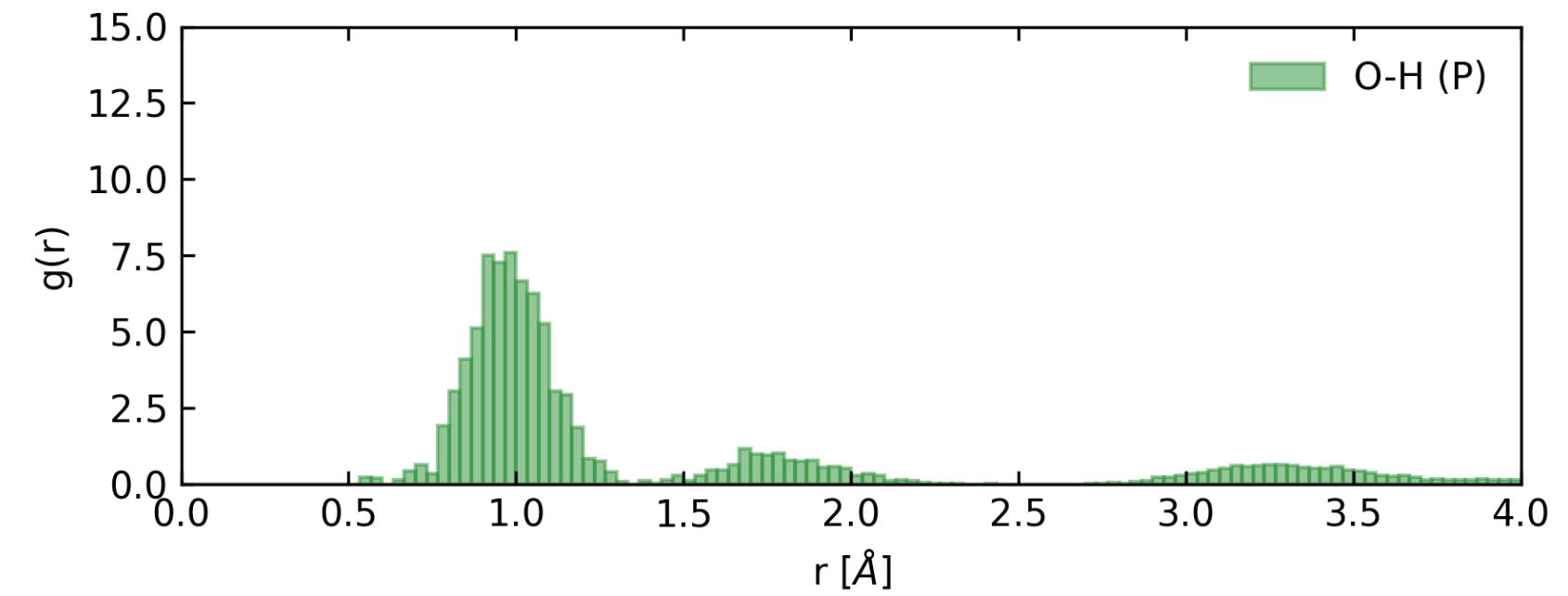
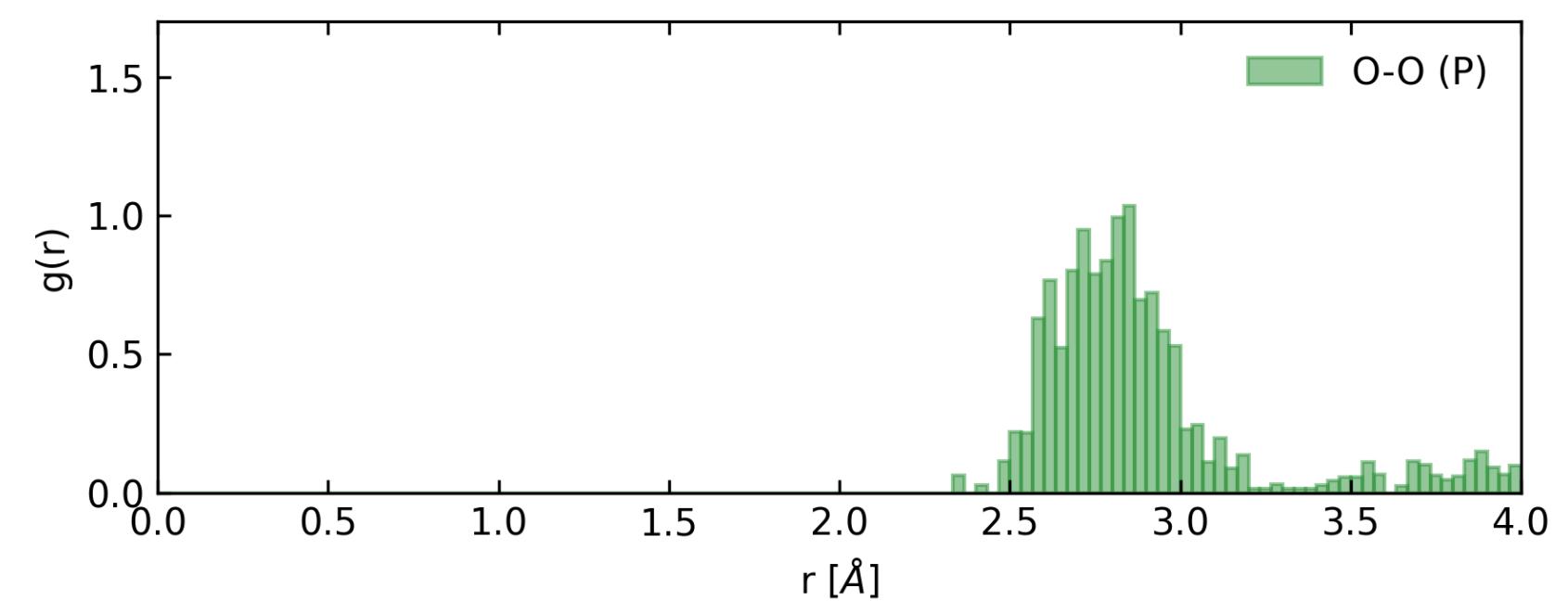
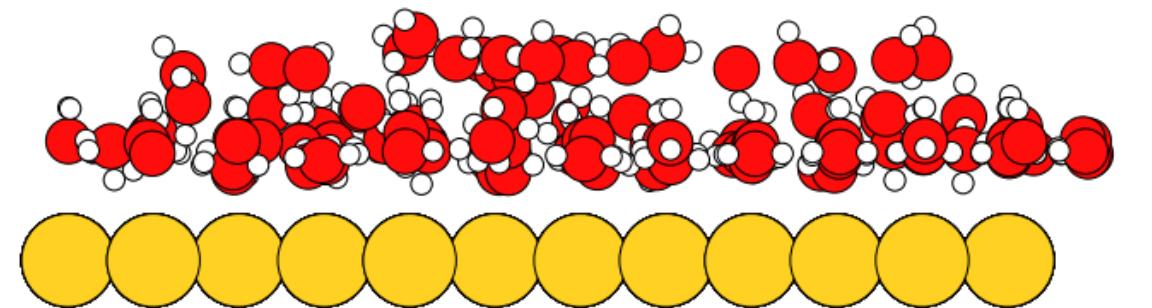
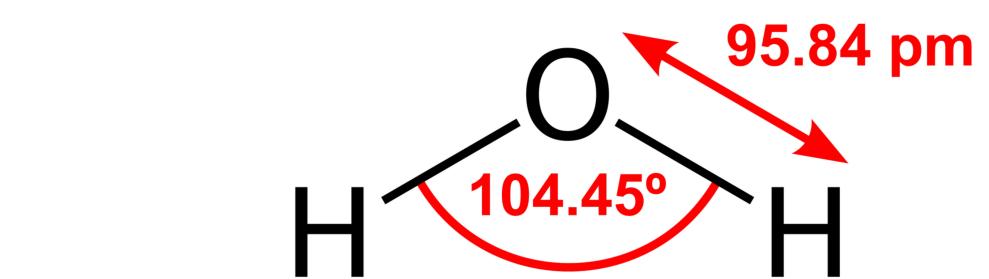
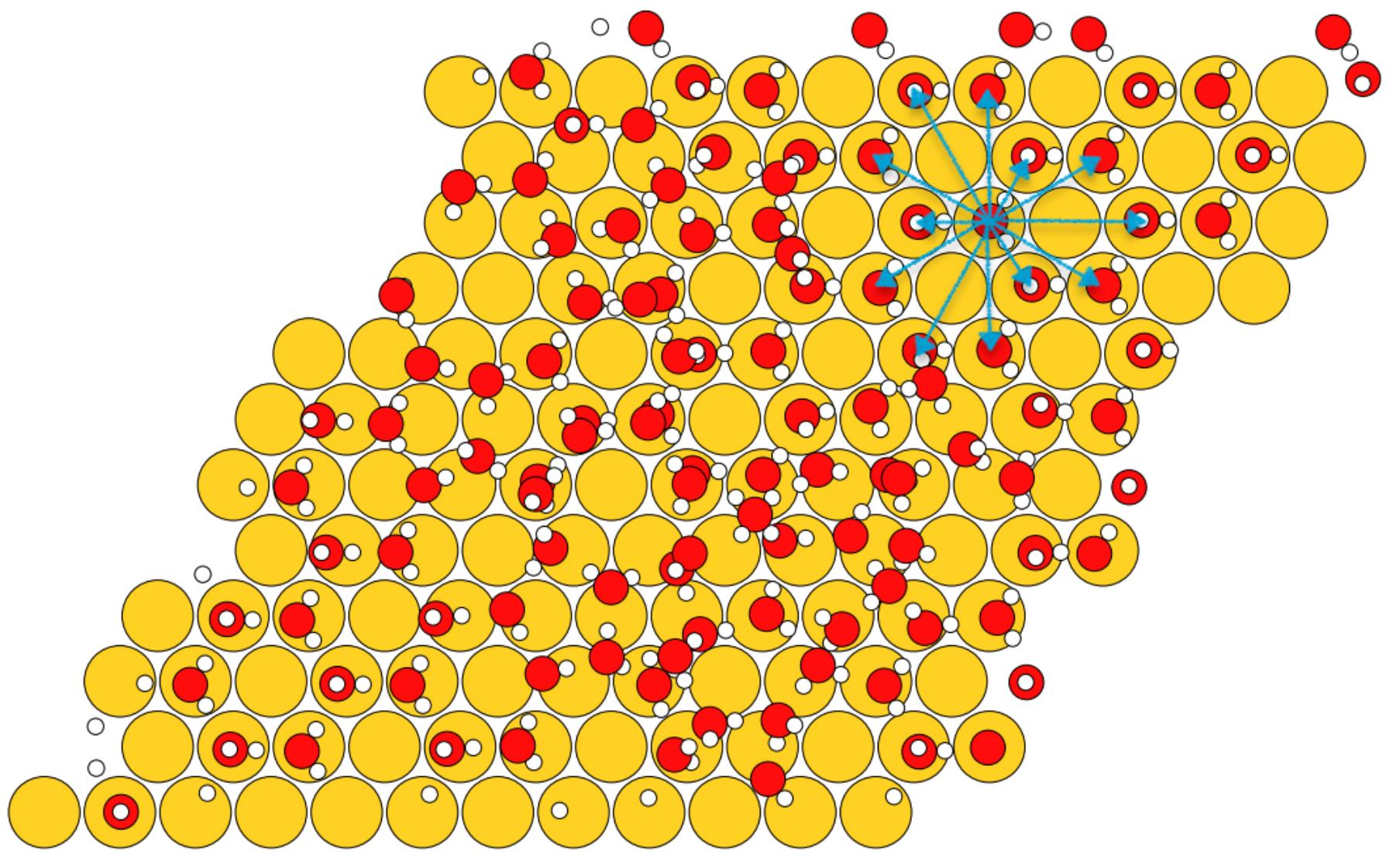
# Performance evaluation

