

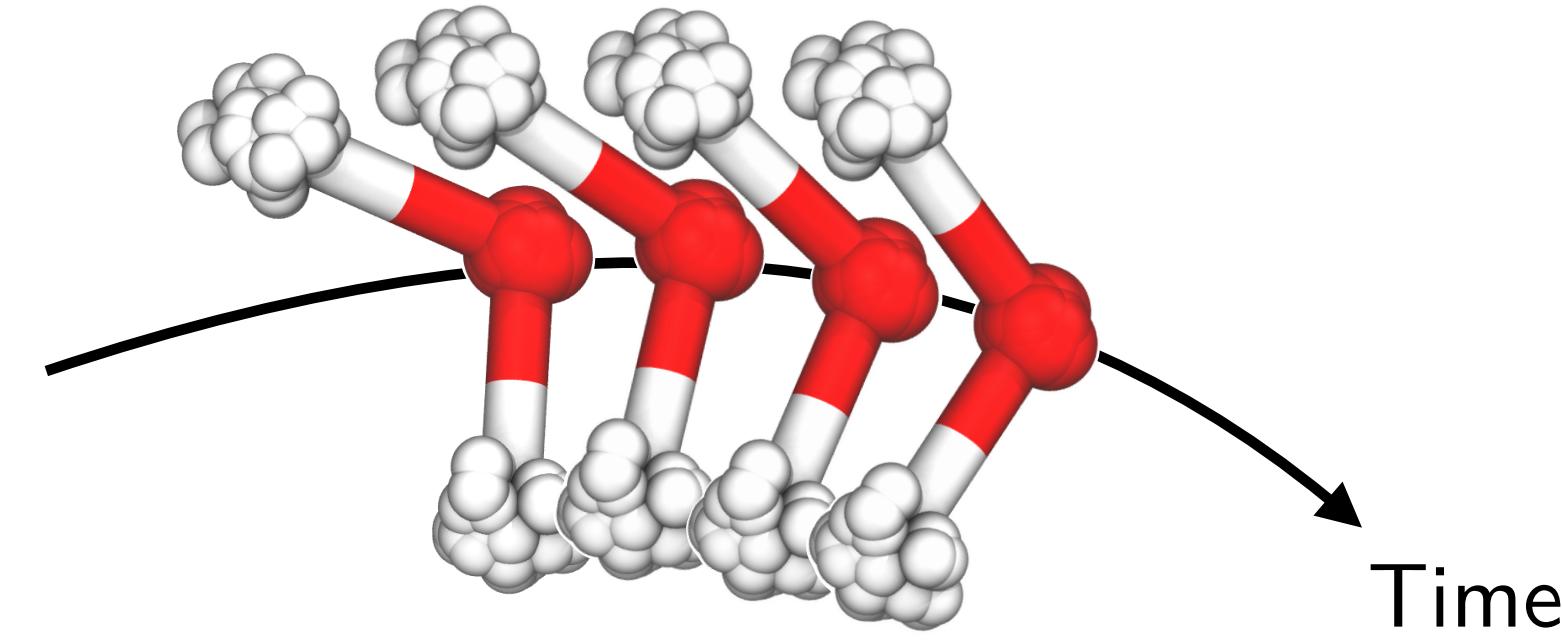
# Machine learning interatomic potentials – active learning techniques

**Elia Stocco, Shubham Sharma, Krystof Brezina and Mariana Rossi**

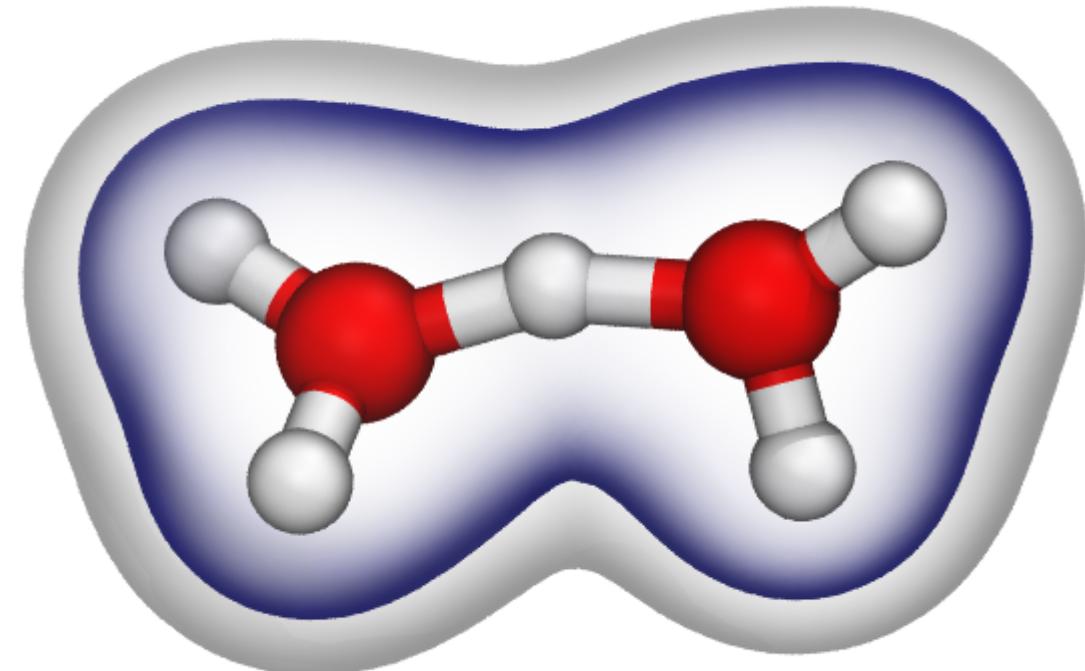
Max Planck Institute for the Structure and Dynamics of Matter  
Hamburg, Germany

CNPEM/Illum – Max Planck Meeting on Electronic Structure and Materials Informatics  
Campinas, São Paulo, Brazil  
July 3, 2025

