

Machine Learning in Molecular Sciences

My Background



2005 - 2009
B.S.
Chemical Engineering



Northwestern
University

2010 - 2015
Ph.D.
Chemical Engineering



PRITZKER SCHOOL OF
MOLECULAR ENGINEERING
THE UNIVERSITY OF CHICAGO

2015 - 2018
Postdoctoral Training
Molecular Eng.



2019 - present
Assistant Professor
Chemical & Biomol. Eng.

Undergraduate Research at ND



Fulbright Scholarship, Spain, 2009



Born and raised in Puerto Rico



Research: Computational Molecular Engineering

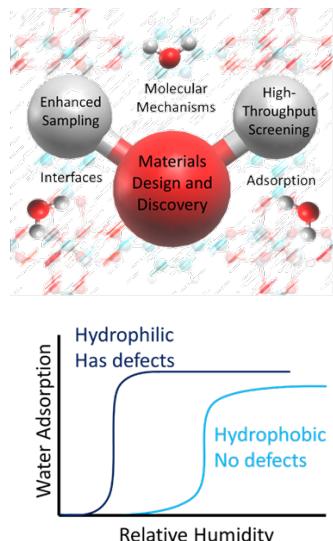
Research Methods

- Molecular Modeling
- Statistical Mechanics
- Materials Science
- Machine Learning and Data Science

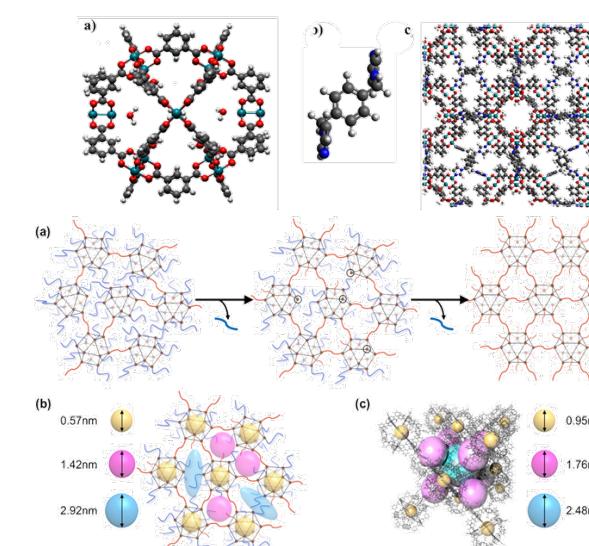
Application Areas

- Energy Storage
- Gas adsorption and Separations
- Molecular Design of Porous Materials
- Water farming, Water security
- Physical Properties of Novel Materials
- Materials Design and Discovery
- Self-Assembly of Porous Materials
- Interfacial Engineering

