

Overview

Data-centric informatics has promoted materials discovery and design

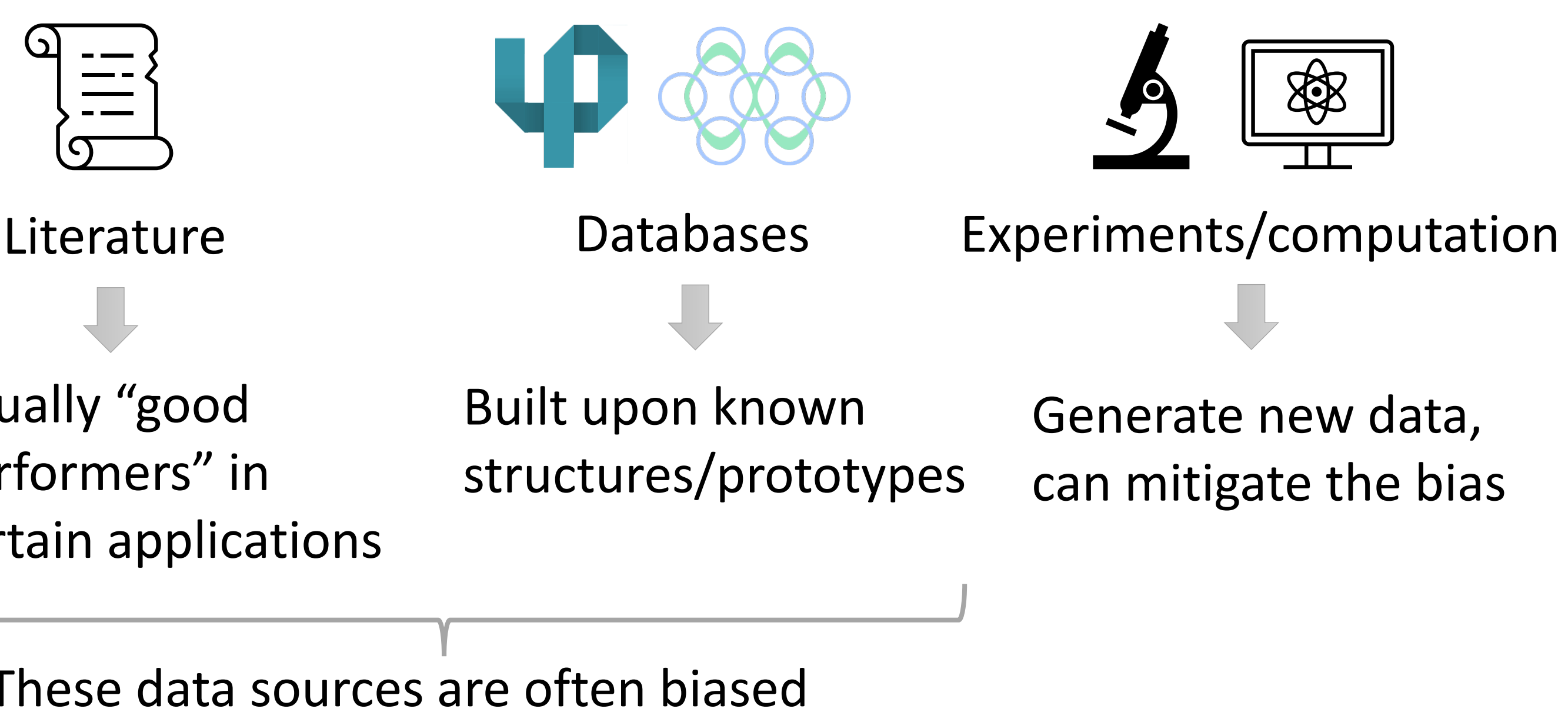
- Models, e.g., machine learning (ML), draw much attention
- Quality of data is equally important but less studied

This work focuses on data bias: lower bias \rightarrow better coverage of design space \rightarrow better generalizability of models