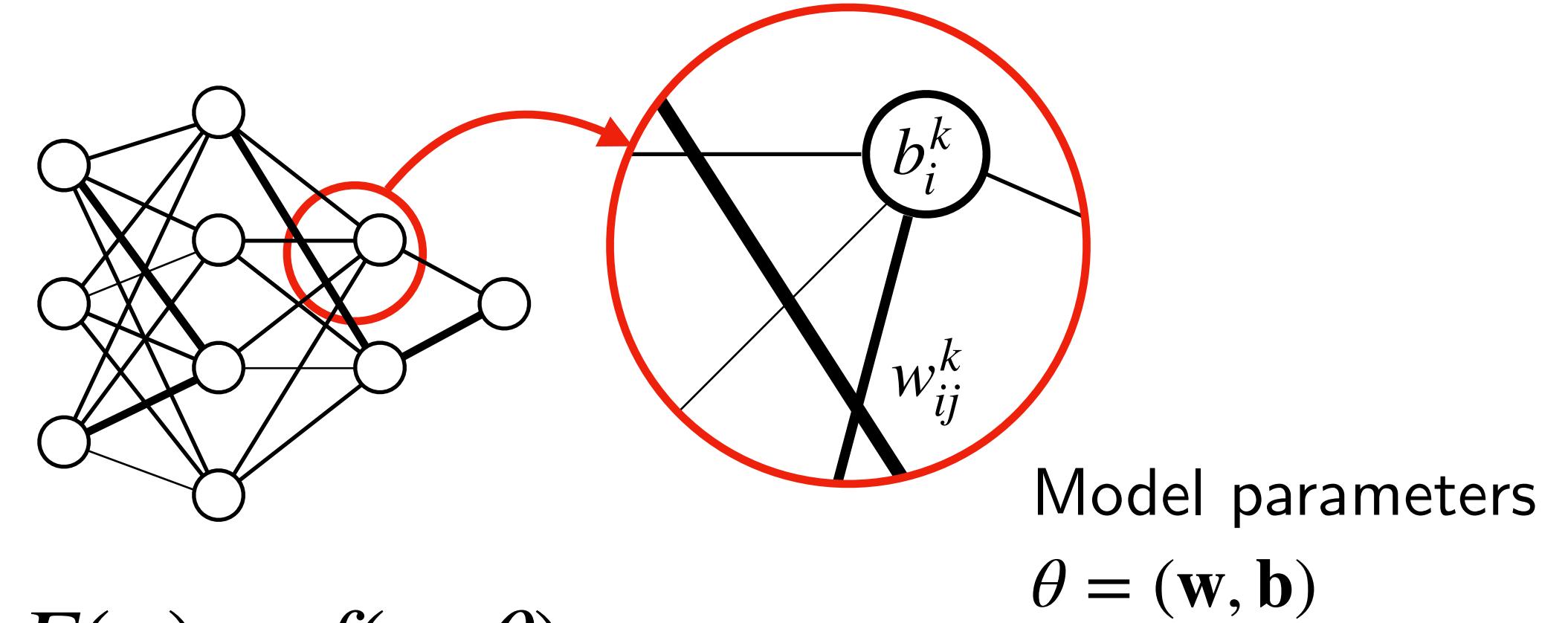
# Machine-learning interaction potentials (MLIPs)



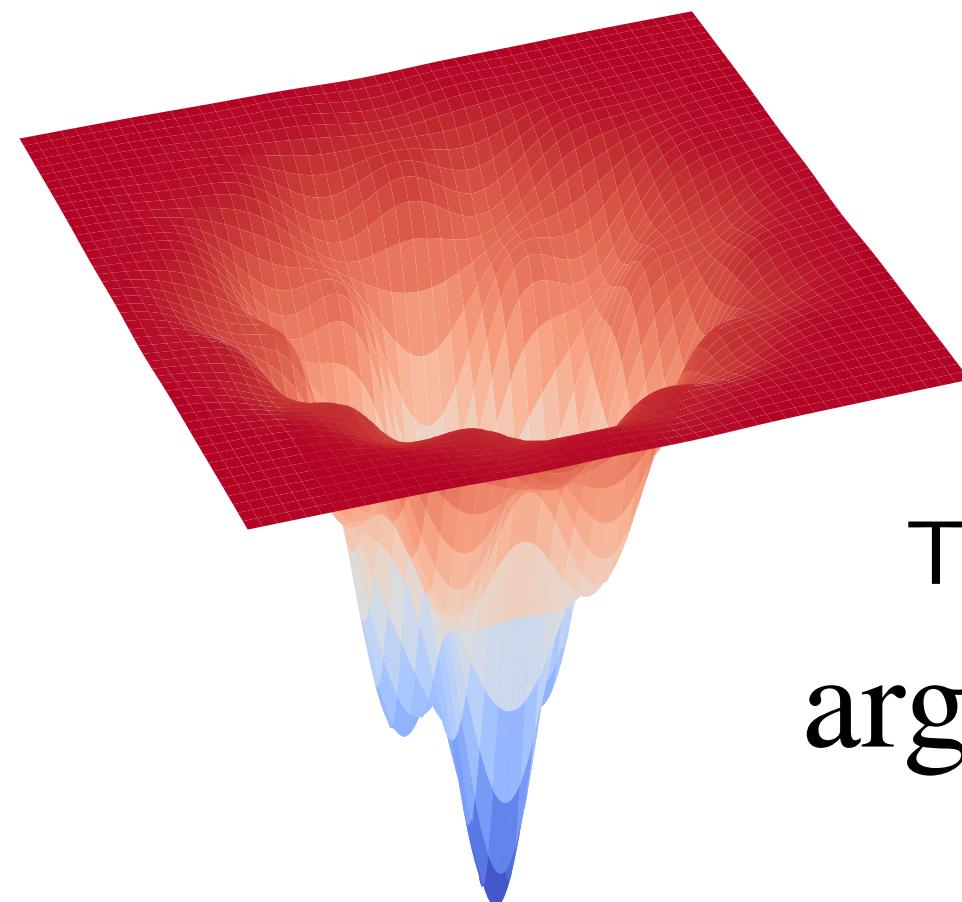
$$\mathbf{M}\ddot{\mathbf{q}} = - \nabla E(\mathbf{q})$$



$$E(\mathbf{q}) = \langle \Psi_0 | \hat{H}_{\text{electronic}}(\mathbf{q}) | \Psi_0 \rangle$$



$$E(\mathbf{q}) = f(\mathbf{q}; \theta)$$



$$\text{Training} \\ \arg \min_{\theta} \mathcal{L}(D, \theta)$$

# What do we require from the training data set $D$ ?

- **Uniform** and **dense enough** coverage of the thermally available configuration space
- **Interpolation** vs. **extrapolation** prediction regime
- **Compact** and **robust** data sets

Tent-pole analogy of **good** training sets:

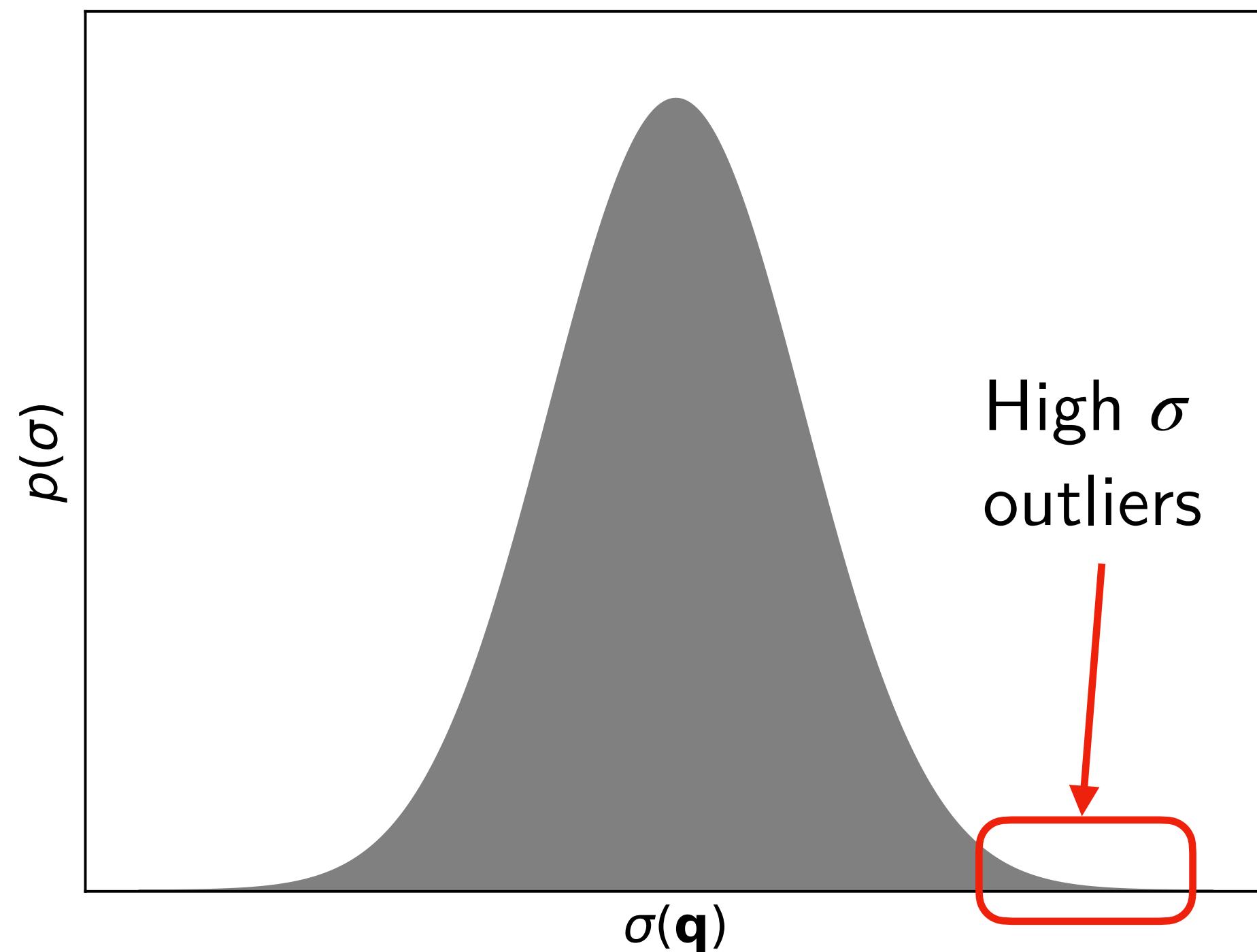


# Active learning

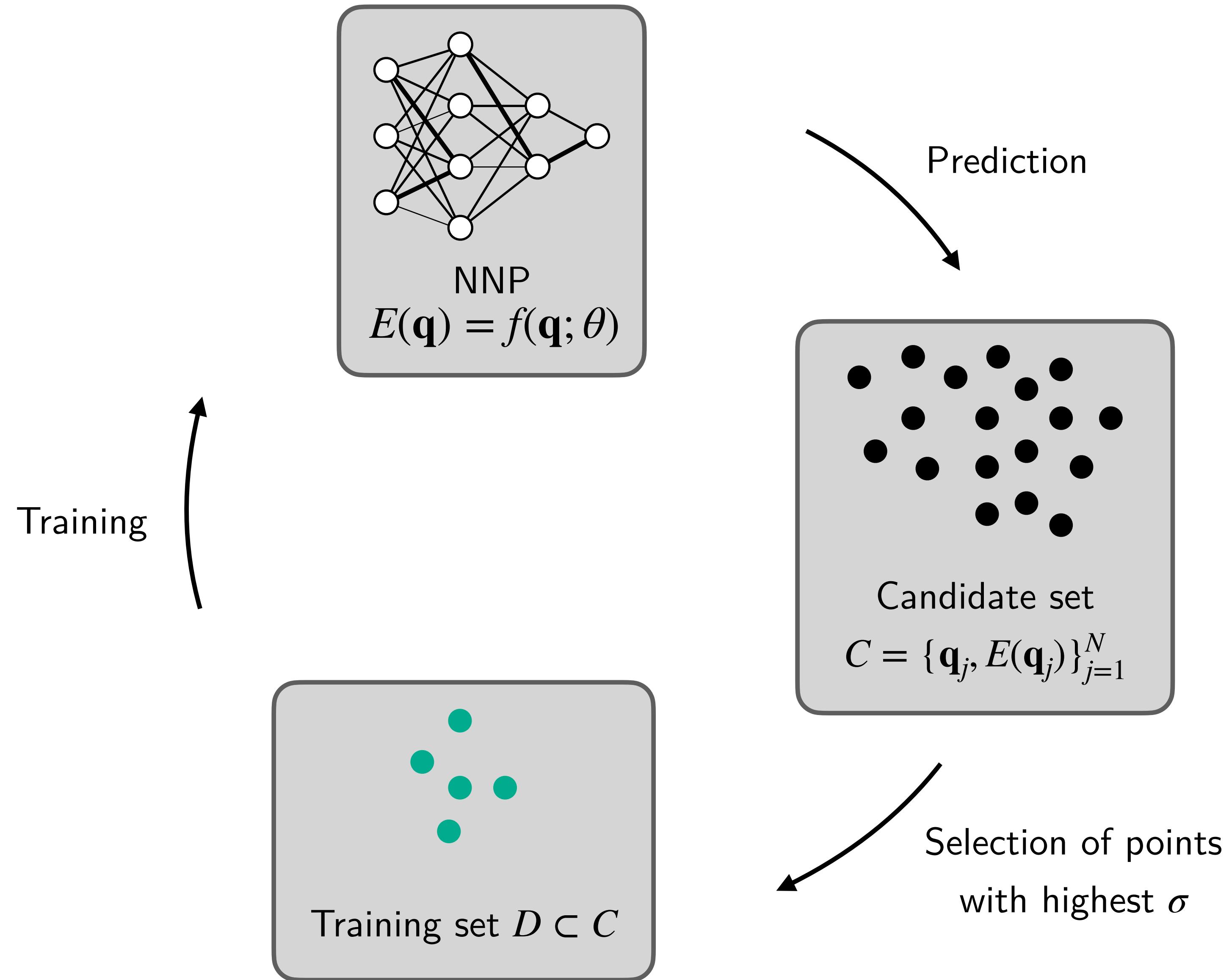
**Active learning:** selection of new data points based on some *uncertainty* measure  $\sigma$