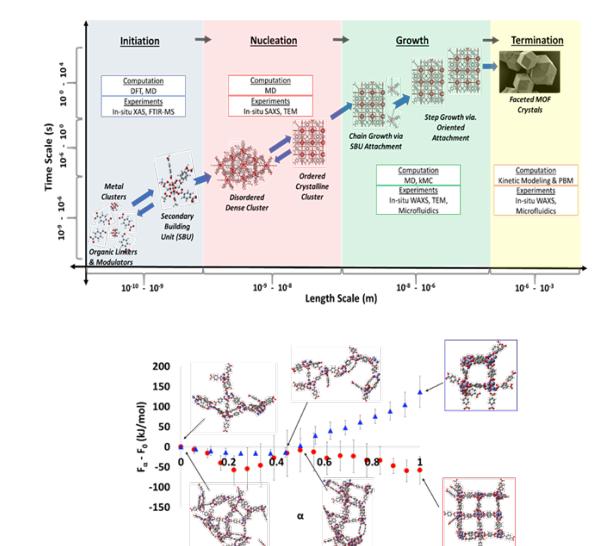
Water-Energy Nexus



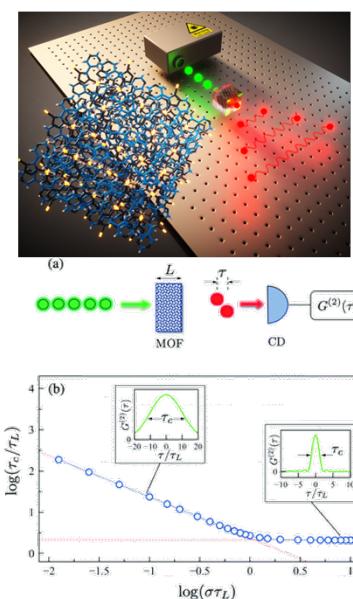
Soft Porous Coordination Polymers



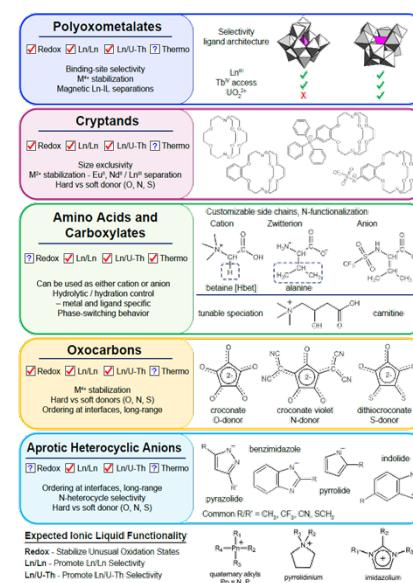
Self-Assembly of Porous Materials



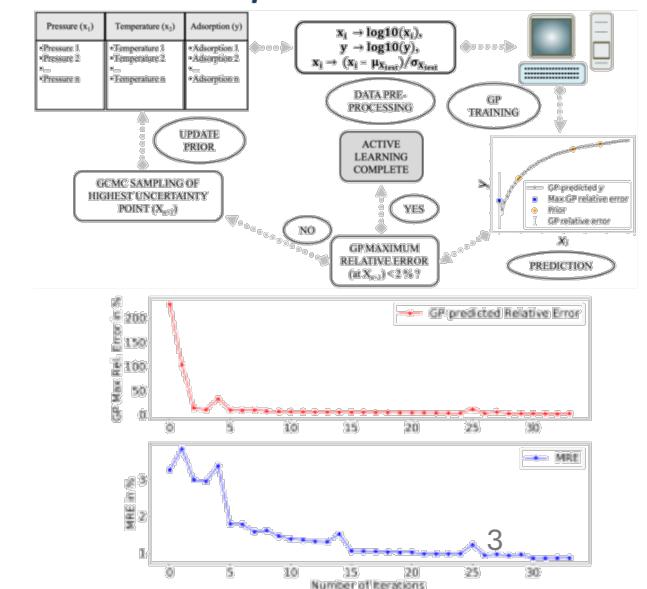
Quantum Technologies



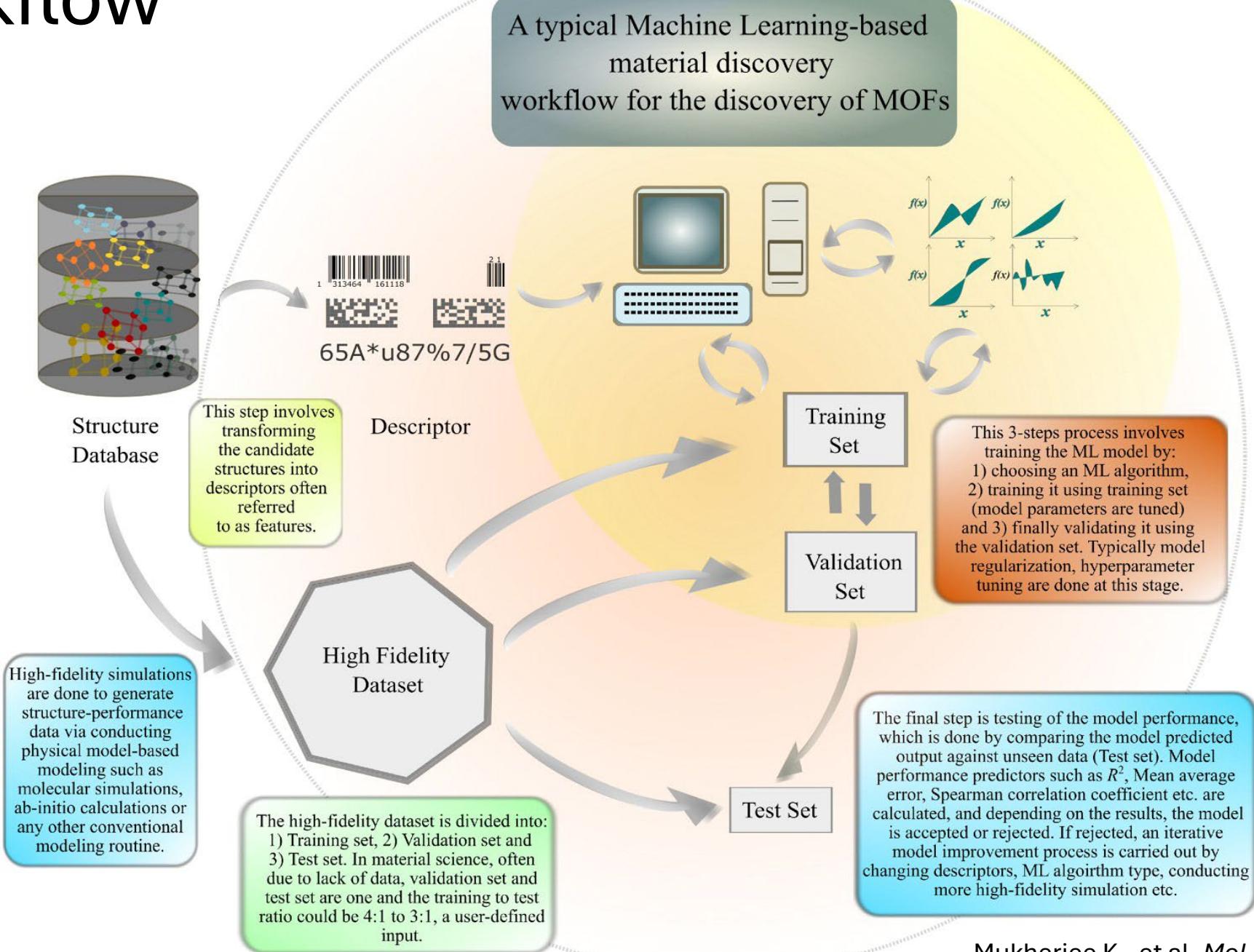
ILs for REE separations



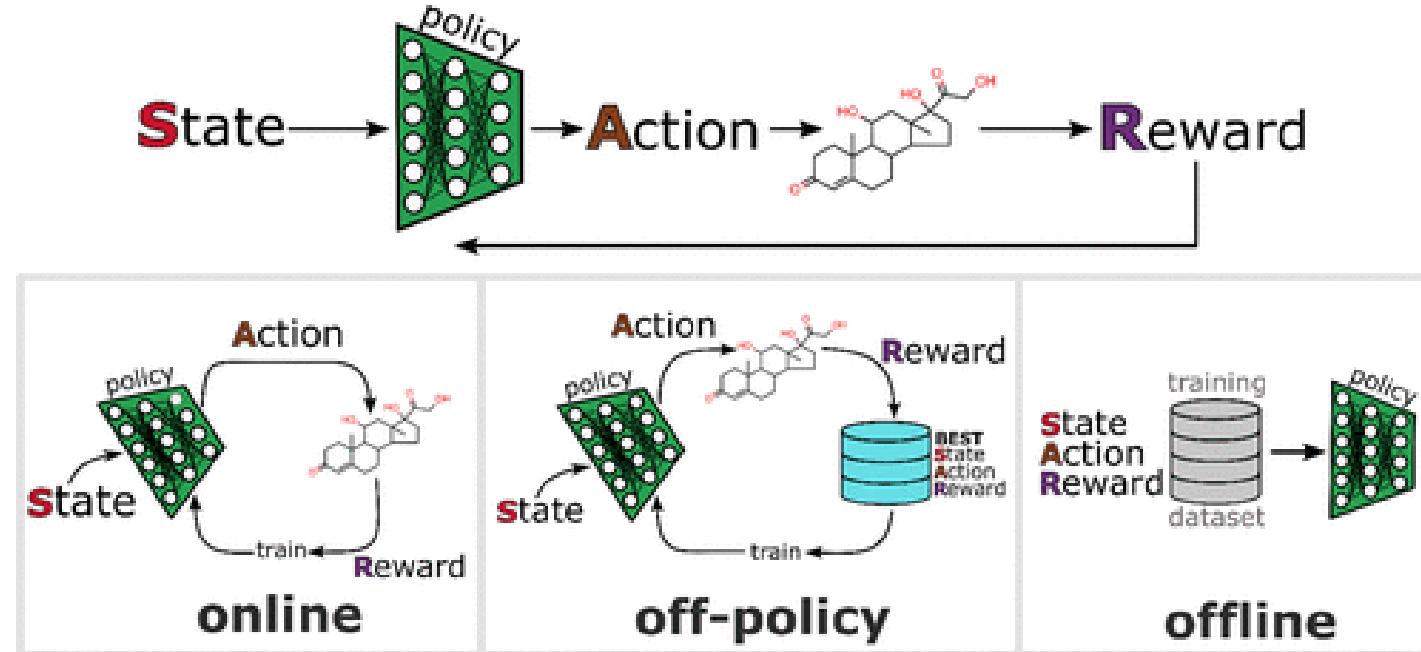
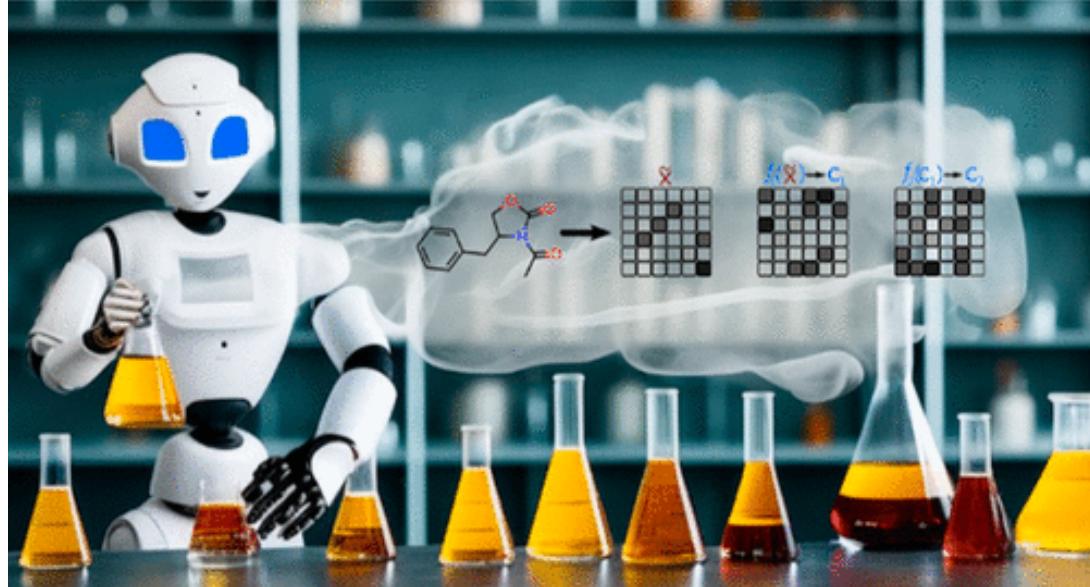
Machine Learning for Phase Equilibria Predictions