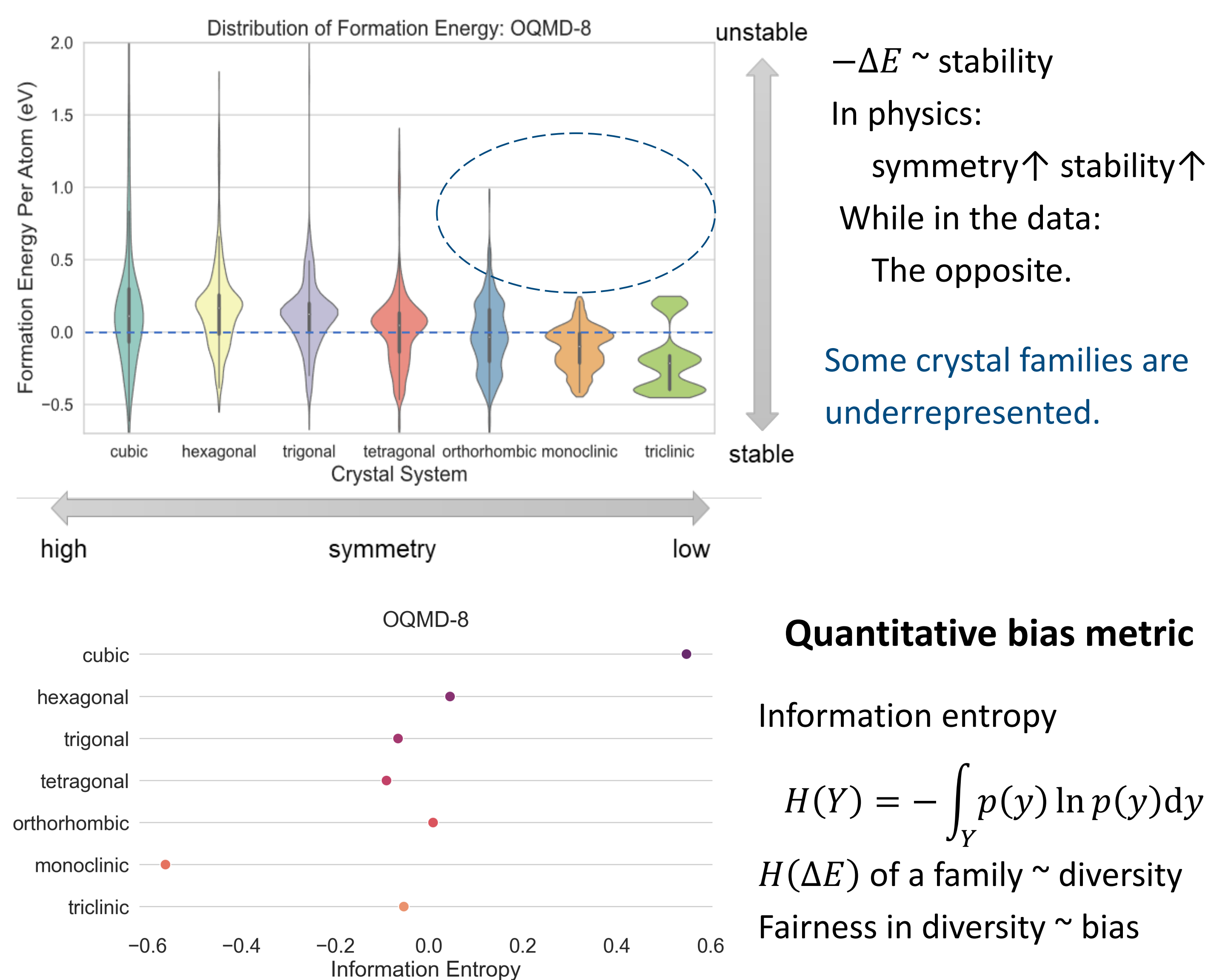
- A metric for measuring bias in materials data
- A framework guiding data acquisition to mitigate the bias