**Regime 1:** filtering a *labeled* candidate set

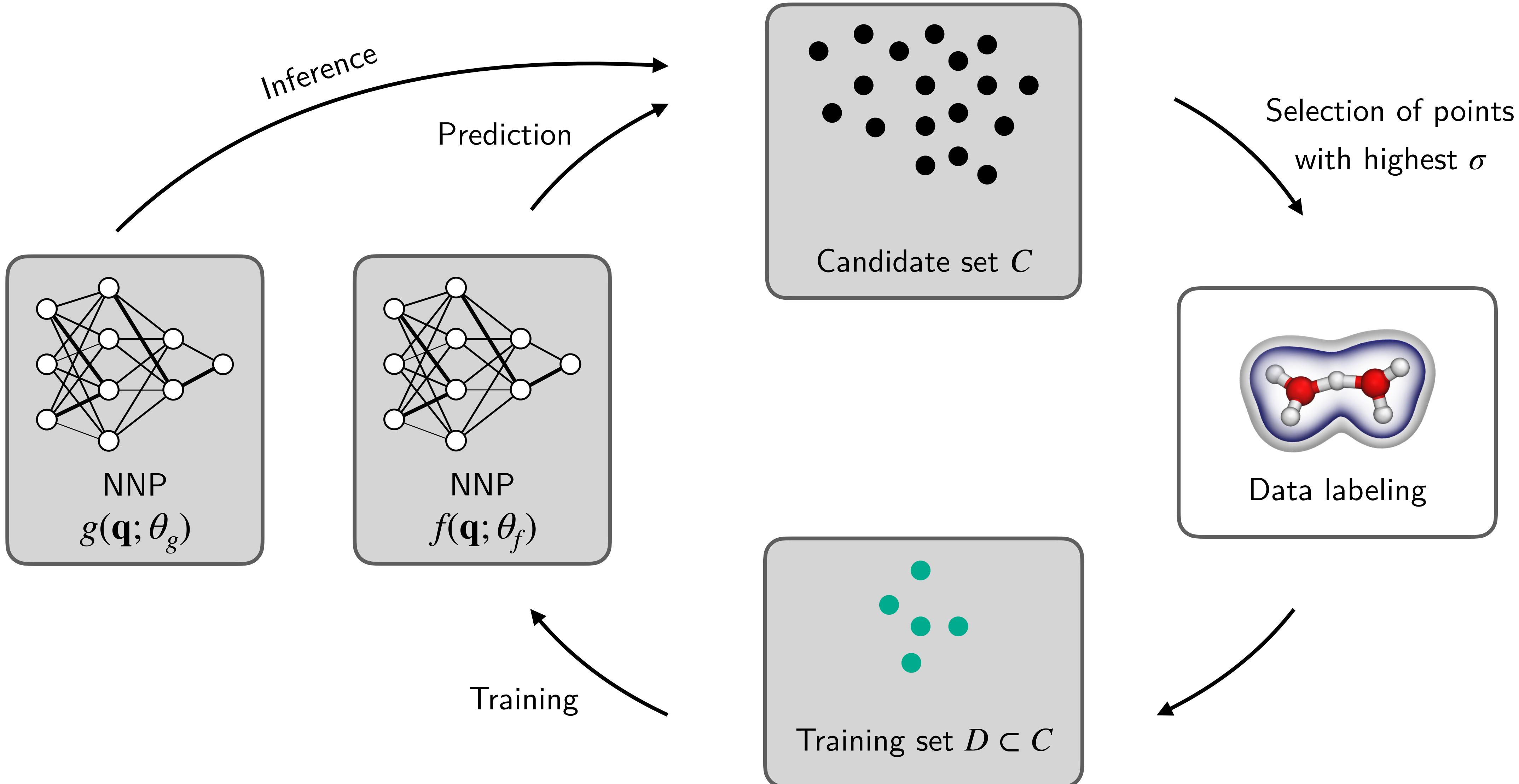
**Regime 2:** selecting relevant data from an unseen, exploratory candidate set (*unlabeled*)