ML workflow

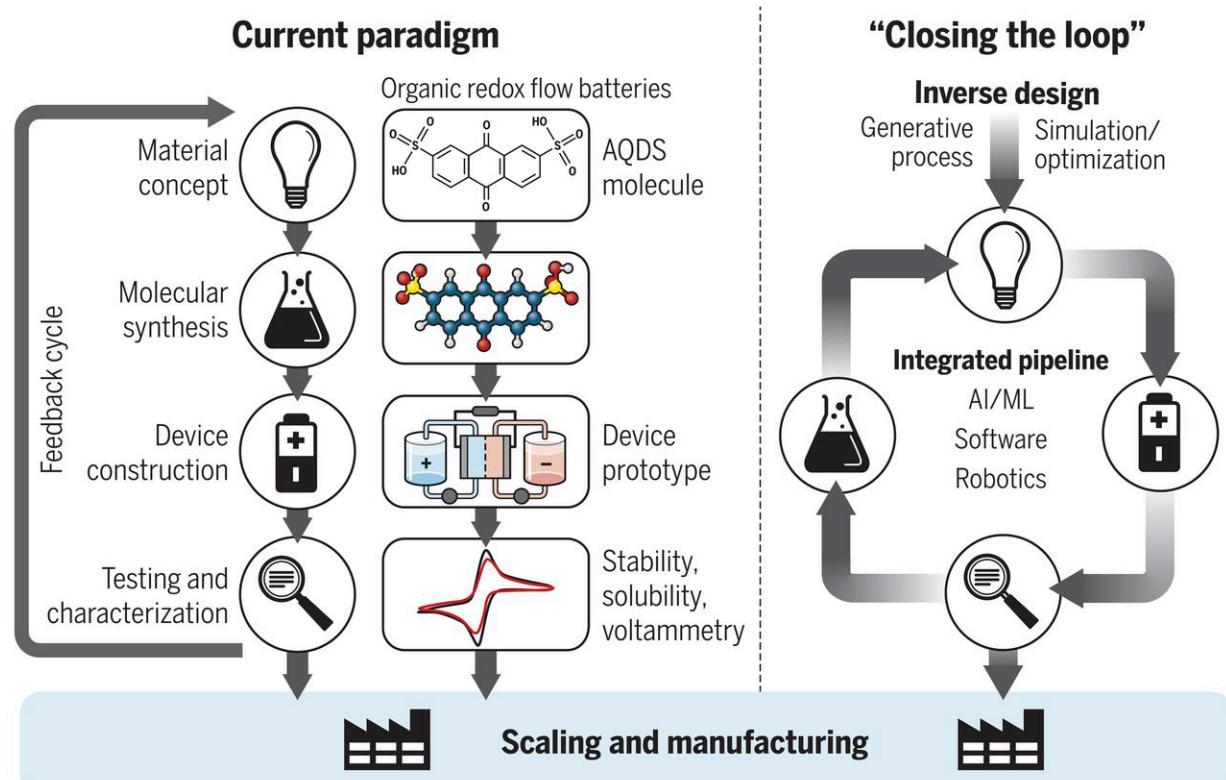
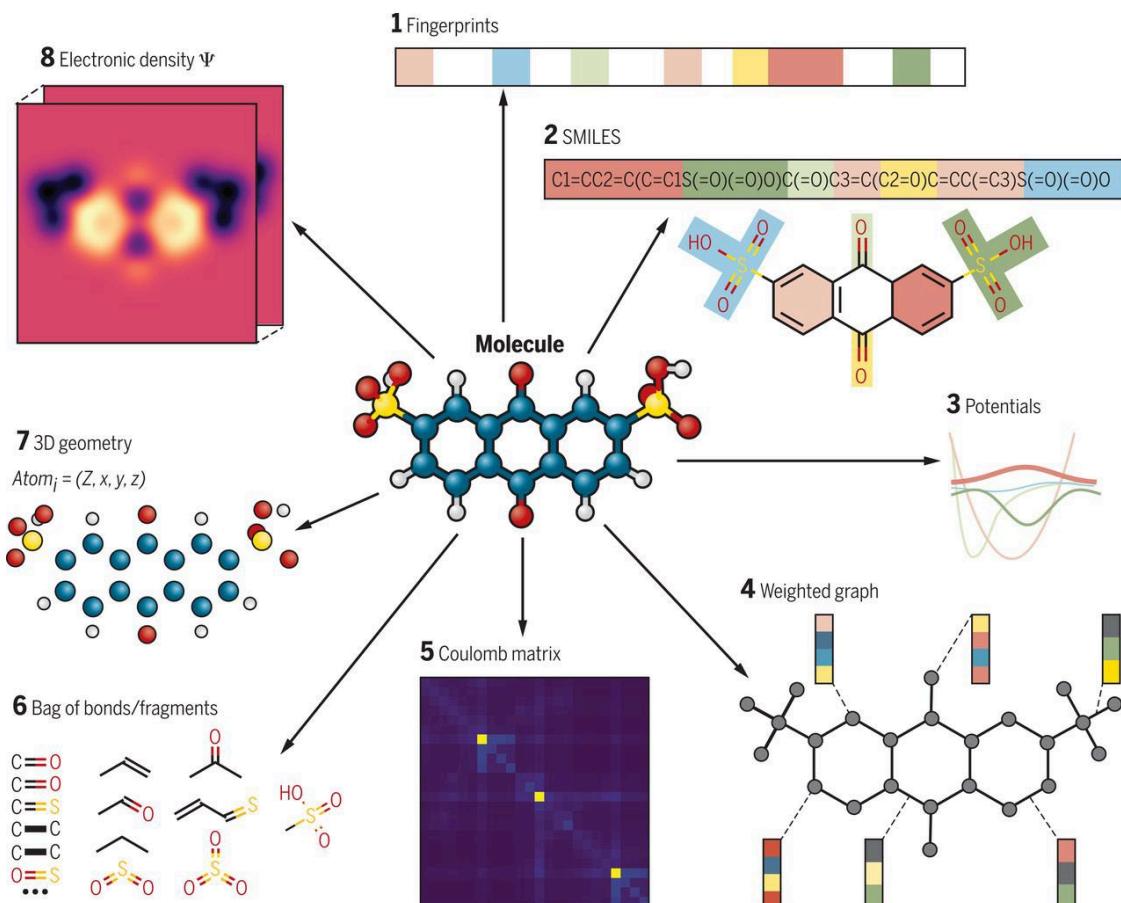


Generation of Chemistry



Anstine, D. M.; Isayev, O., Generative Models as an Emerging Paradigm in the Chemical Sciences. *Journal of the American Chemical Society* **2023**, 145 (16), 8736-8750.

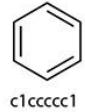
Molecular Descriptors and Inverse Design



Digital Chemical Spaces

(a)

SMILES input

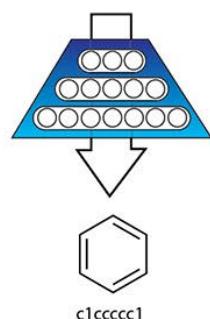


(b)

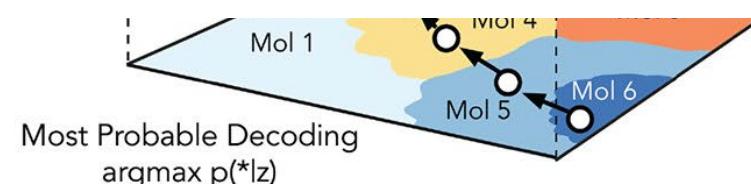


Two autoencoder systems were trained: one with 108 000 molecules from the QM9 data set of molecules with fewer than 9 heavy atoms (31) and another with 250 000 drug-like commercially available molecules extracted at random from the ZINC database. (32) We performed random optimization over hyperparameters specifying the deep autoencoder architecture and training, such as the choice between a recurrent or convolutional encoder, the number of hidden layers, layer sizes, regularization, and learning rates. The latent space representations for the QM9 and ZINC data sets had 156 dimensions and 196 dimensions, respectively.