Problem Statement

Where materials informatics researchers get data

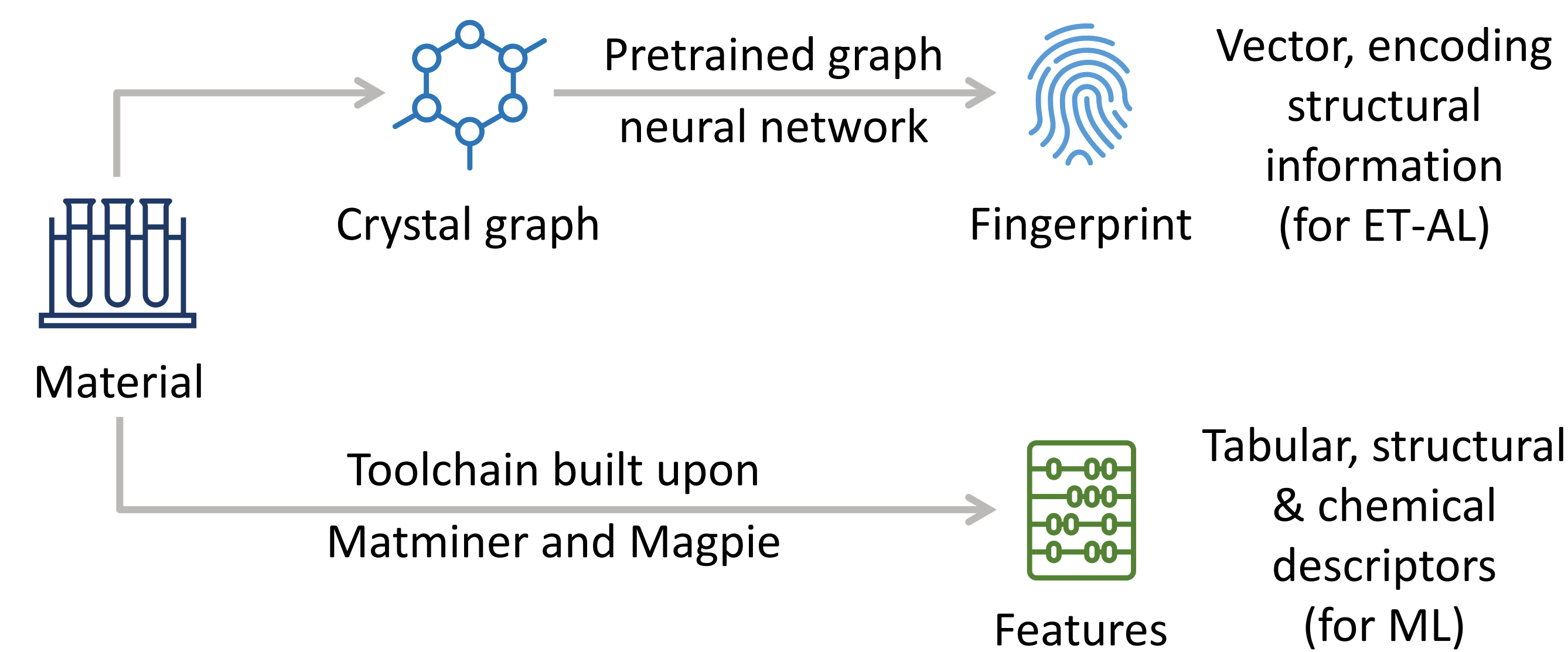