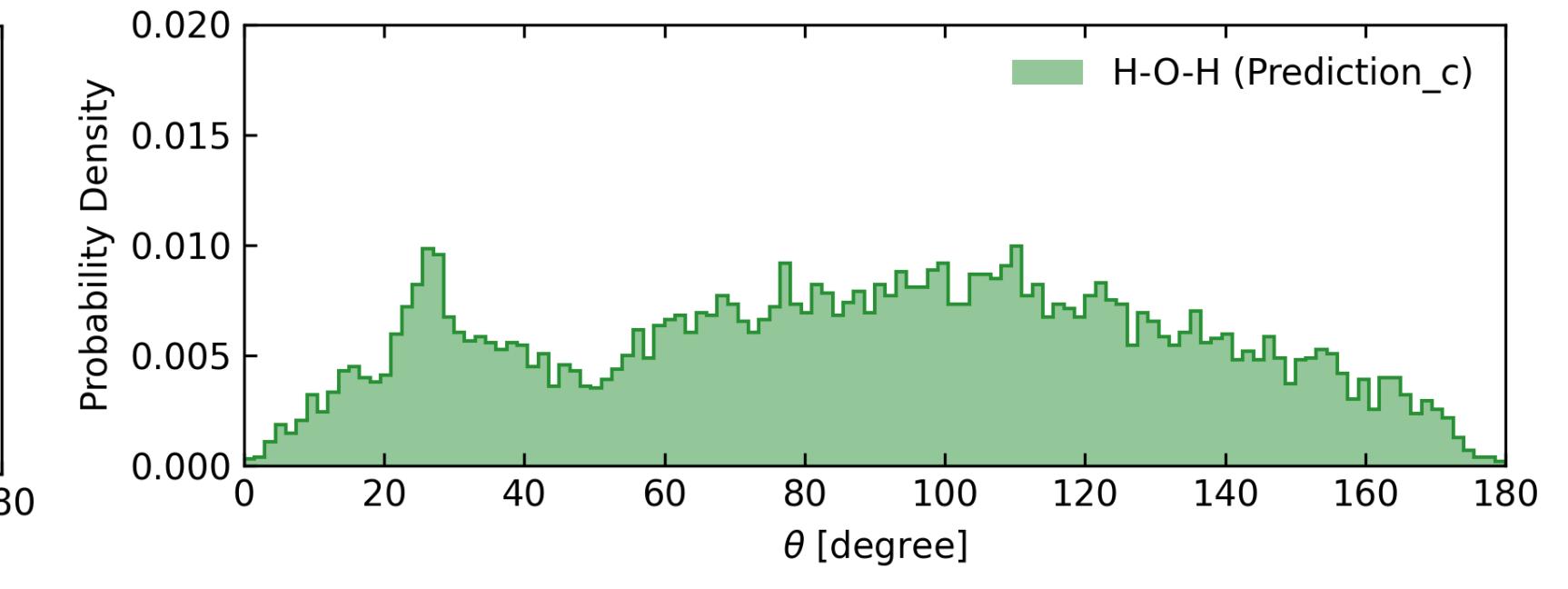
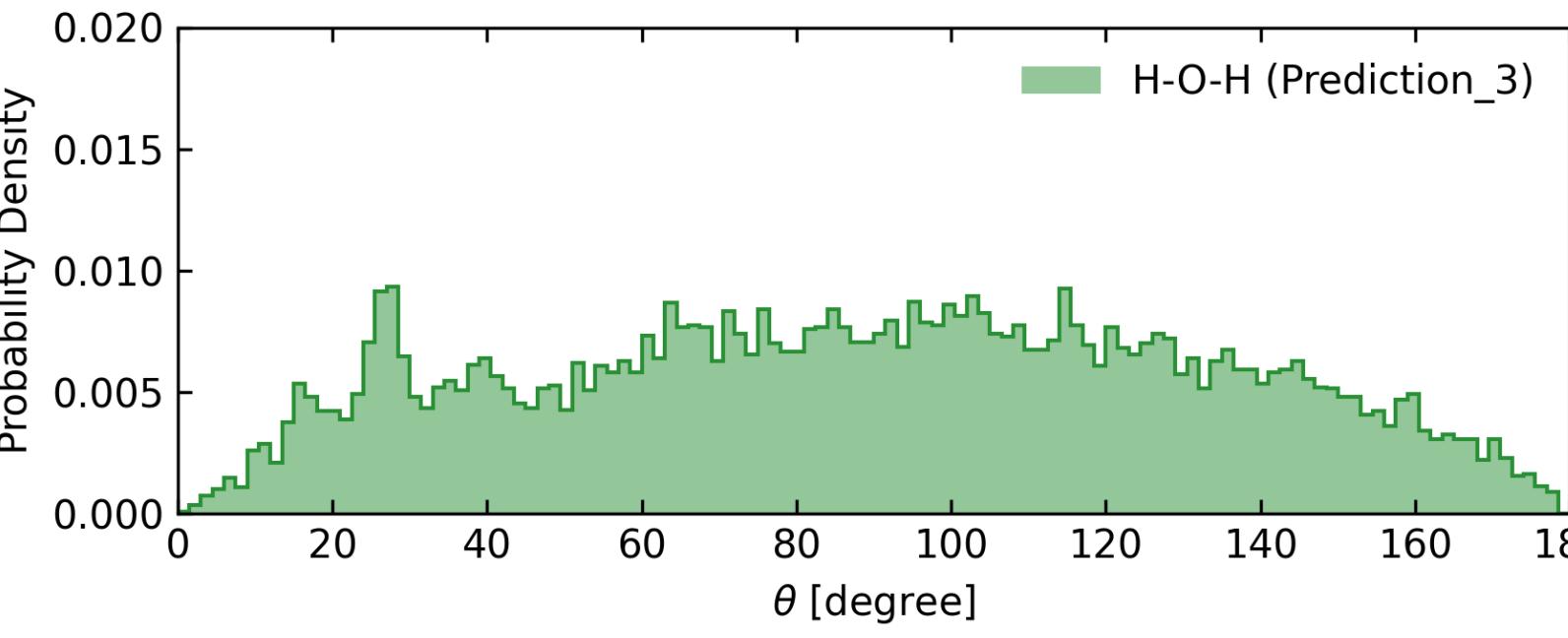