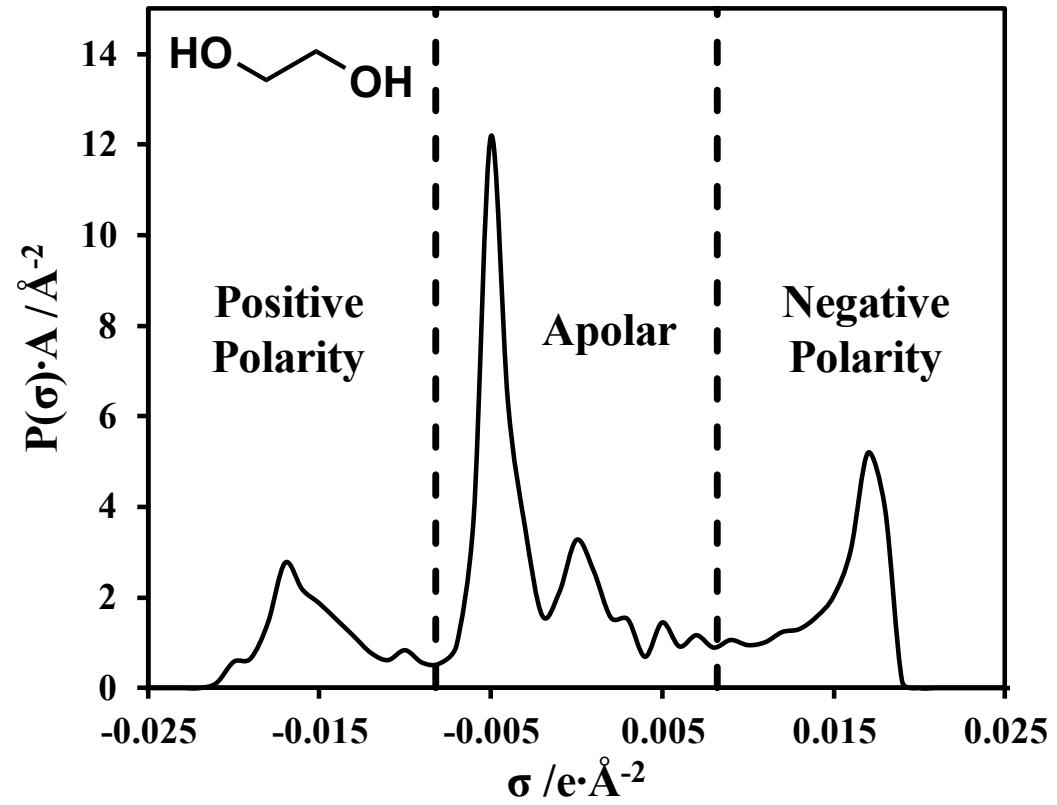
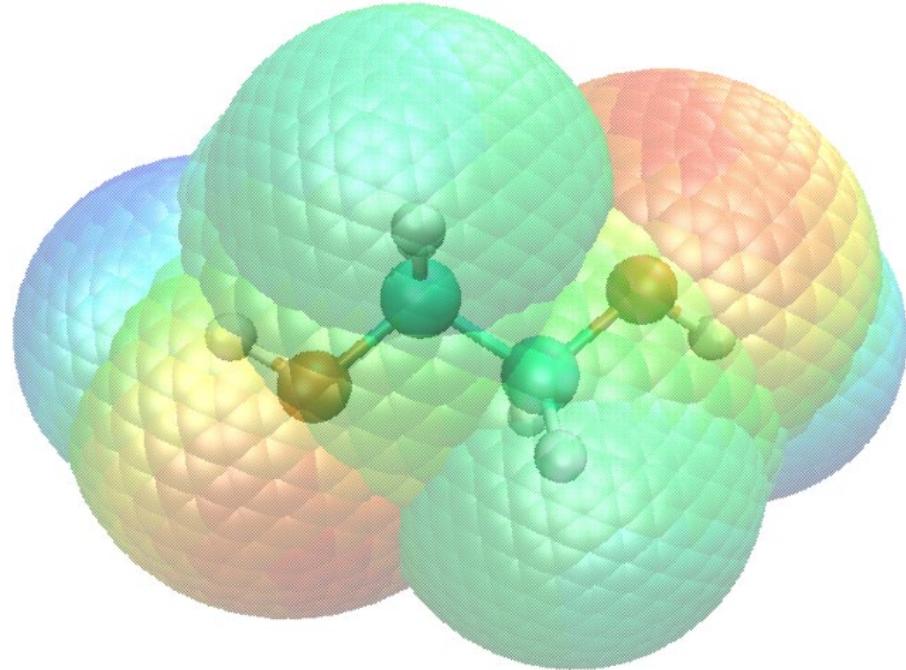
DECODER
Neural Network



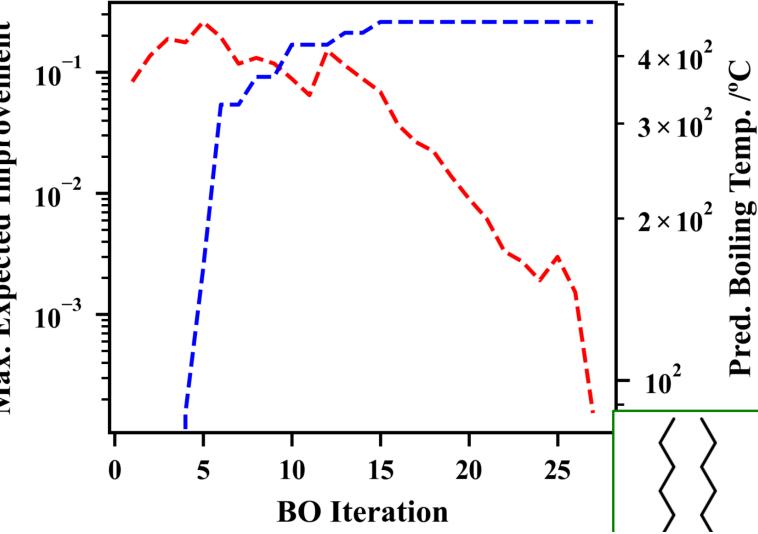
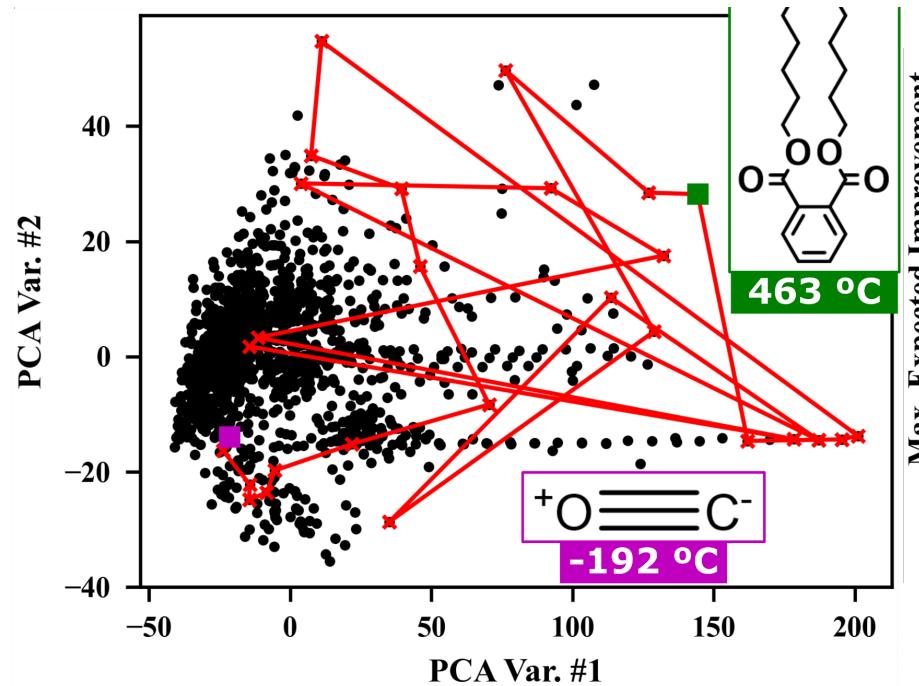
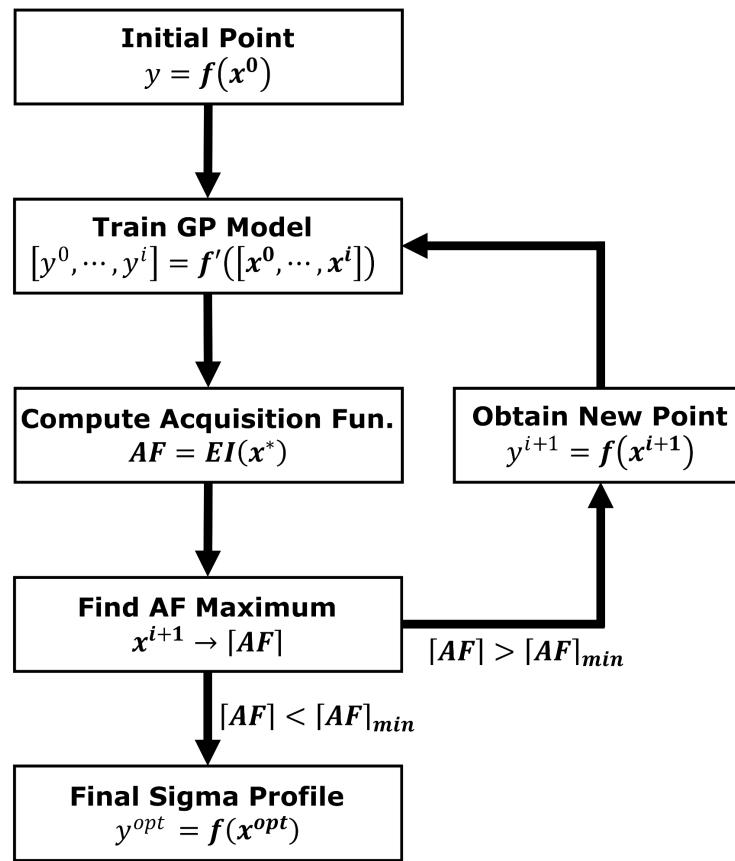
SMILES output