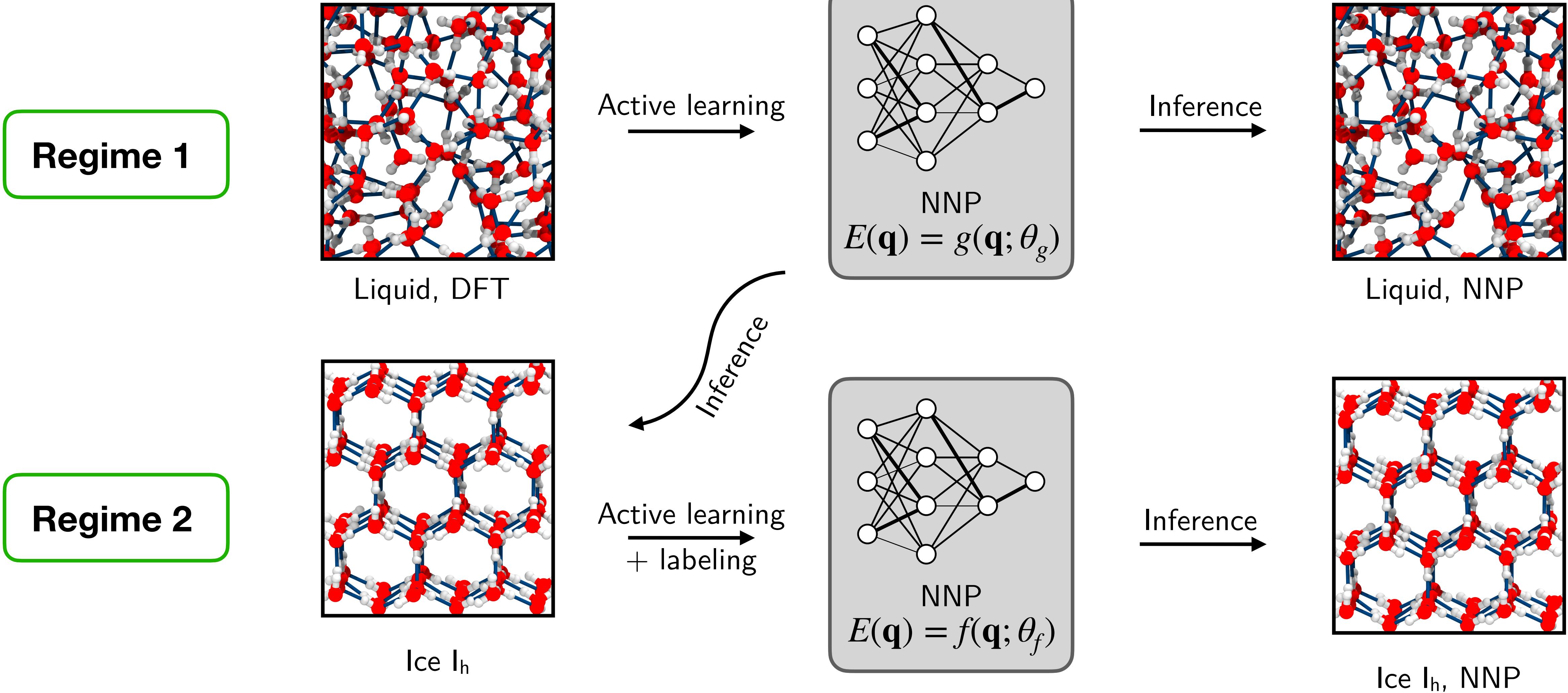
# Regime 1: labeled candidates



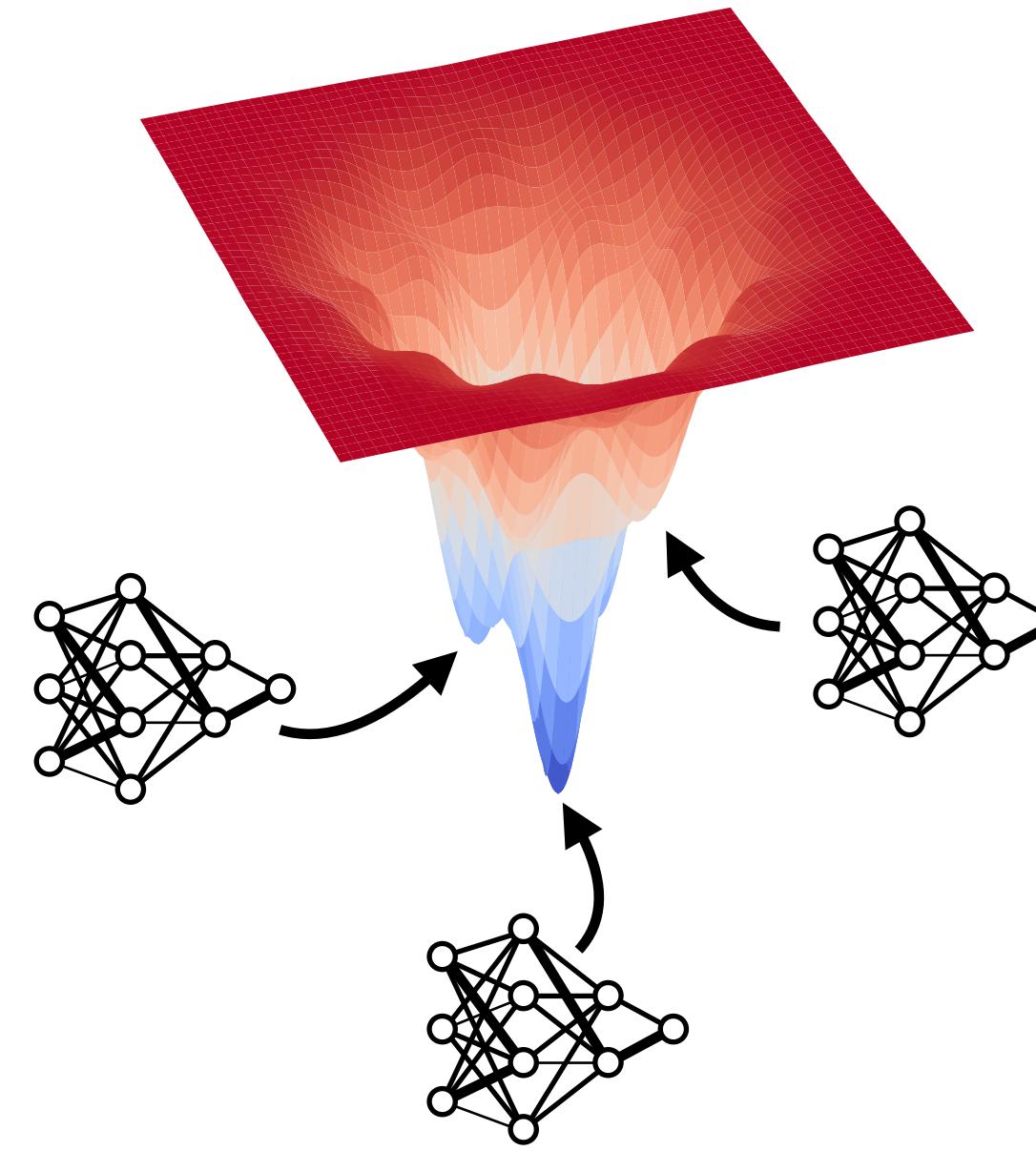
# Regime 2: refinement using unlabeled candidates