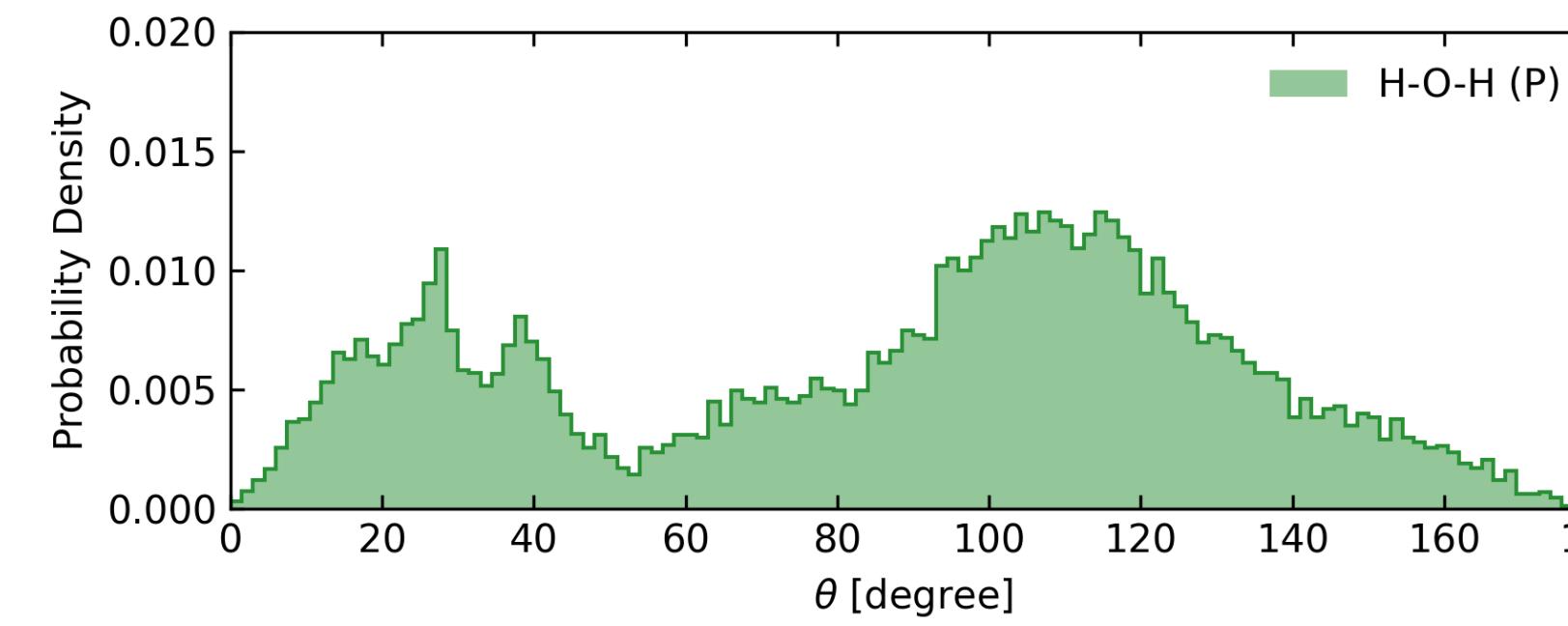
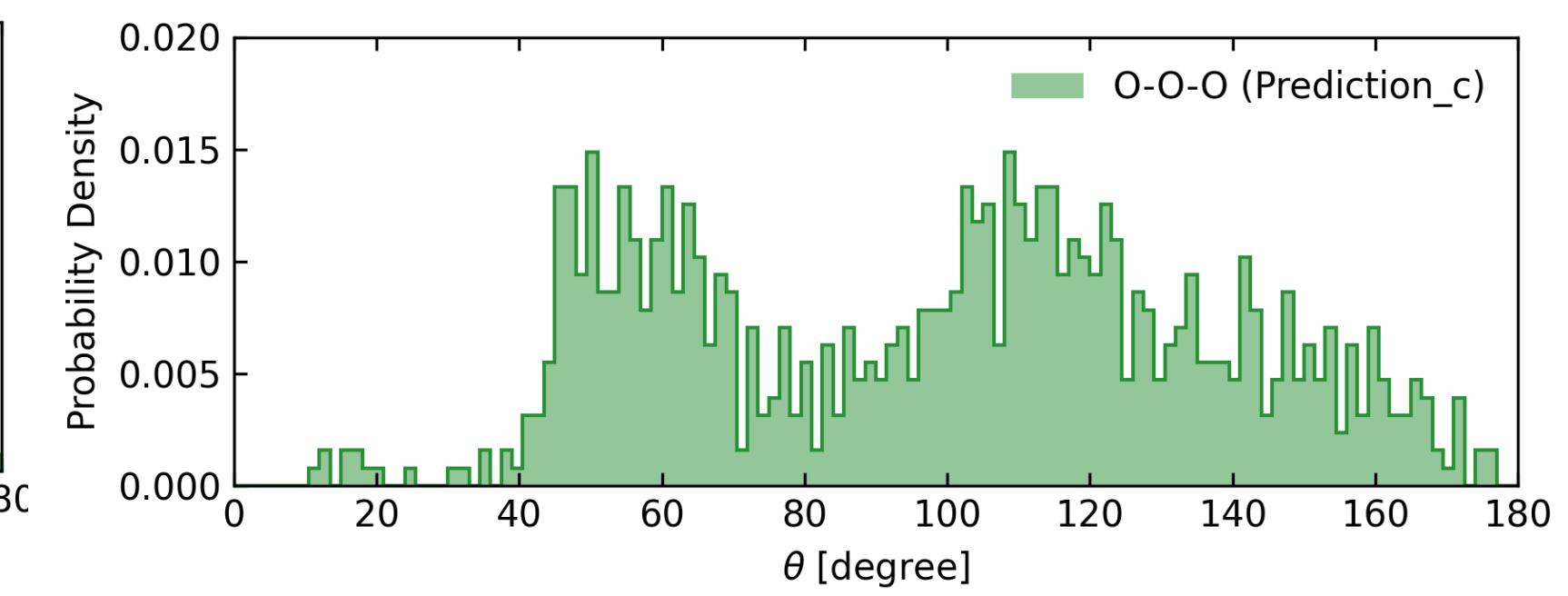
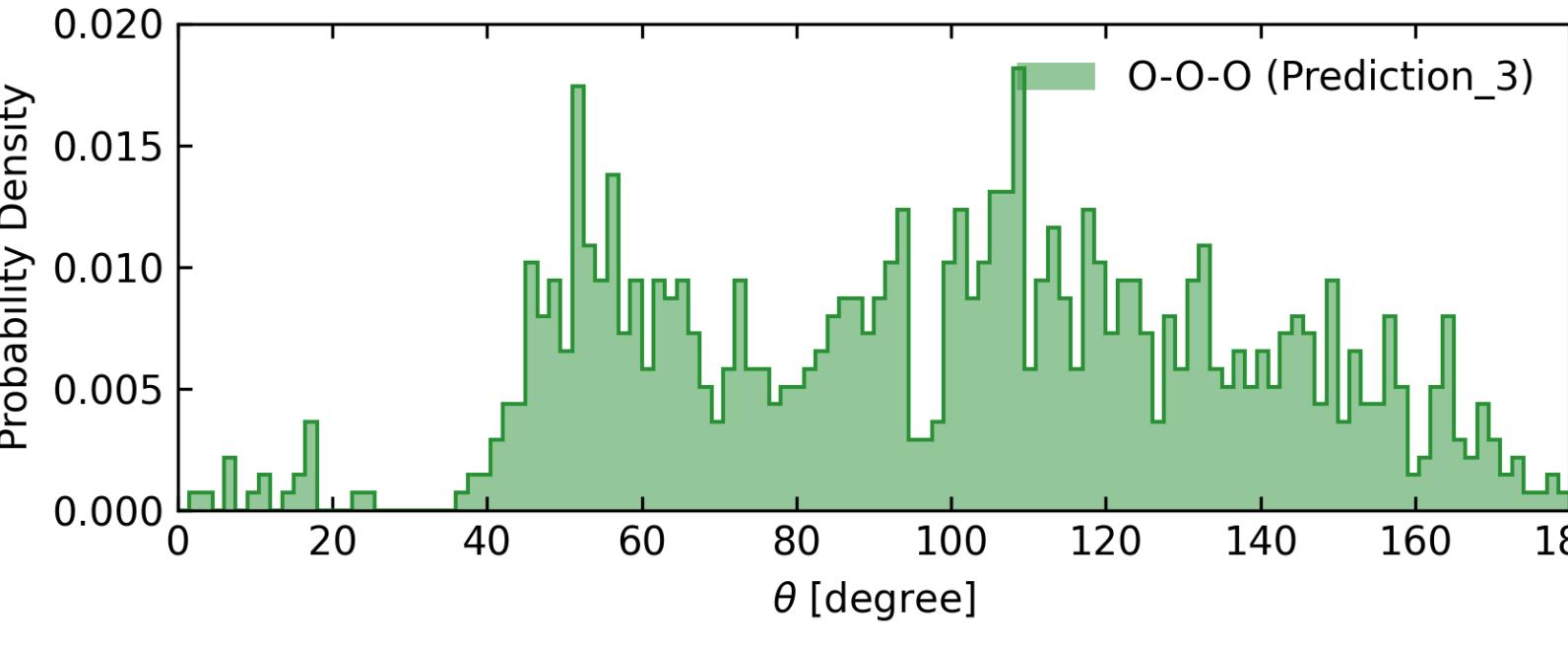