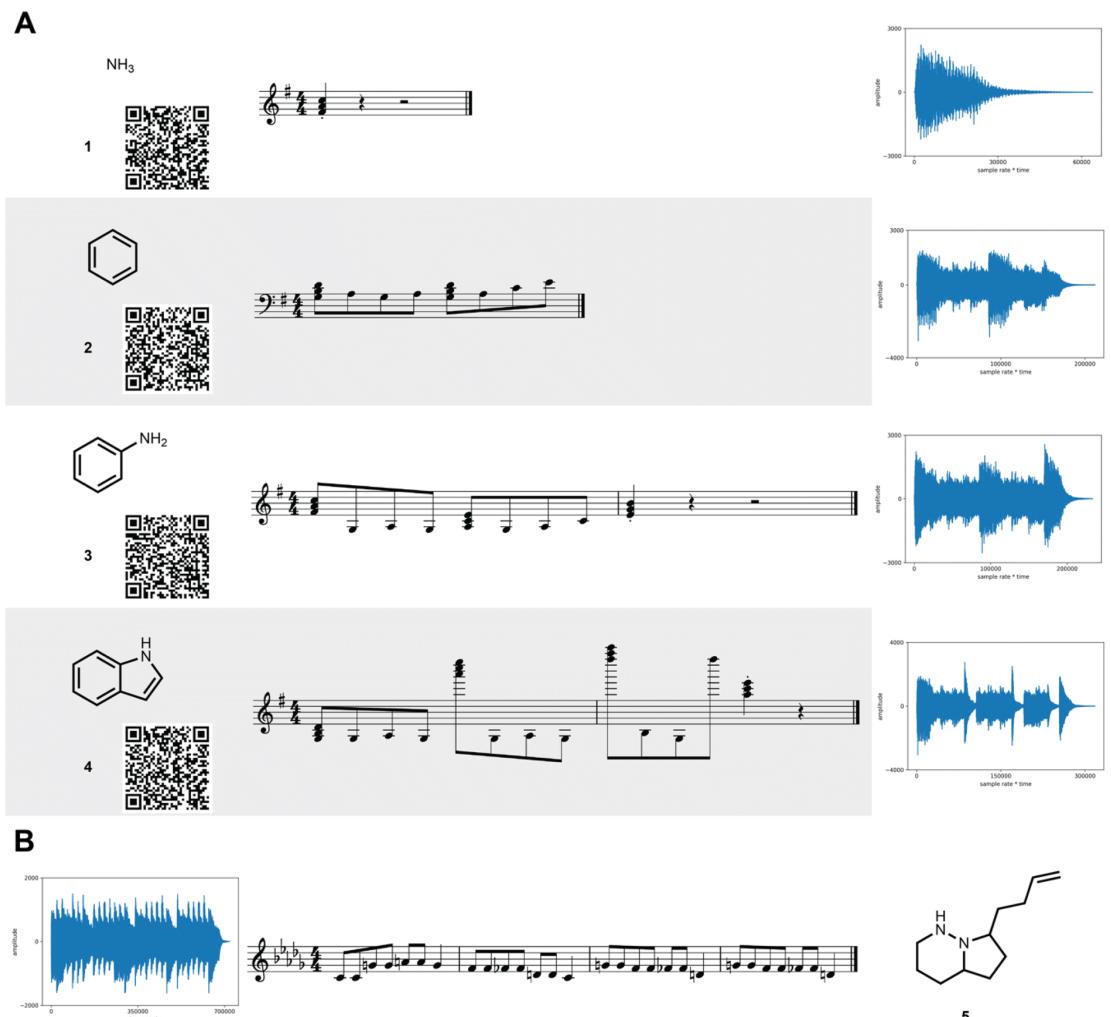
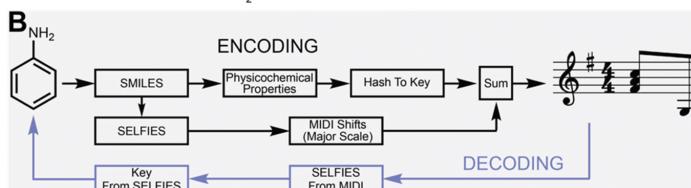
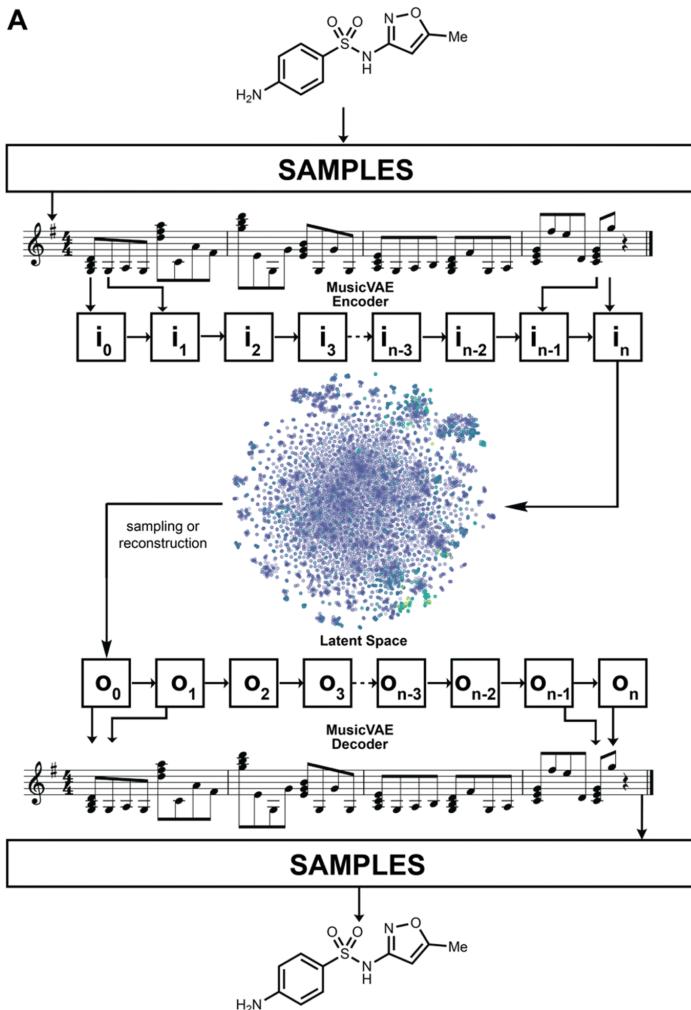
New Chemical Descriptors



Navigating Chemical Digital Spaces



Molecular Sonification



Mahjour, B.; Bench, J.; Zhang, R.; Frazier, J.; Cernak, T., Molecular sonification for molecule to music information transfer. *Digital Discovery* **2023**, 2 (2), 520-530.