# A real-world example: thermodynamic state points of water

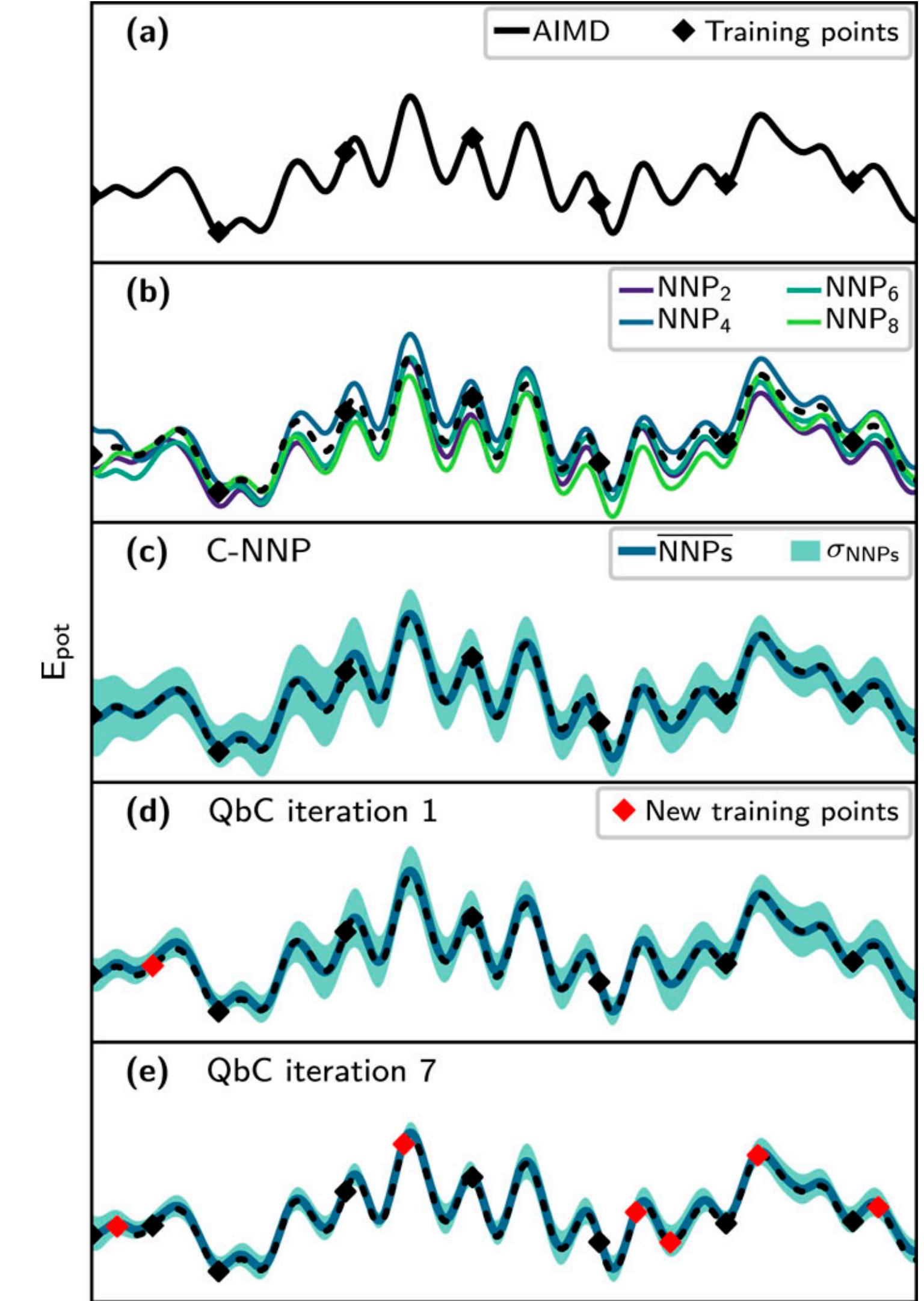
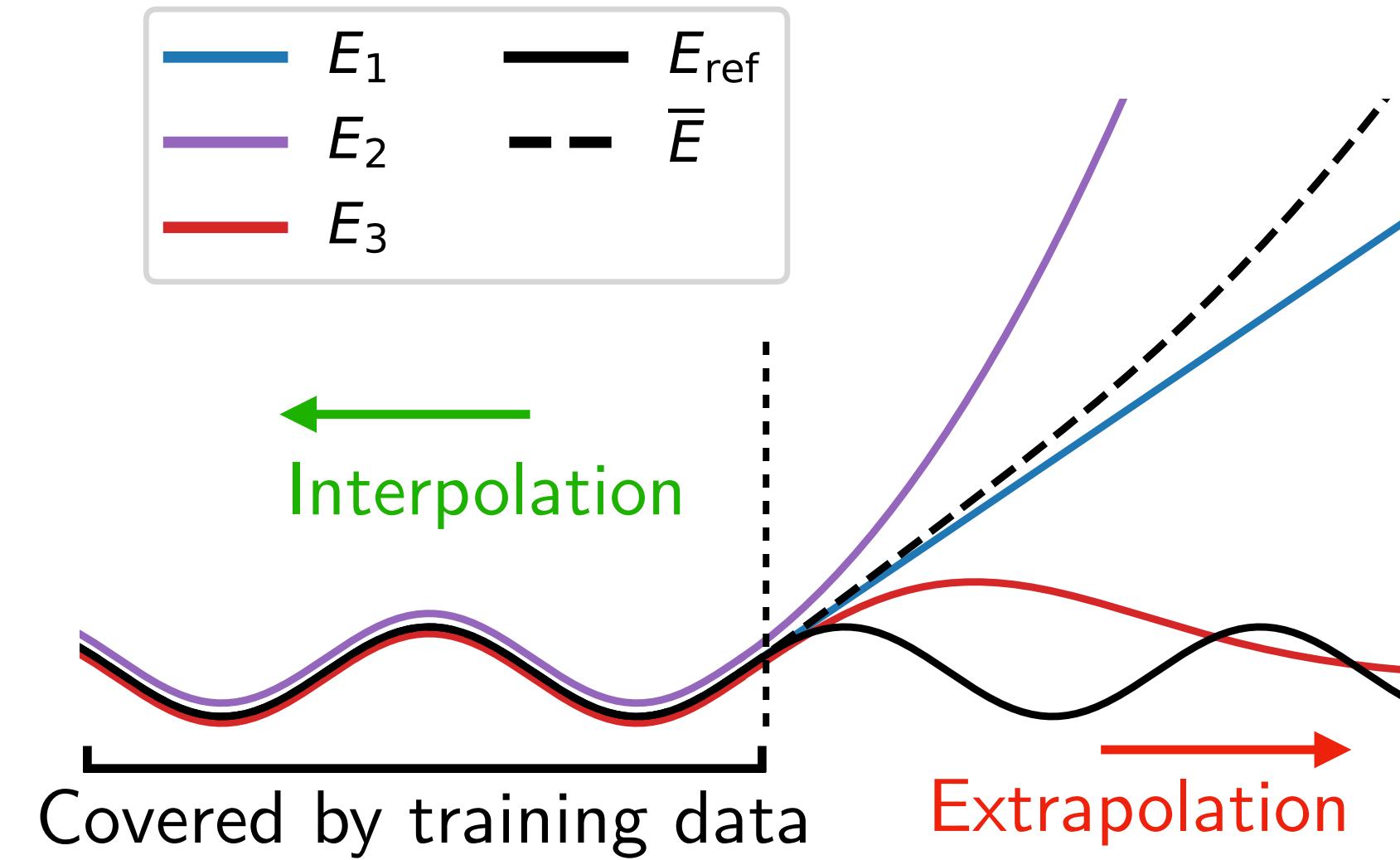