Example: structure–stability bias in OQMD

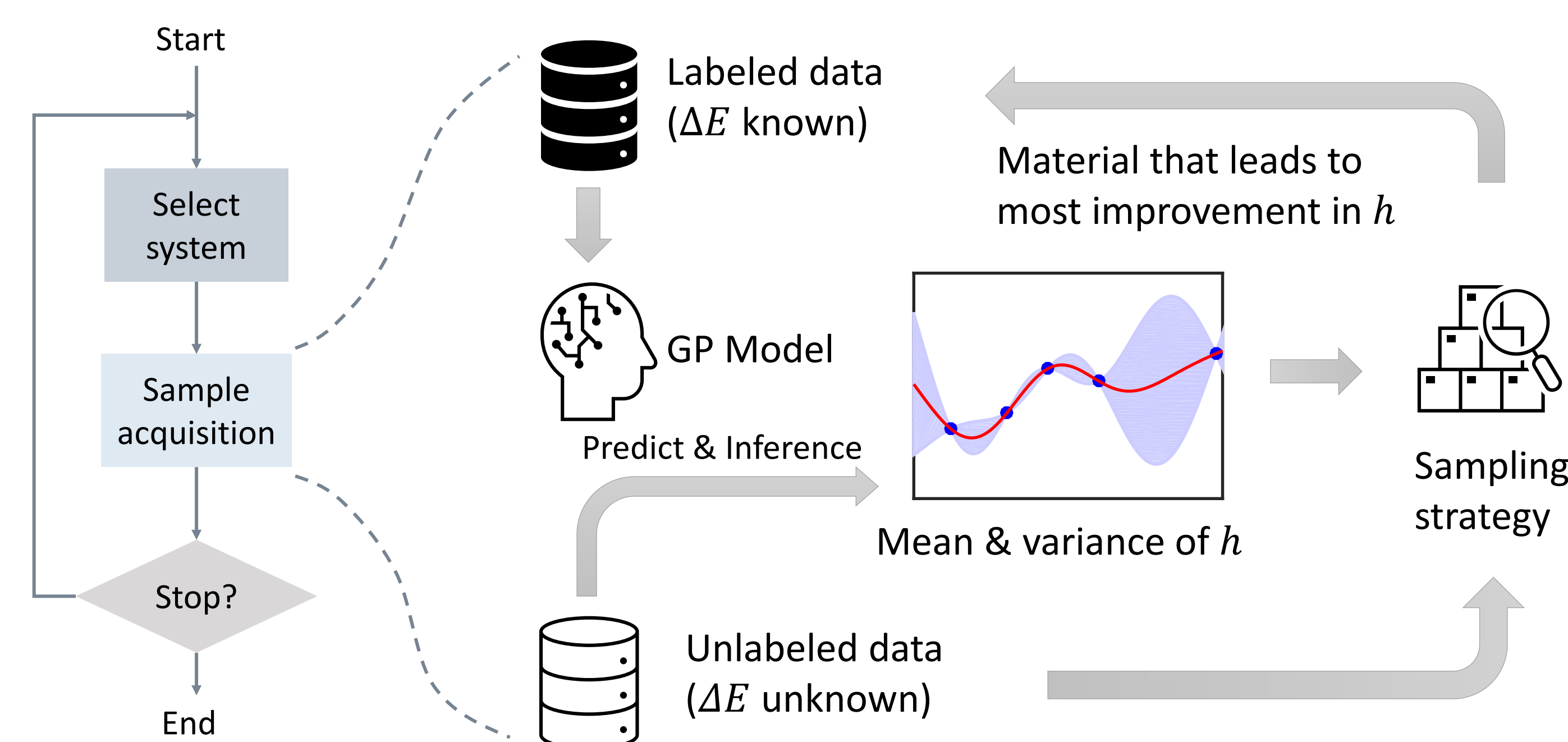