Featurization

Physics-Inspired Structural Representations for Molecules and Materials

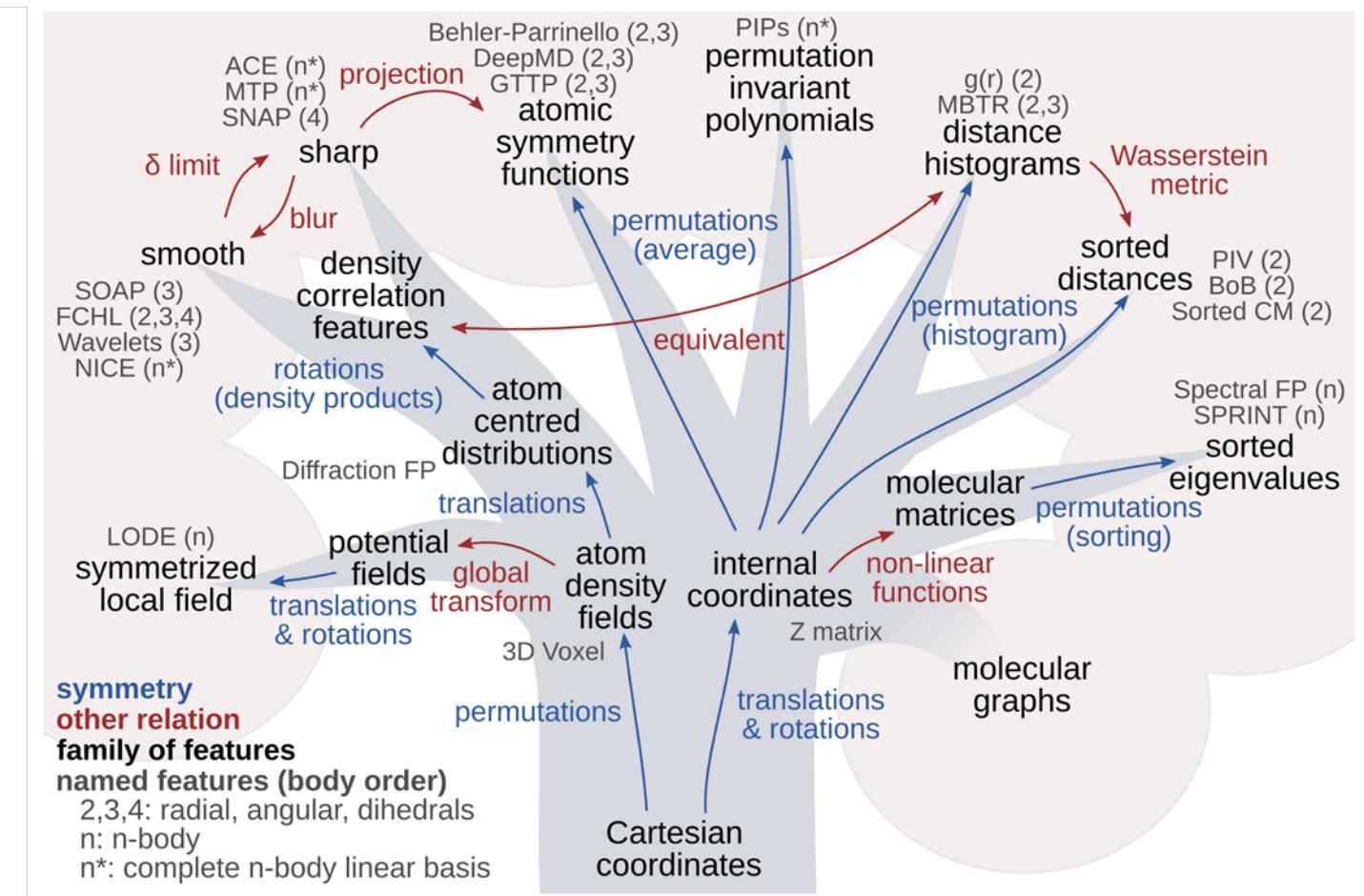
Felix Musil, Andrea Grisafi, Albert P. Bartók, Christoph Ortner, Gábor Csányi, and Michele Ceriotti*