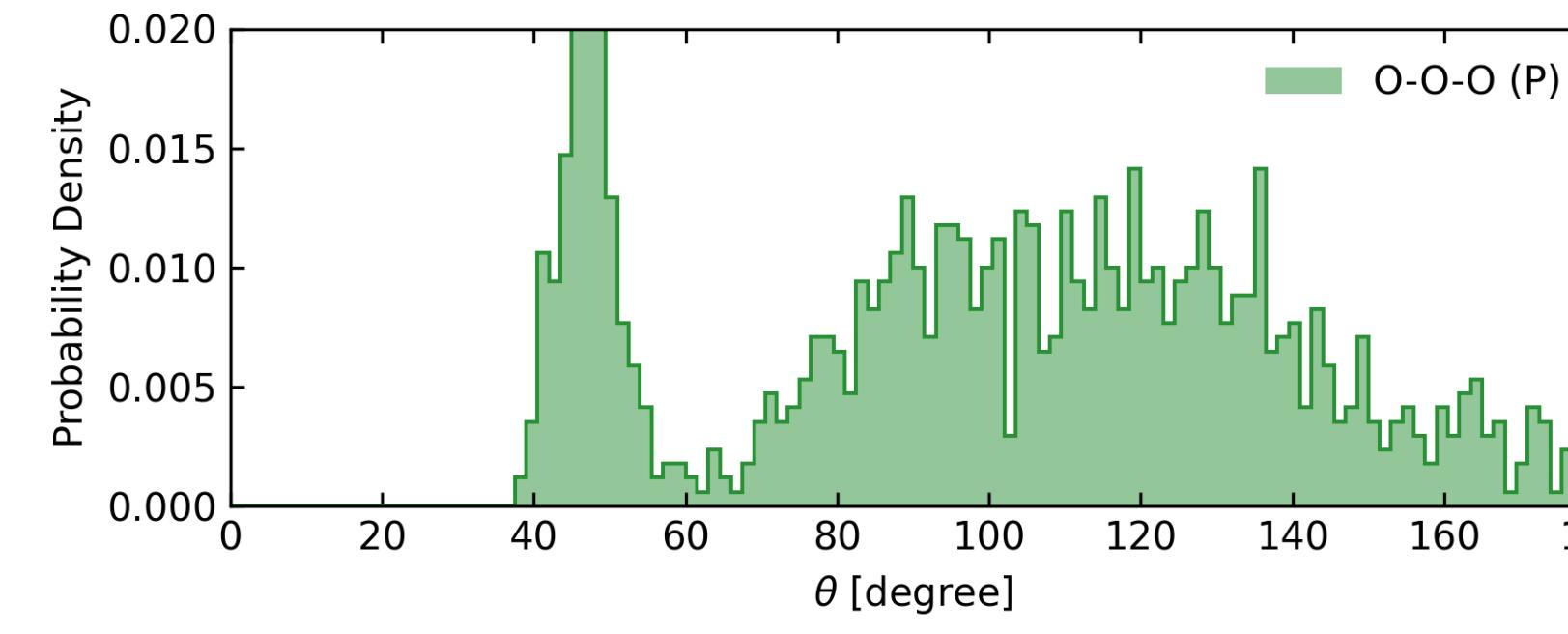
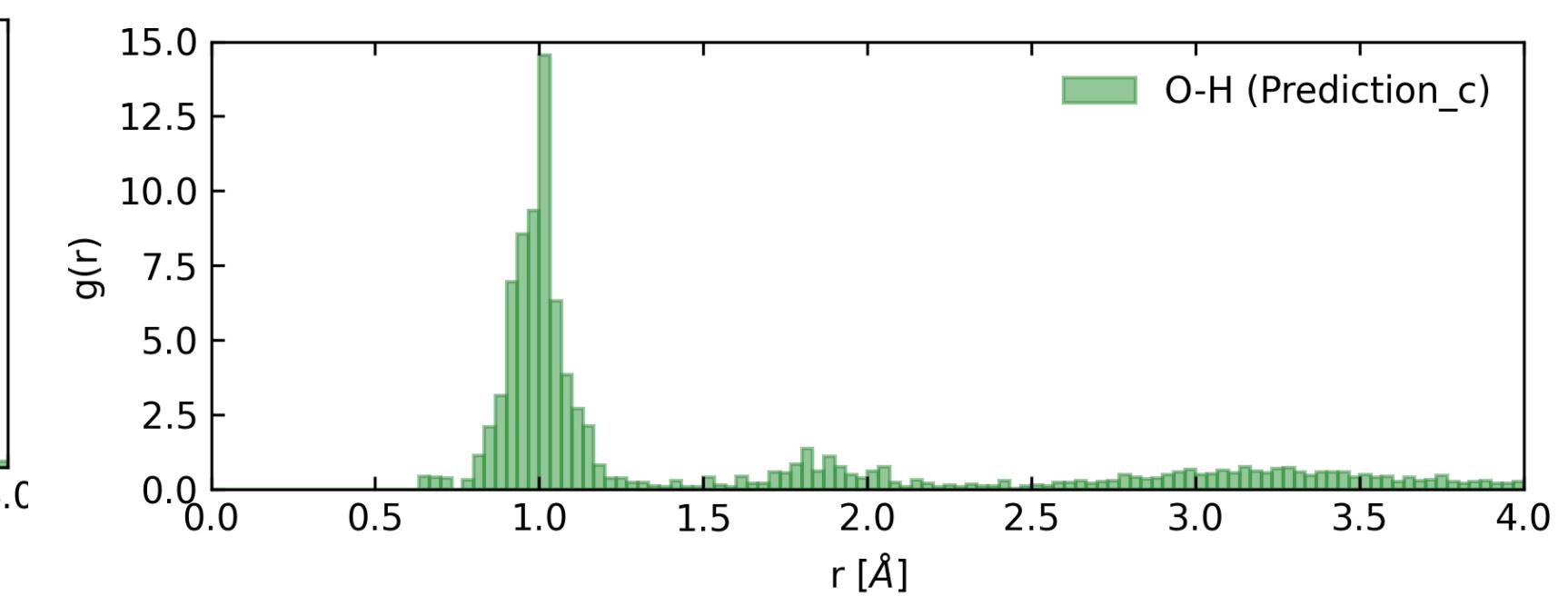