Methods

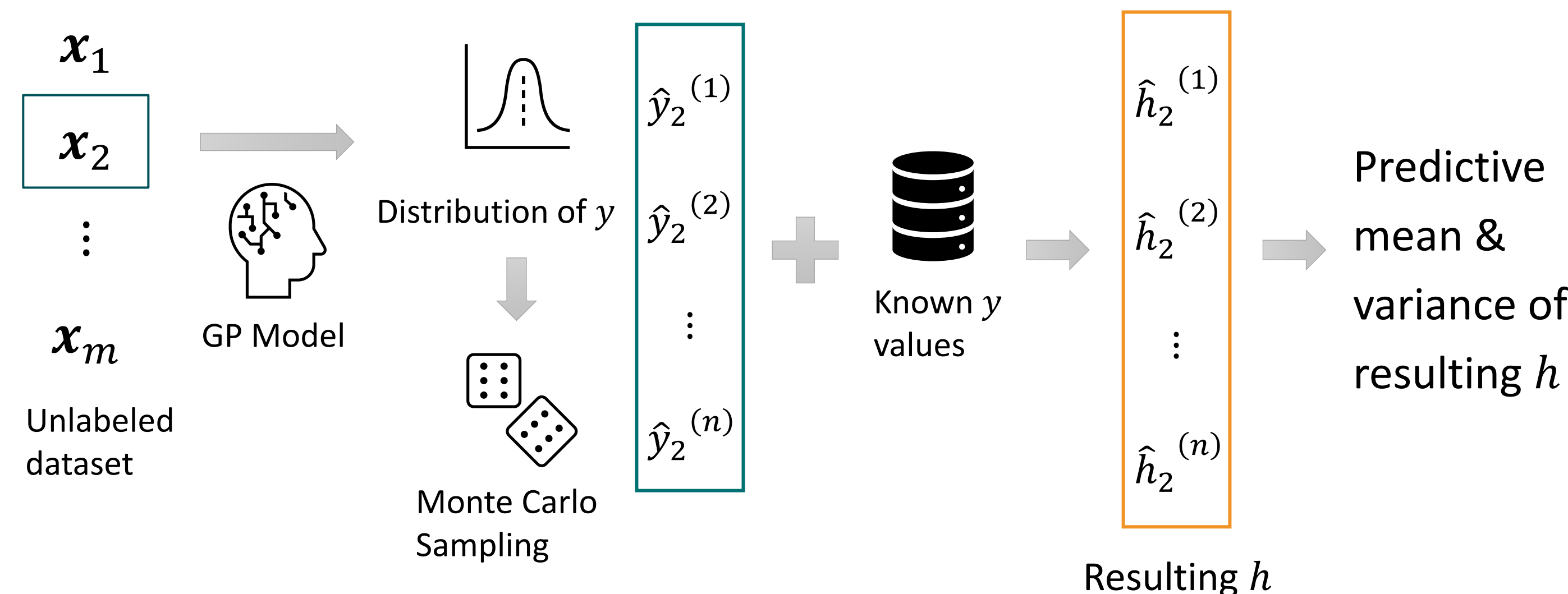
Representations of materials (input to the models)