# Committee models & Query by committee



Low disagreement → low generalization error

High disagreement → high generalization error



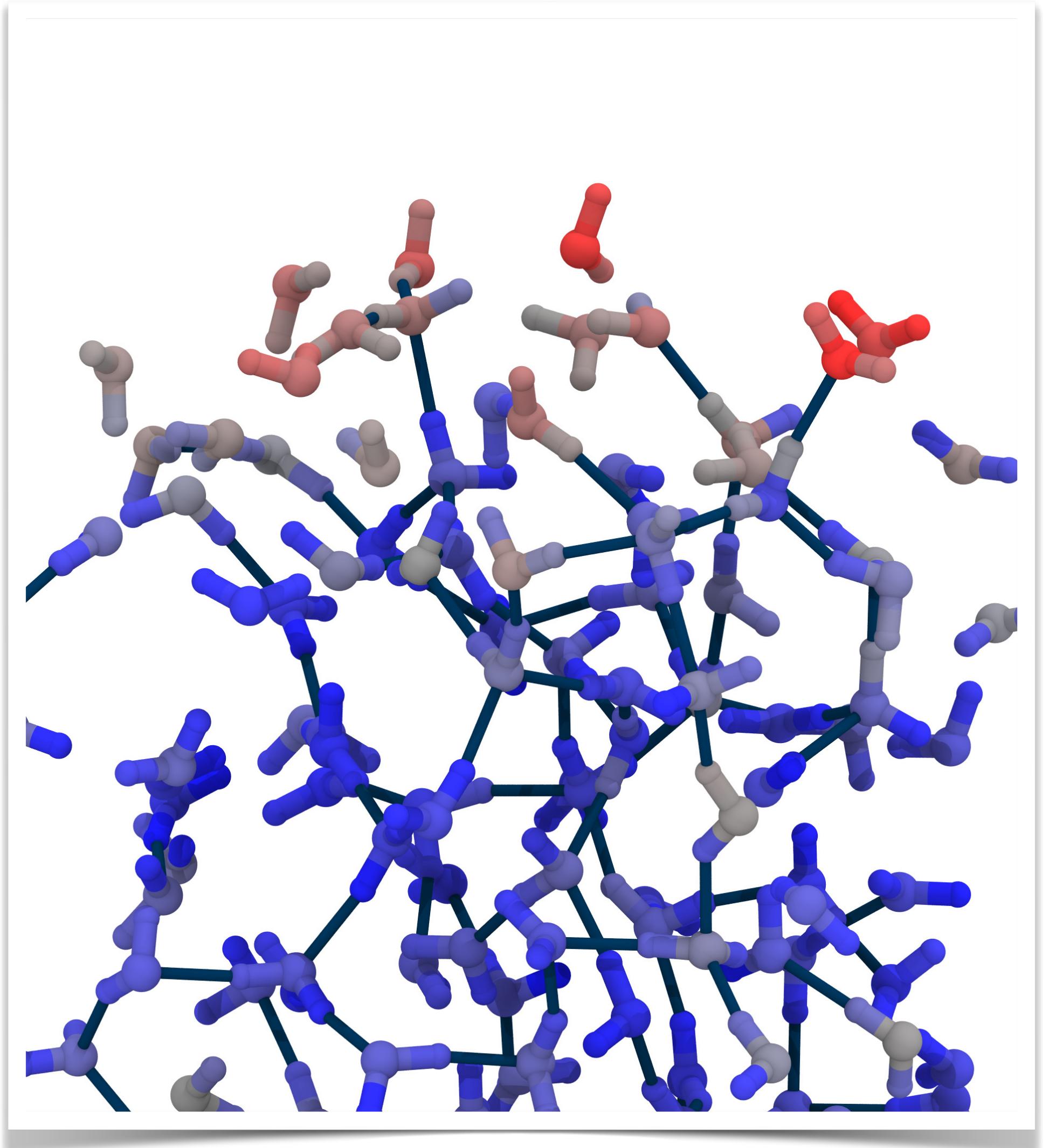
# Disagreement: energies or forces?

$$\sigma_E(\mathbf{q}) = \sqrt{\frac{1}{N_c} \sum_{n=1}^{N_c} [E_n(\mathbf{q}) - \bar{E}(\mathbf{q})]^2}$$

$$\bar{\sigma}_F(\mathbf{q}) = \frac{1}{N} \sum_{\alpha=1}^N \sqrt{\frac{1}{N_c} \sum_{n=1}^{N_c} [\mathbf{F}_{\alpha n}(\mathbf{q}) - \bar{\mathbf{F}}_\alpha(\mathbf{q})]^2}$$

Atoms      Committee members

Force disagreement is typically used in active-learning applications for practical reasons.