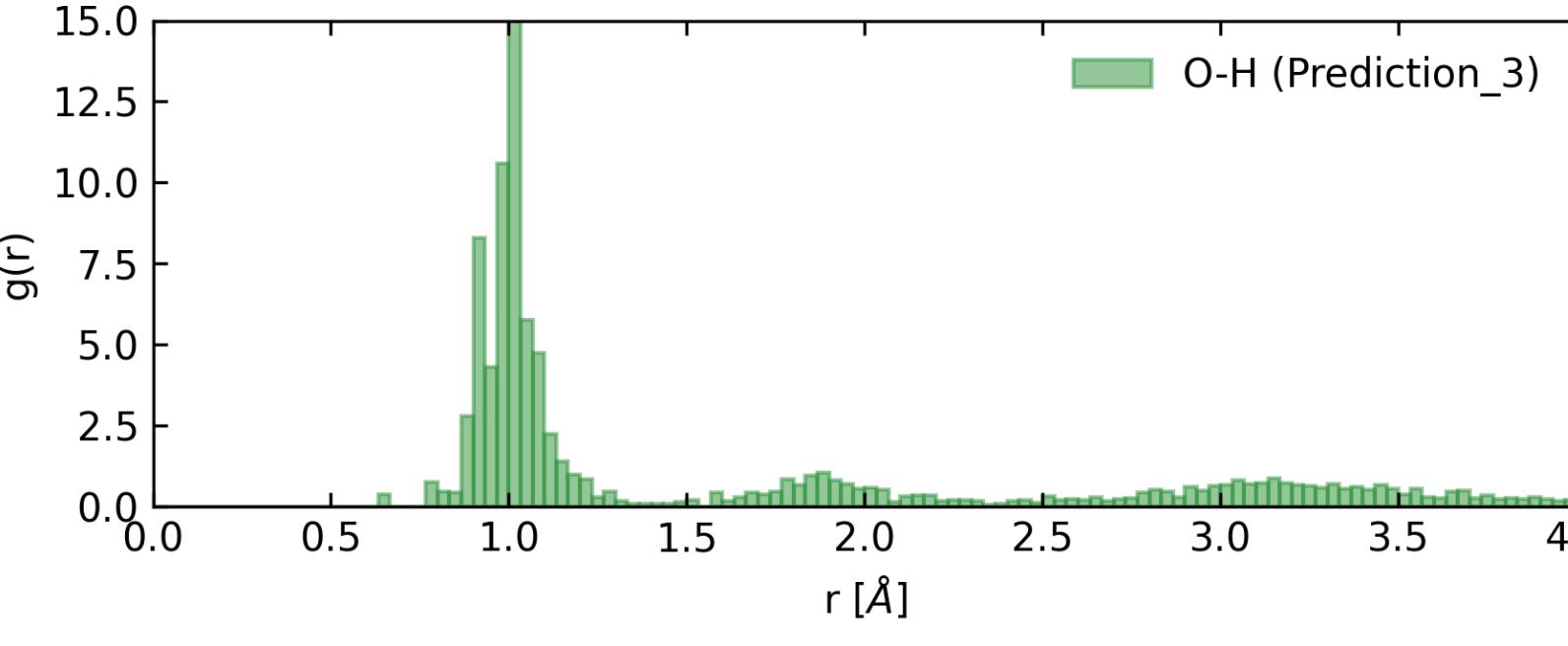
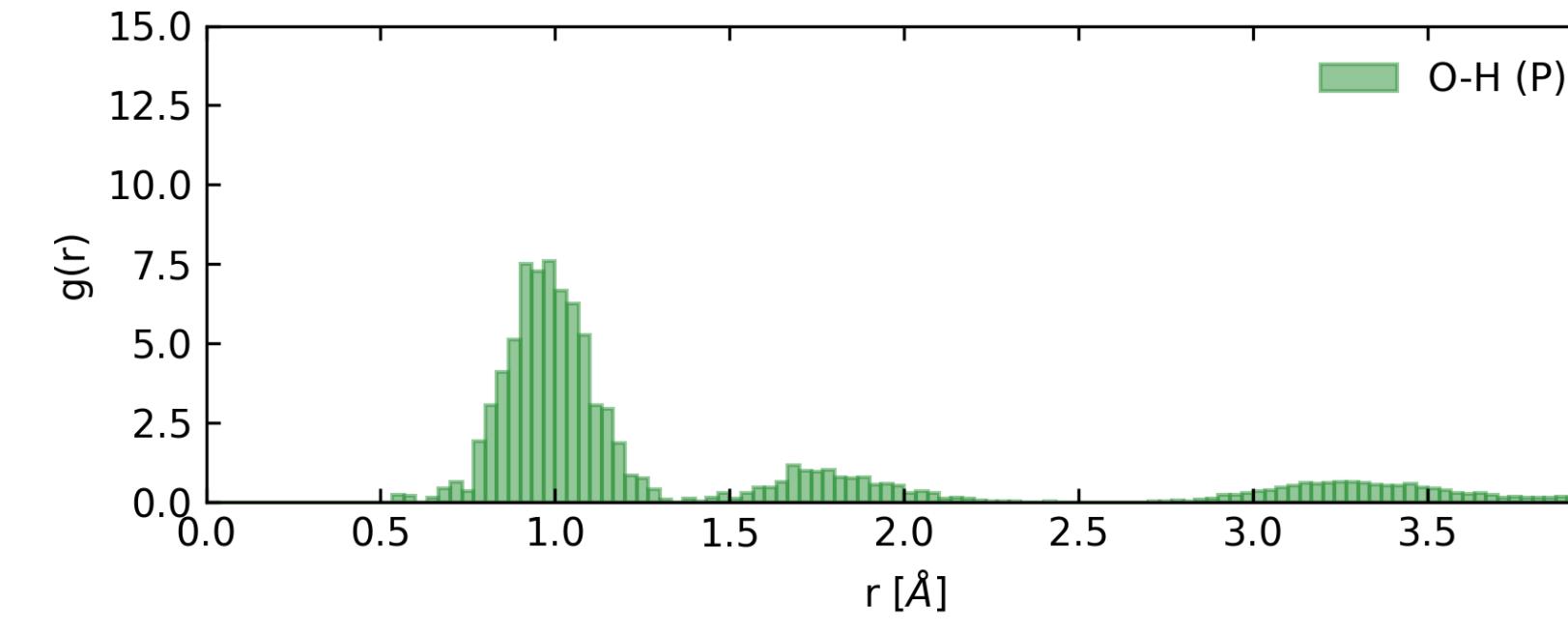
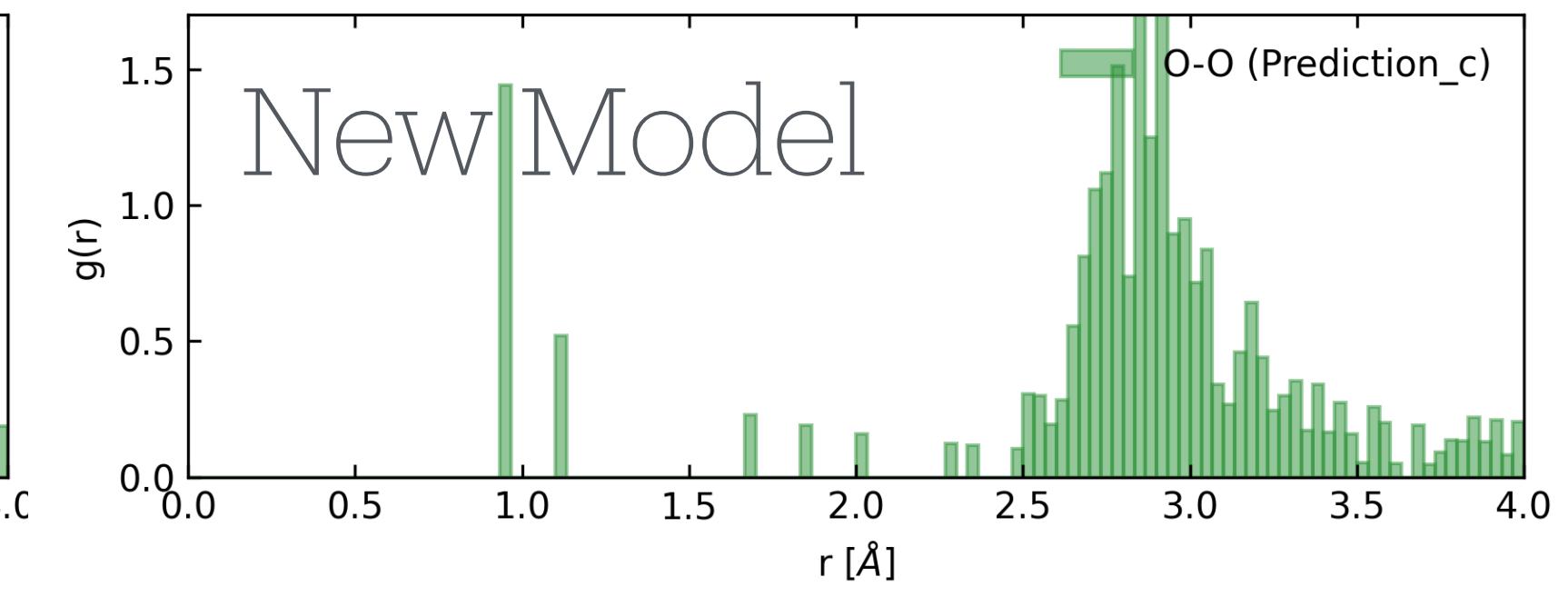