ET-AL framework: maximize H of underrepresented systems



Monte Carlo inference for uncertainty in h



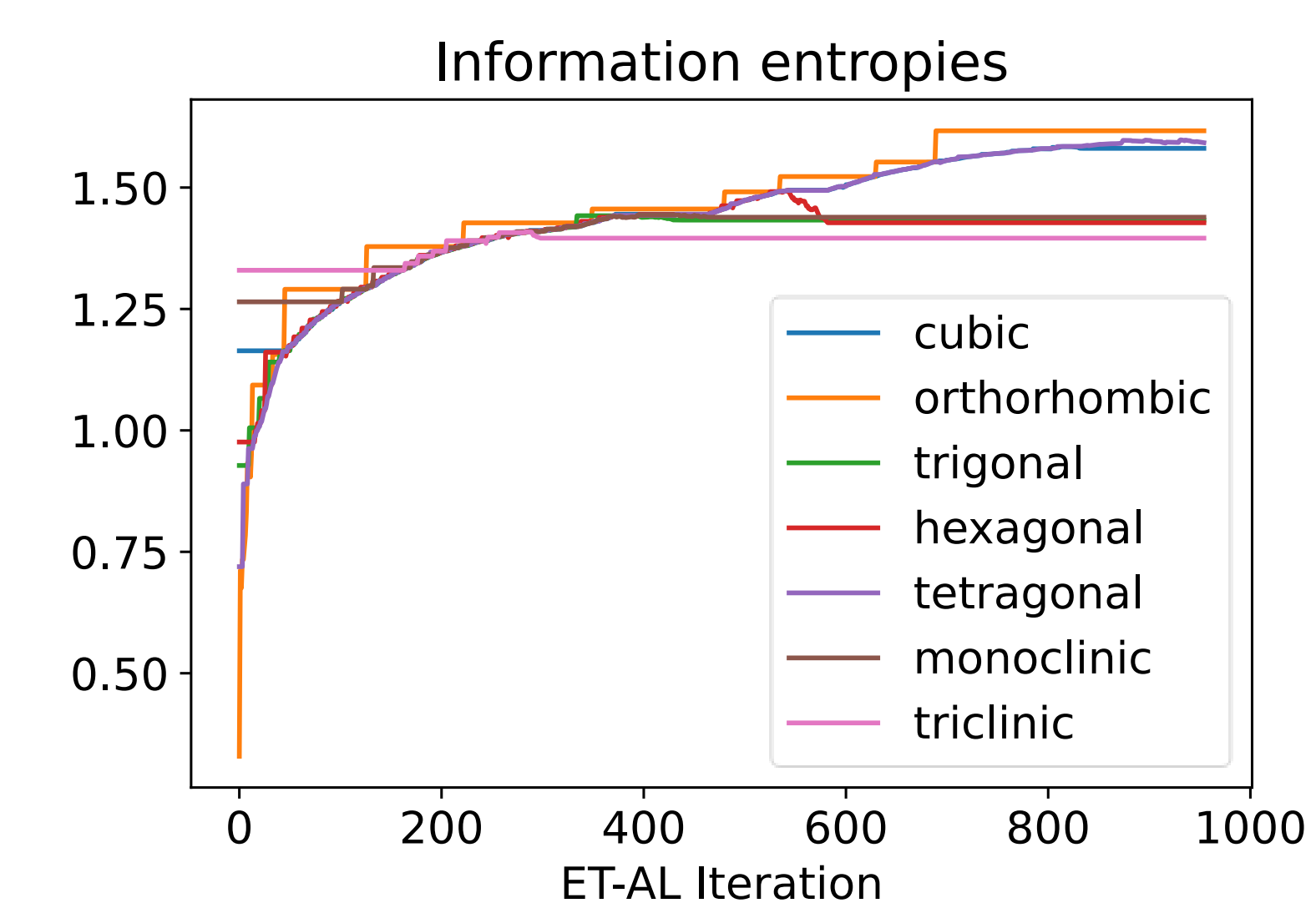
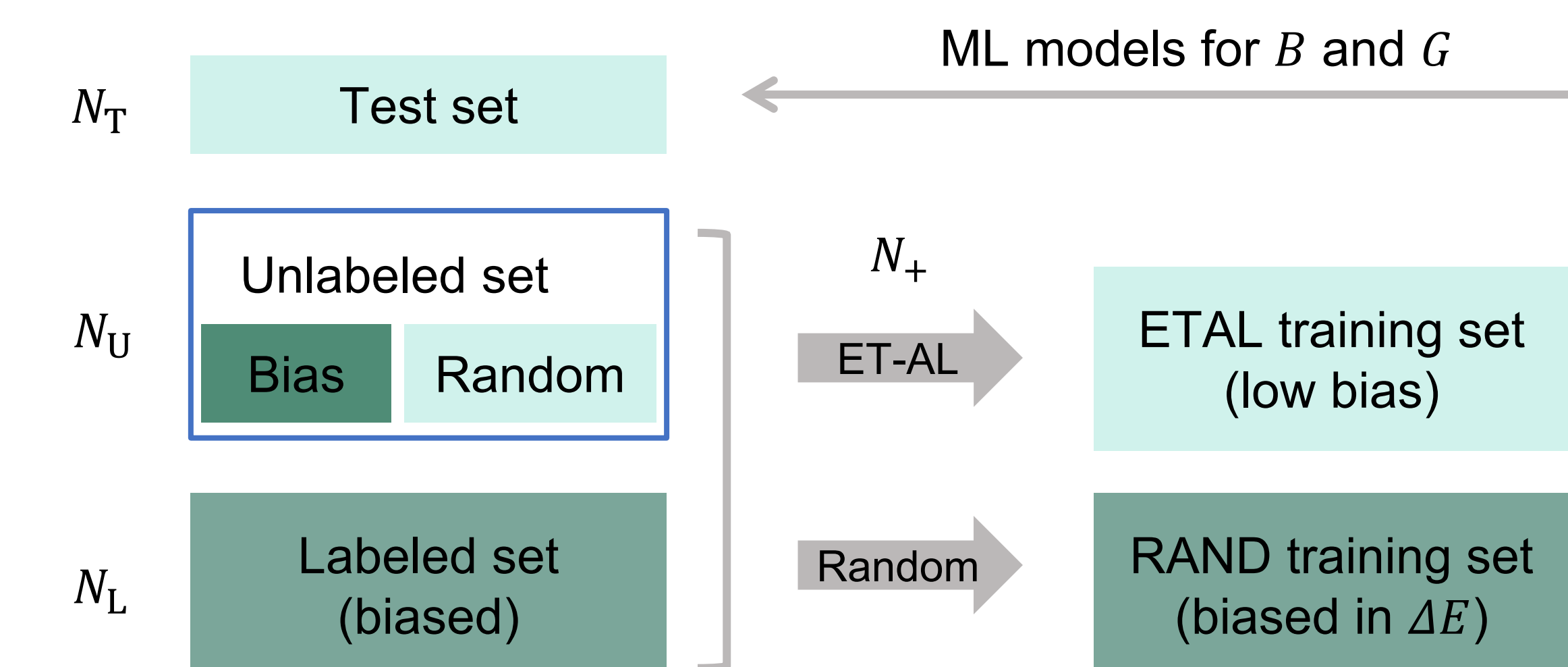
Outlook