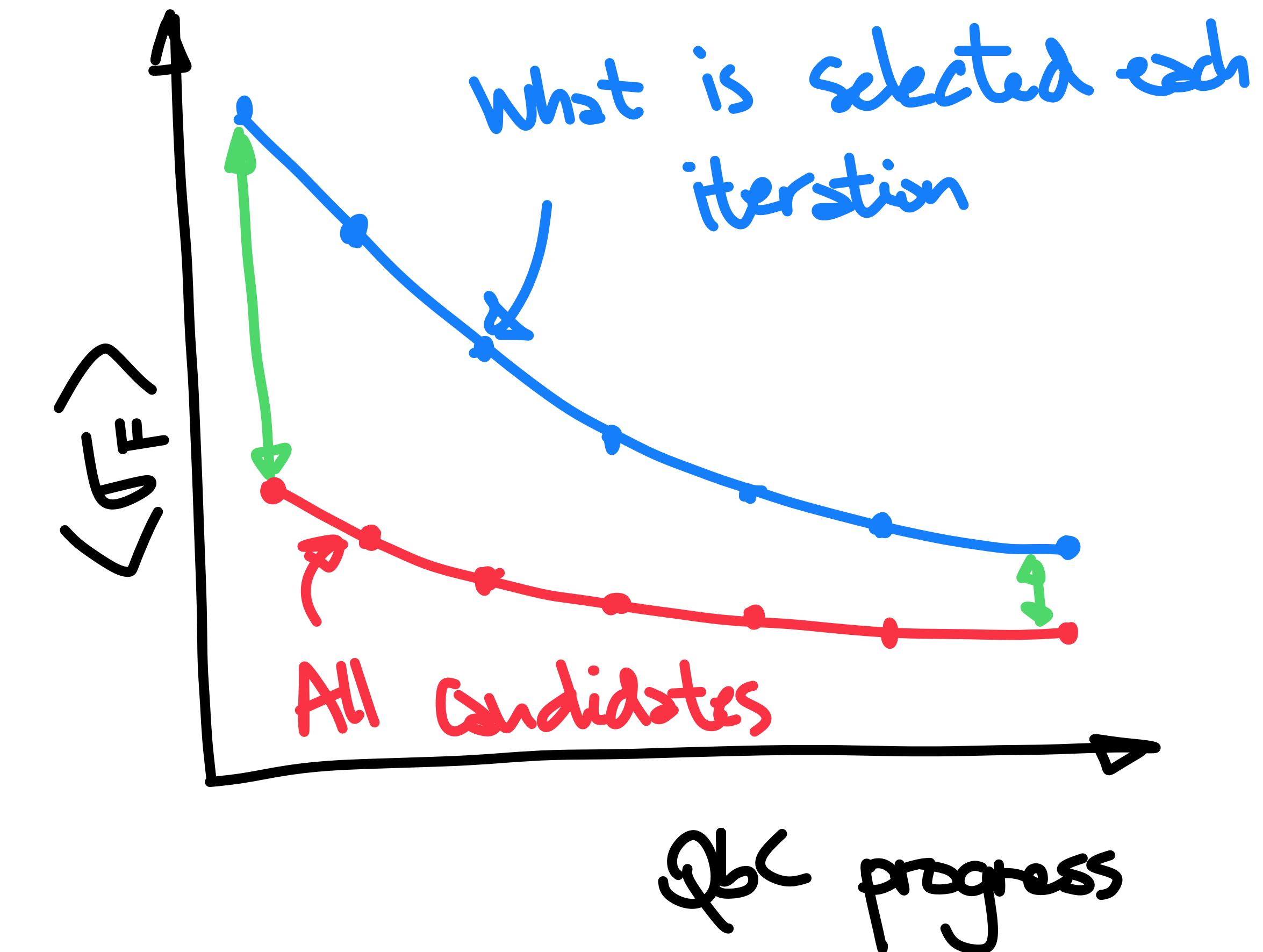


# Takehome messages: controlling QbC convergence

**QbC goal:** extract as much valuable **unlearned information** from the candidate set in as **few samples** as possible

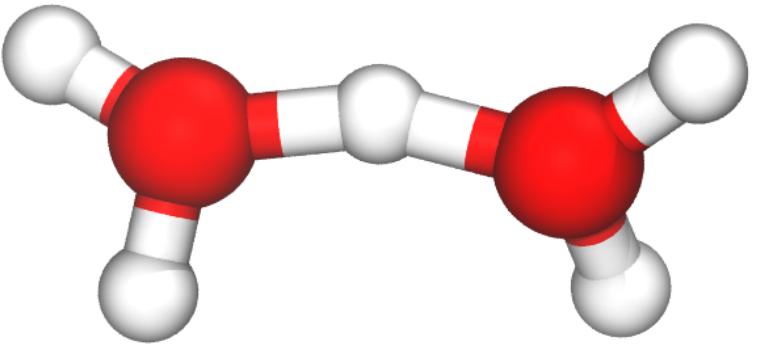
**Convergence:** this happens when the  $\sigma$  for the newly selected batch of structure **plateaus** and **approaches** the  $\sigma$  for the remaining candidates

**Generalization error:** acceptable error over an independent **test set**

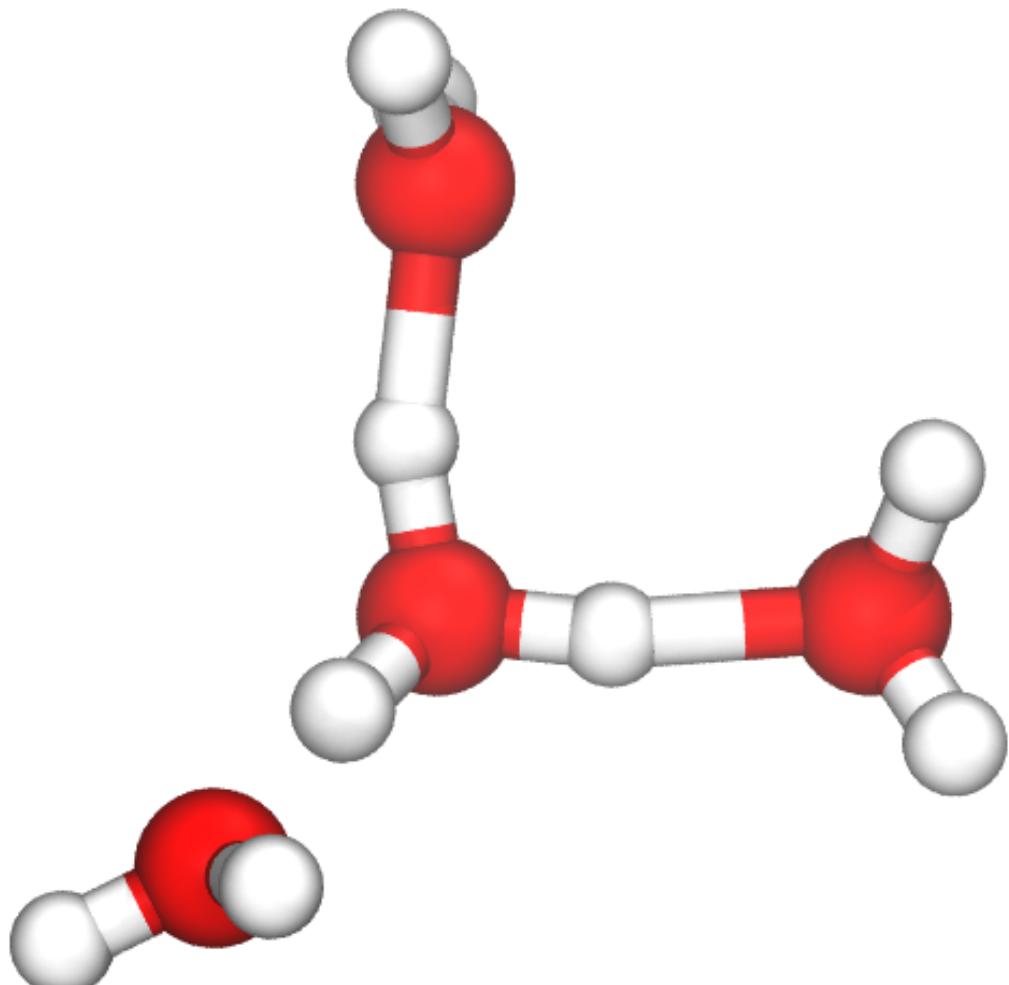


# Goals of the this tutorial

- **Actively select labeled data points** for the Zundel cation using **QbC** with a committee of **MACE [1]** models
- Explore the properties and the evolution of **committee disagreement** over the course of the QbC cycle
- **Compare the quality** of active-learned and random-sampled models
- Perform **exploratory** inference MD using the Zundel-trained model on the Eigen cation
- Use QbC with on-the-fly data labeling with **FHI-aims [2]** to **improve** the initial Eigen model



Zundel cation ( $\text{H}_5\text{O}_2^+$ )



Eigen cation ( $\text{H}_9\text{O}_3^+$ )

[1] Batatia I. et al. NeurIPS, 36, 2023

[2] Abbot J. W. et al. arXiv:2505.00125, 2025