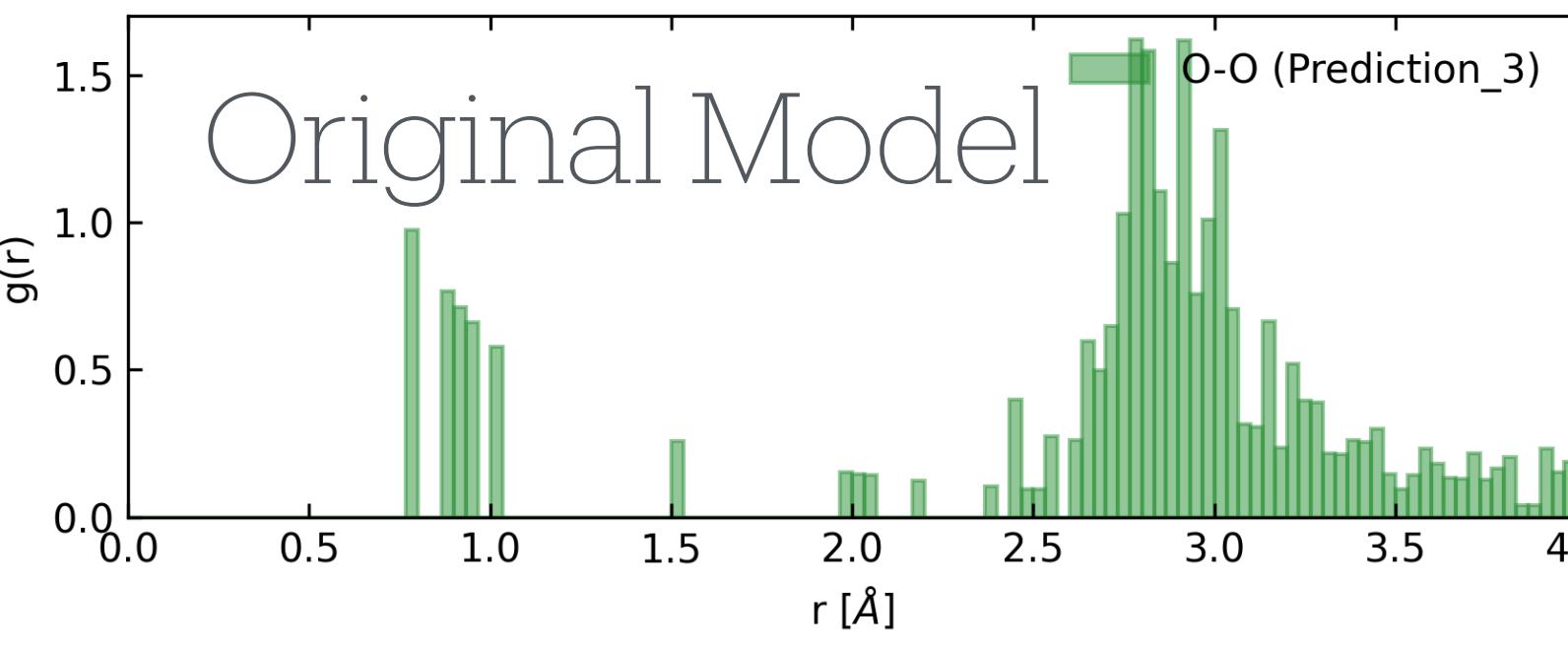
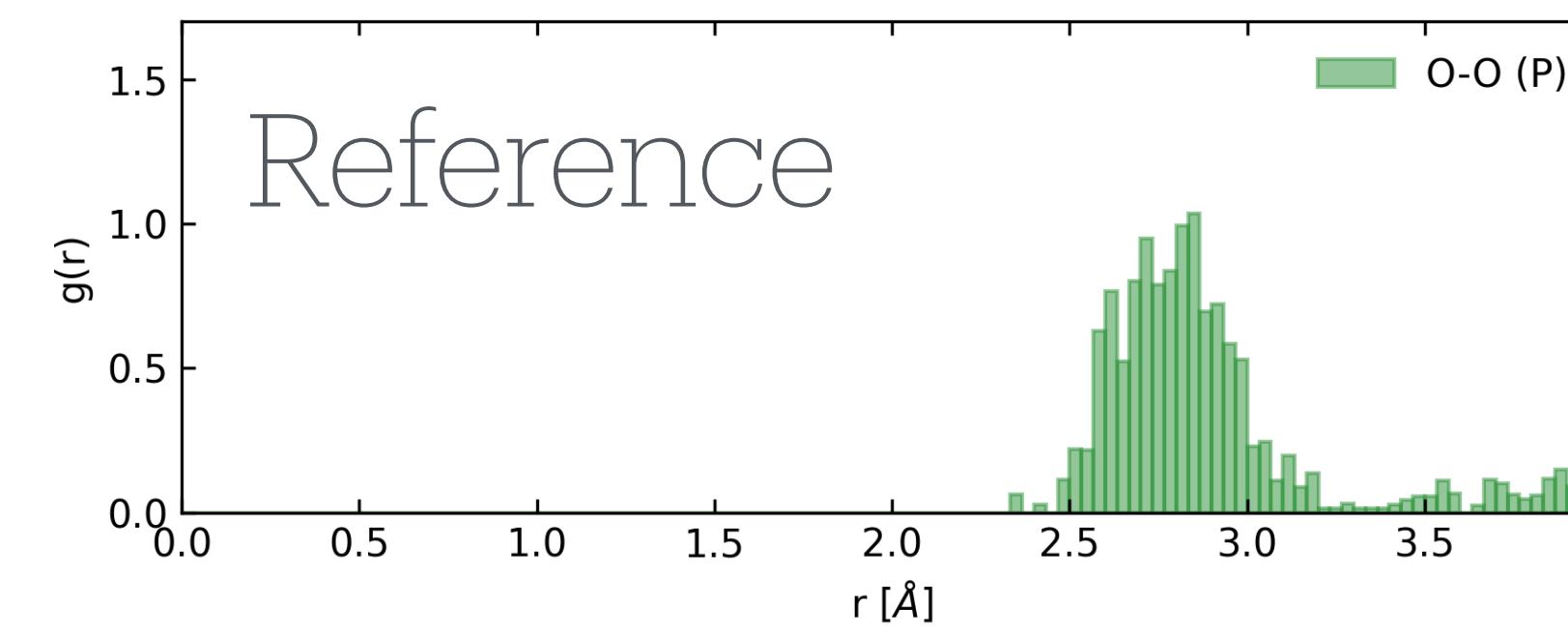