Cite This: *Chem. Rev.* 2021, 121, 9759–9815

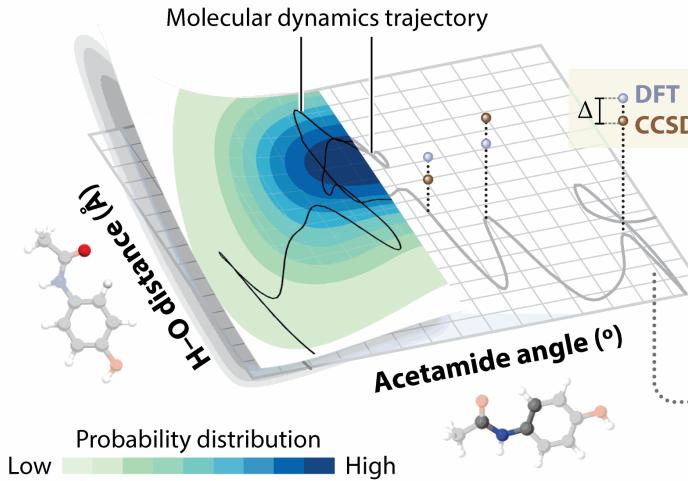


Read Online

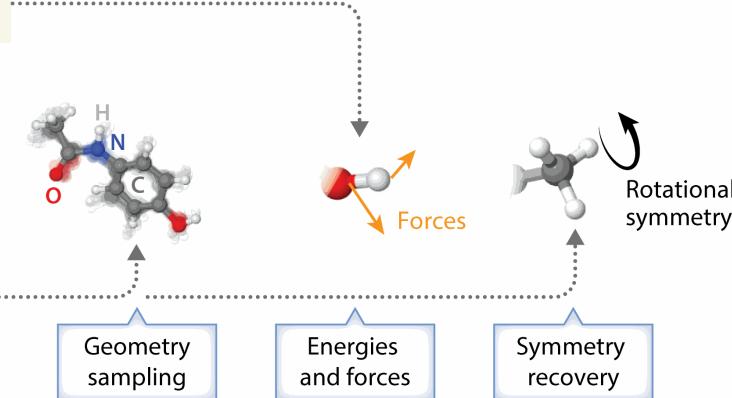


Molecular Simulations

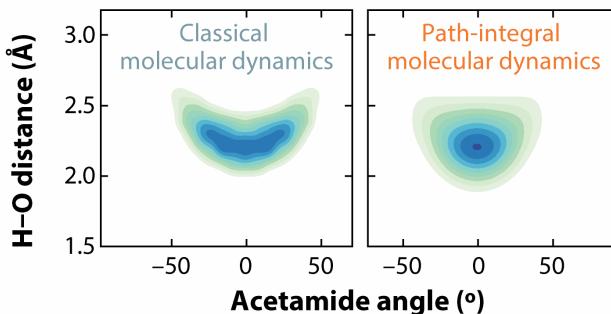
a Molecular dynamics simulation



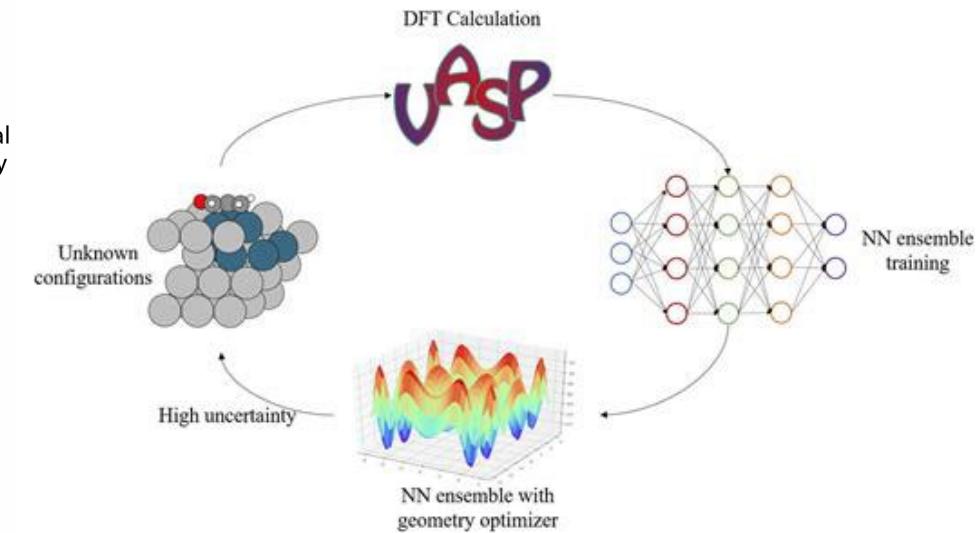
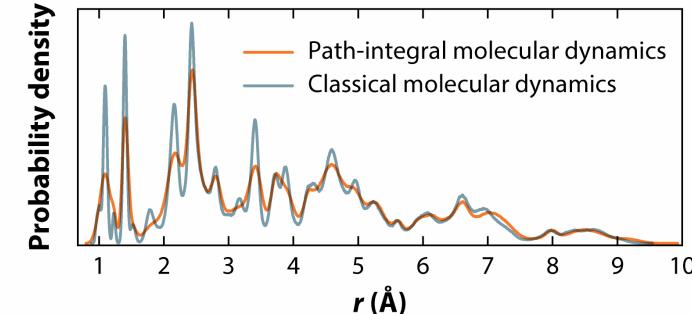
b Constructing machine learning force fields



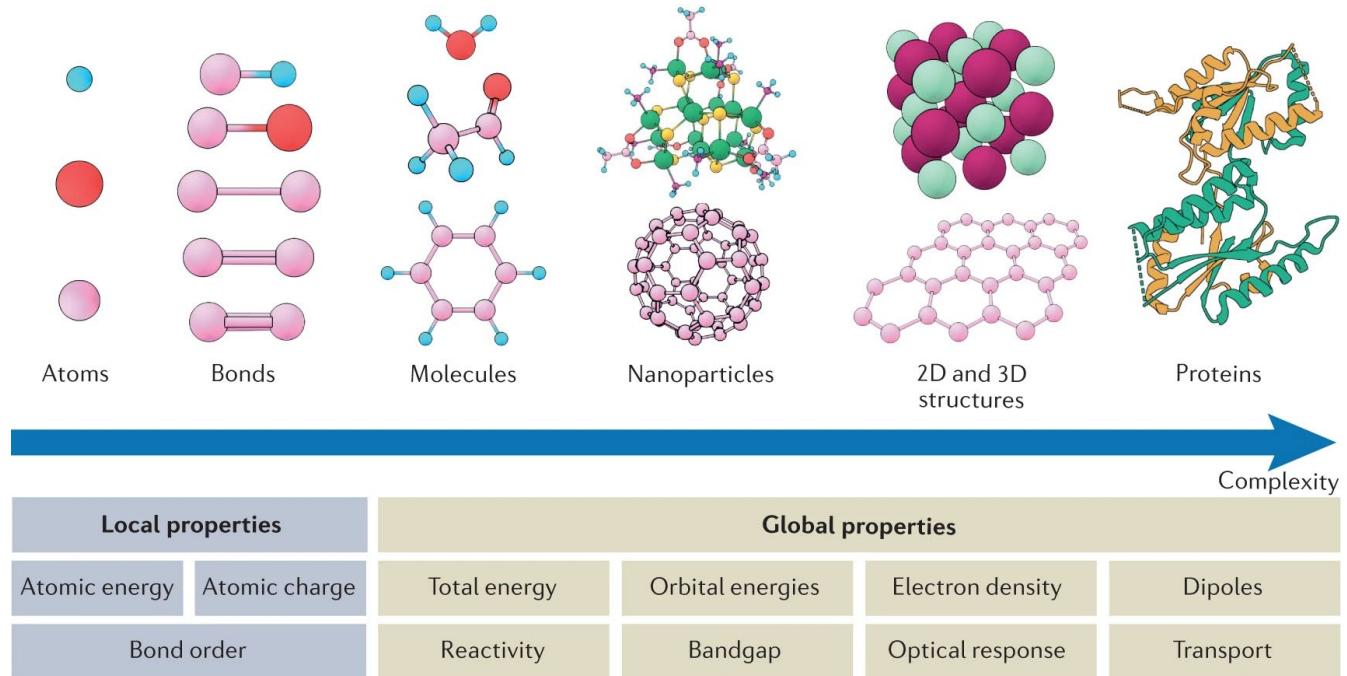
c Applications: free energy surfaces



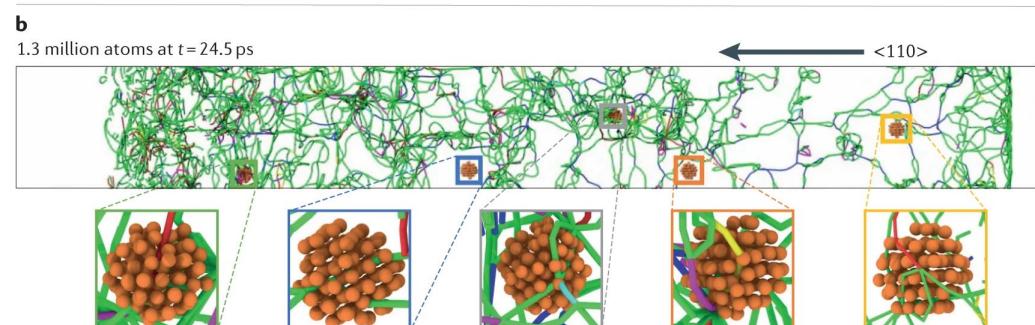
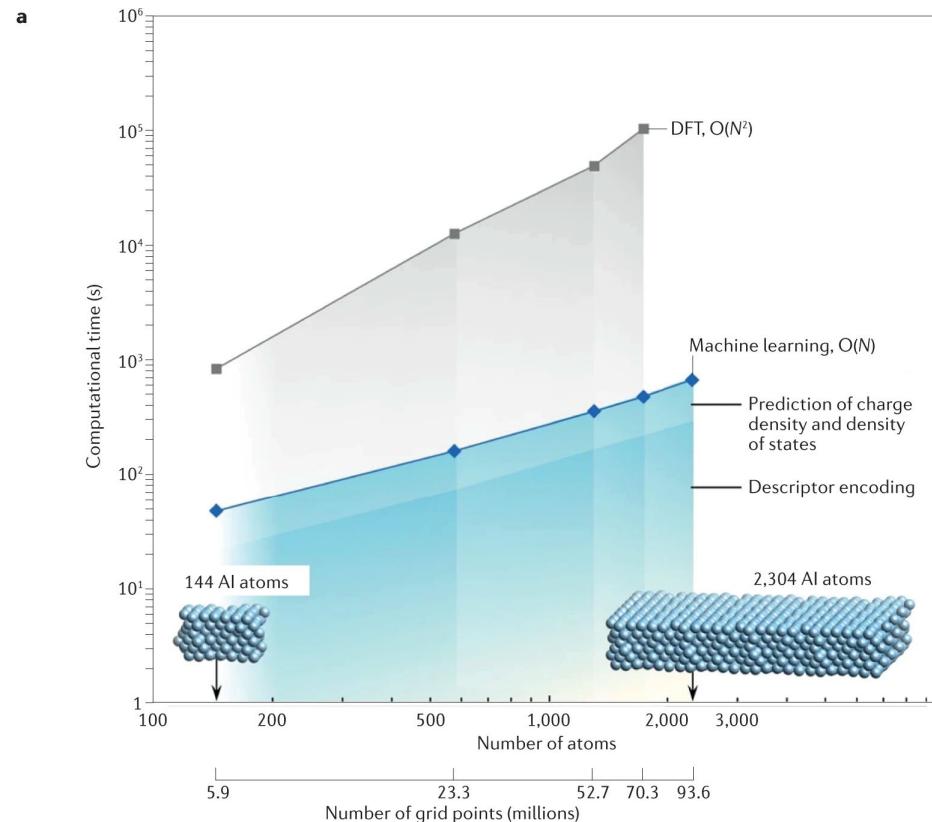
d Interatomic distance distribution