ET-AL can mitigate data bias, thus benefiting data-driven models

- Help researchers improve the quality of datasets
- Guide the construction of materials data platforms
- Applicable to data-centric informatics in other scientific domains

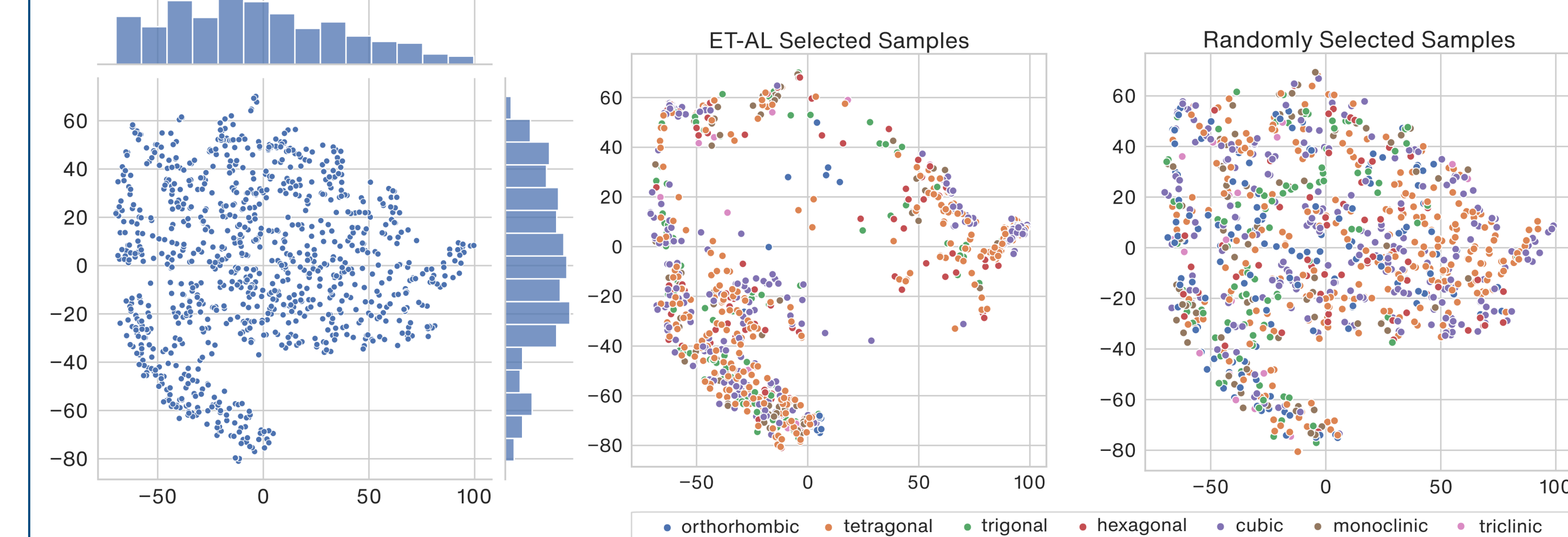
Results

Virtual experiments on Jarvis dataset (size $\sim 12K$)

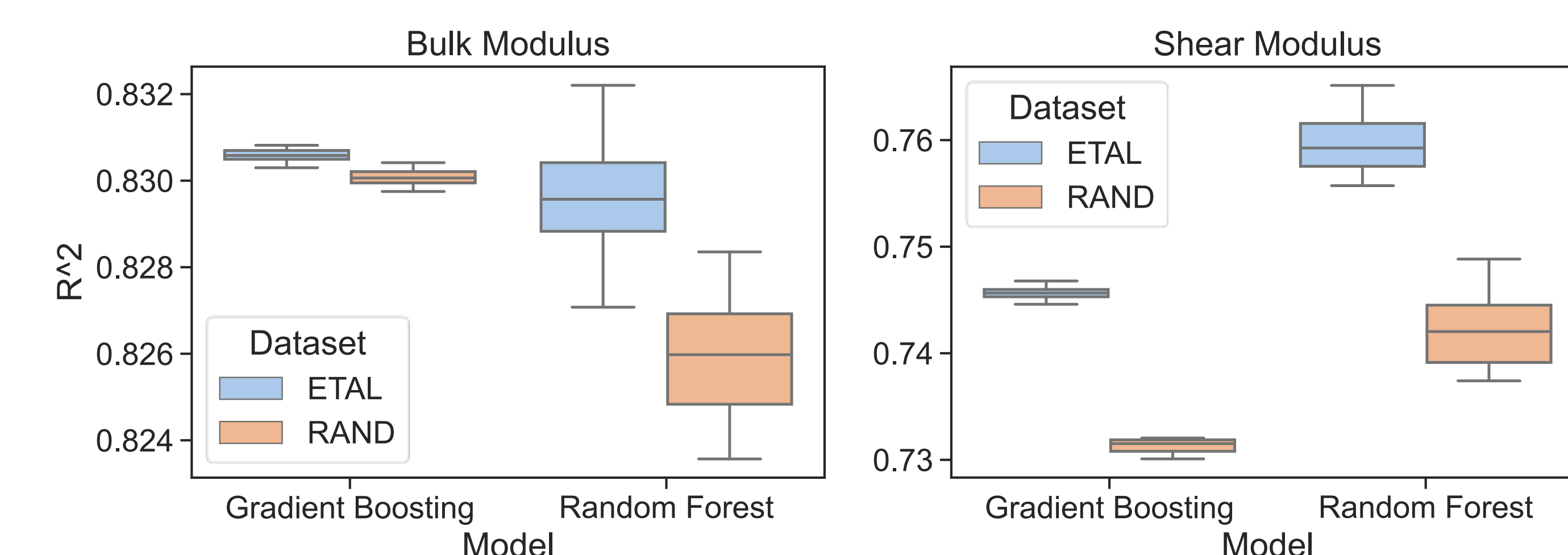


ET-AL successfully fixed the artificially created bias.

ET-AL selects samples in underrepresented regions



Reducing bias benefits ML modeling for materials properties.



Reference: H. Zhang, W. W. Chen, J. M. Rondinelli, & W. Chen. ET-AL: Entropy Targeted Active Learning for Bias Mitigation in Materials Data. *In preparation.*