Assumption:

The predicted structure properties (RDF, ADF) on simulation AFM and real AFM are pretty close.

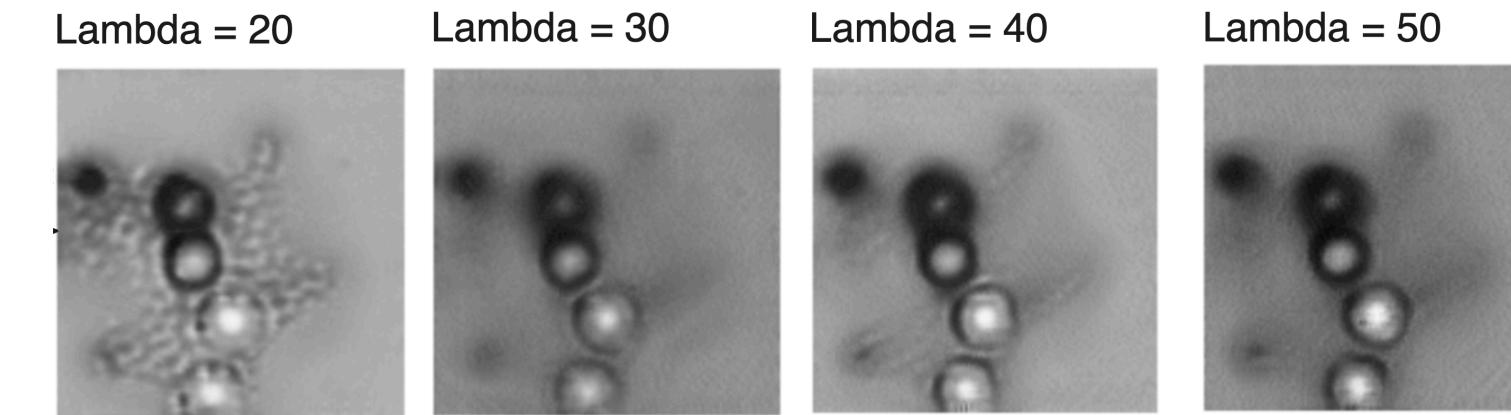
# Structure properties: RDF and ADF



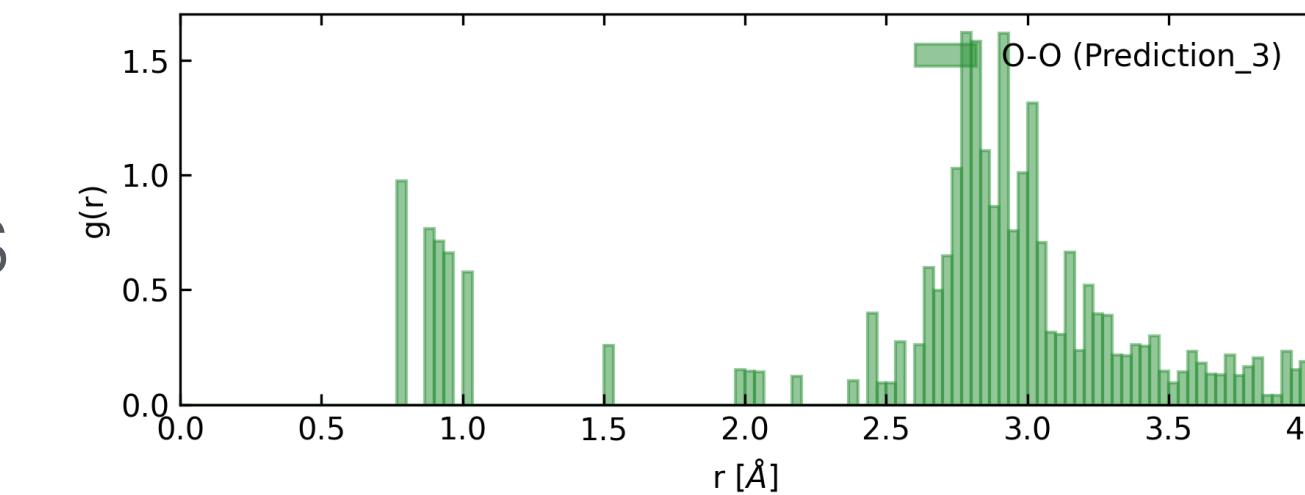


# Possible reasons and directions

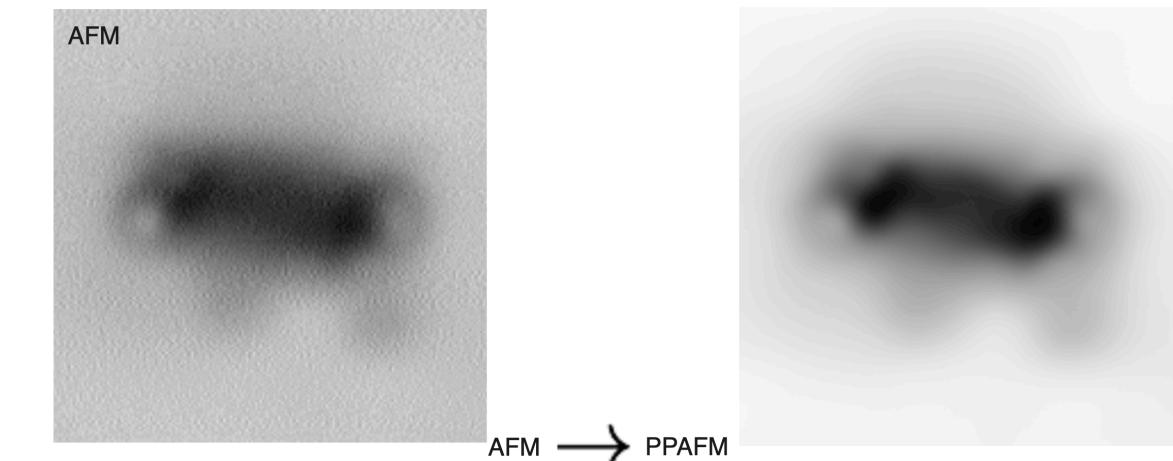
- Hyper-parameter in style translation  
(good in the eyes of human and machine)



- Limitation of data augmentation  
(additional constraint when designing of ML models)

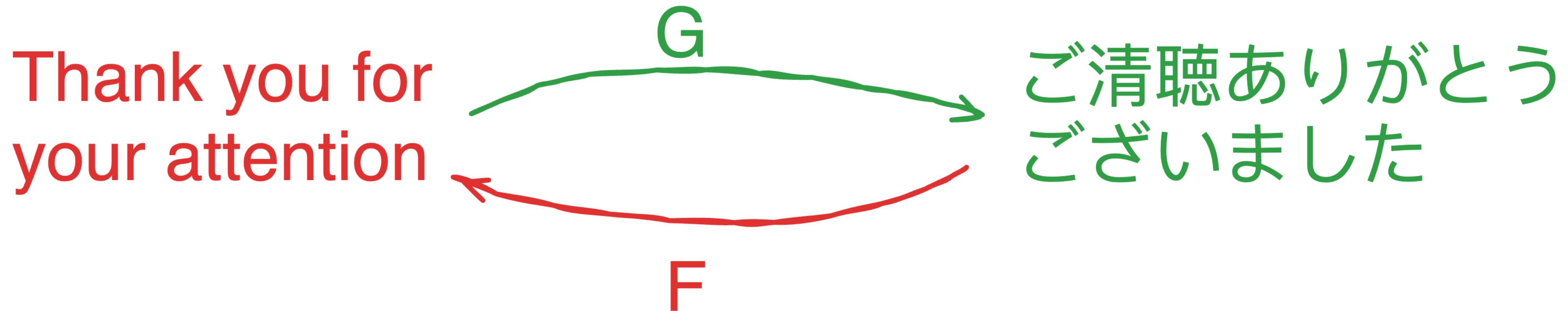


- Using the inverse translation as AFM preprocessing  
(removing is easier than adding, no need to retrain new model)



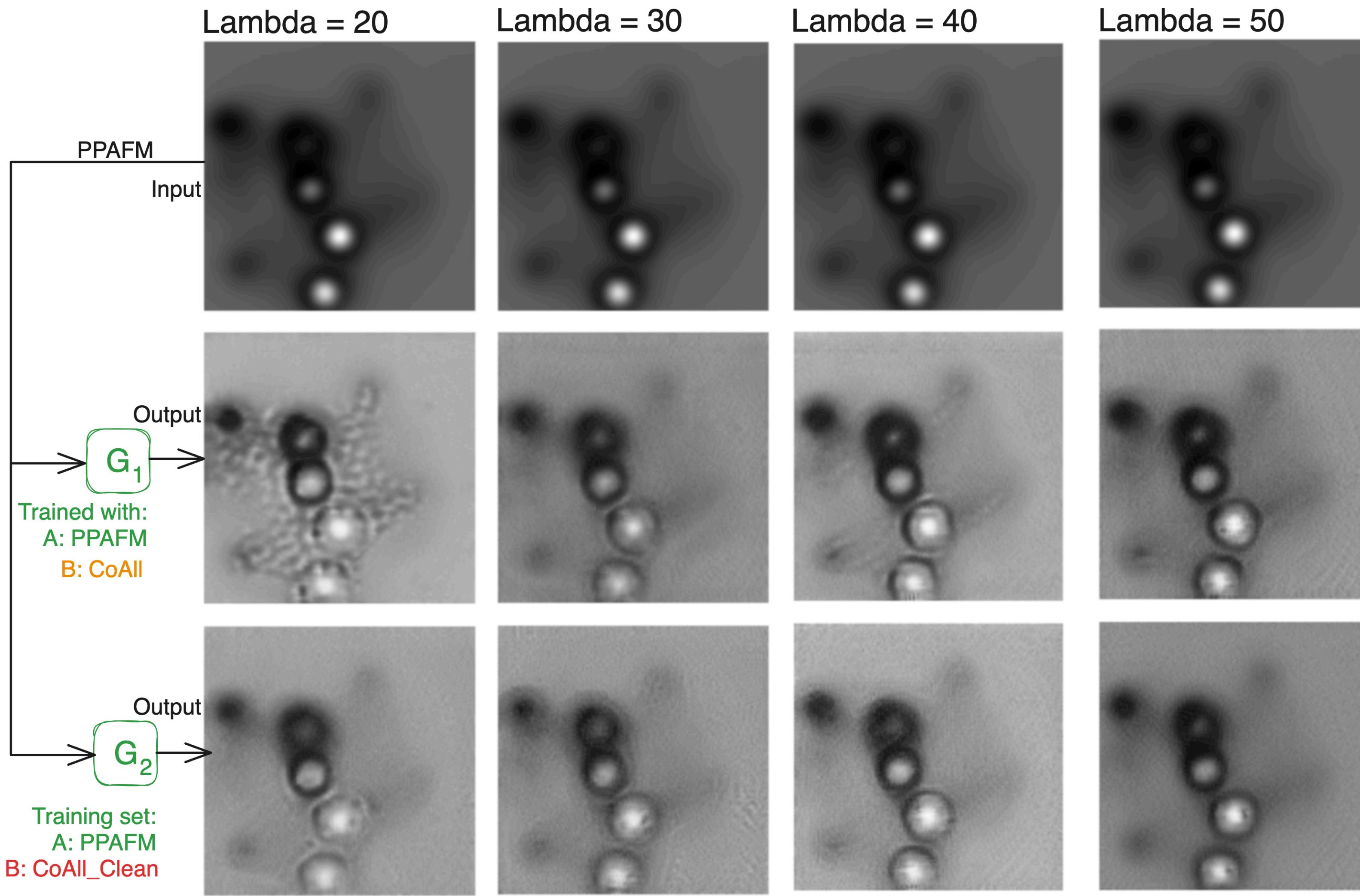
# Summary

- CycleGAN provides an effective solution for style translation between simulation and real AFM images.
- Using style translation as a tool of data augmentation is possible. But whether data augmentation can largely improve the performance of ML models is still uncertain.
- The inverse translation from AFM to simulation AFM could be used as an pre-processing (denoising) tool and it's promising to enhance the ML model performance.



# Appendix

# Observations

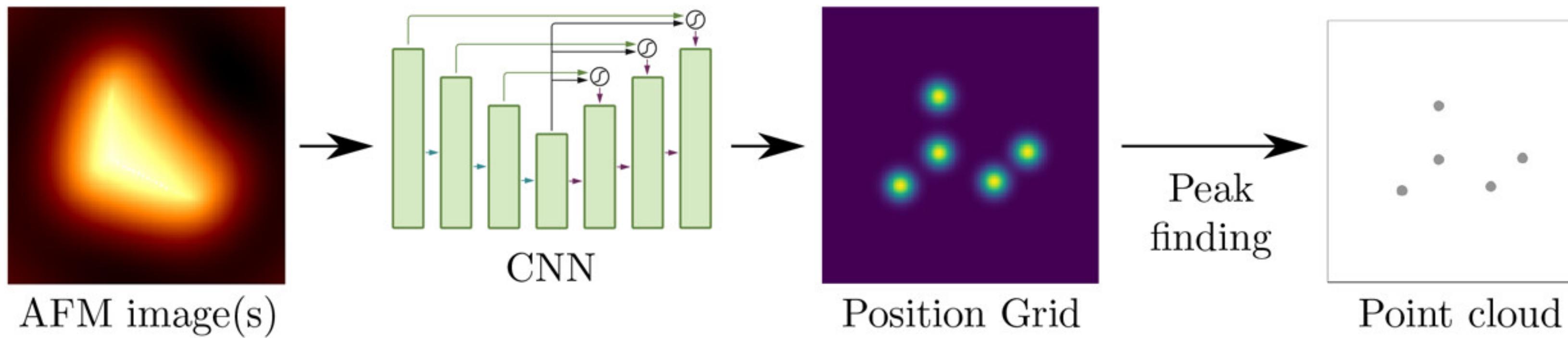


- Keep model structure simple (less parameters)
- Larger Lambda ( $L$ )
- Keep dataset clean

Does the fake AFM dataset good enough?

# Appendix

## 1. Find atom positions



## 2. Construct graph