Molecular Simulations

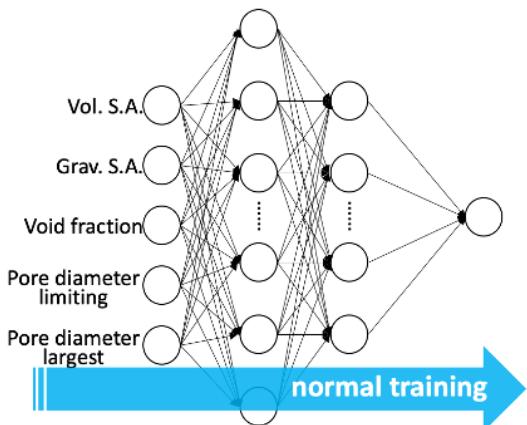


Fedik, N.; Zubatyuk, R.; Kulichenko, M.; Lubbers, N.; Smith, J. S.; Nebgen, B.; Messerly, R.; Li, Y. W.; Boldyrev, A. I.; Barros, K.; Isayev, O.; Tretiak, S., Extending machine learning beyond interatomic potentials for predicting molecular properties. *Nature Reviews Chemistry* **2022**, 6 (9), 653-672.

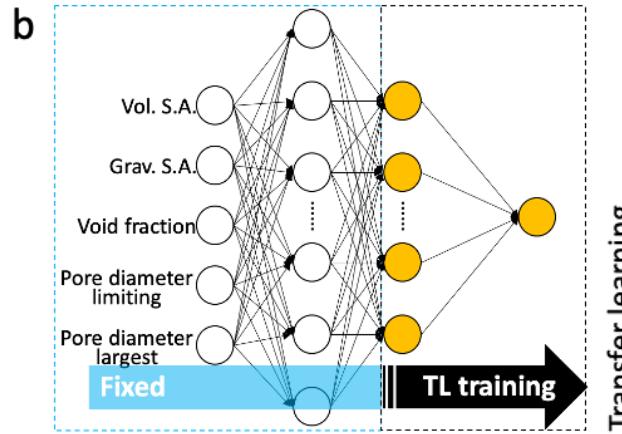


Transfer Learning

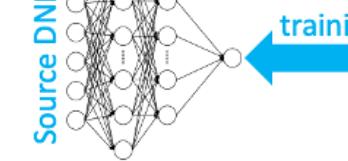
a



b



c Source DNN



Source property: all 13,506 data



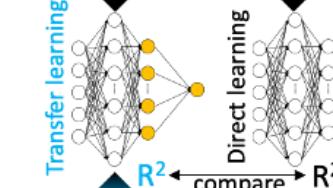
Target property: all 13,506 data



Batch 1

100 random data

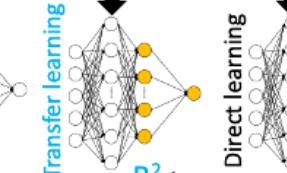
Transfer learning



Batch 2

100 random data

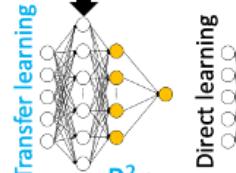
Transfer learning



Batch 1,000

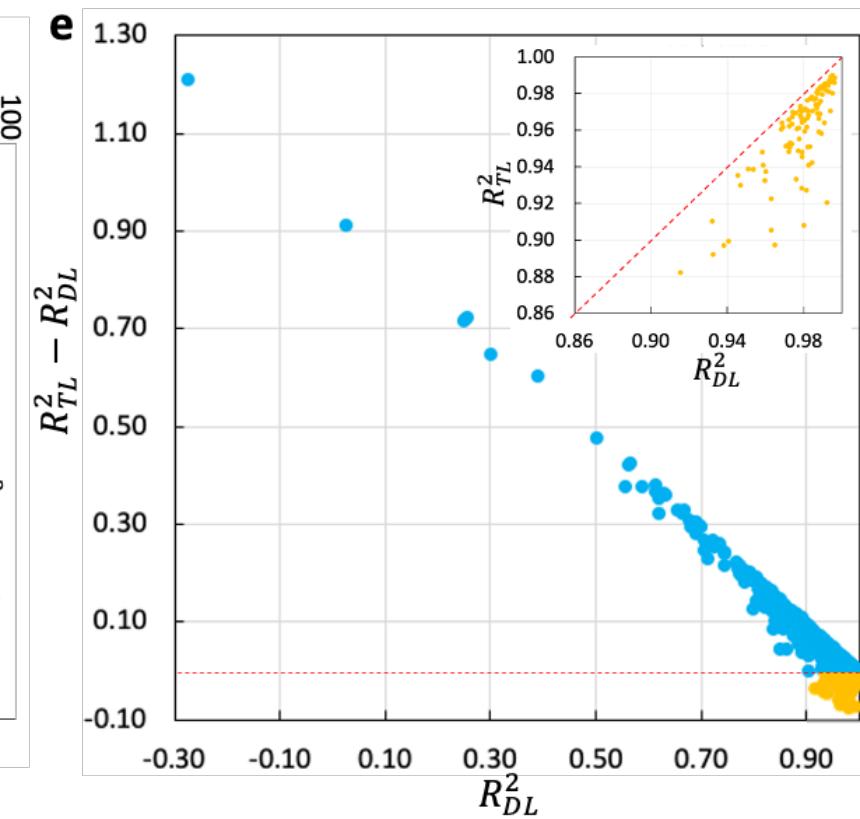
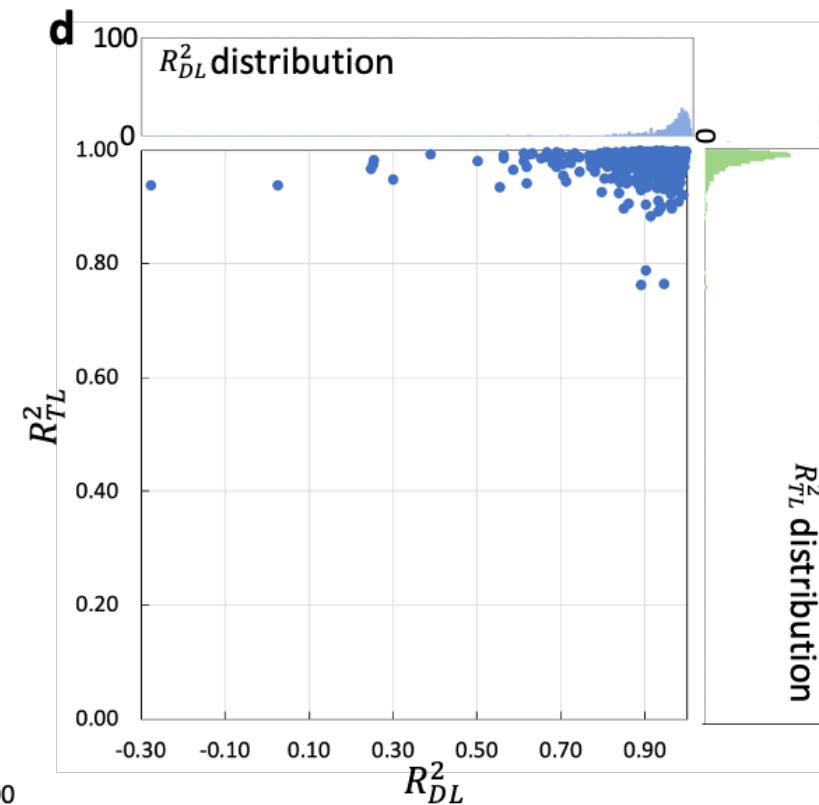
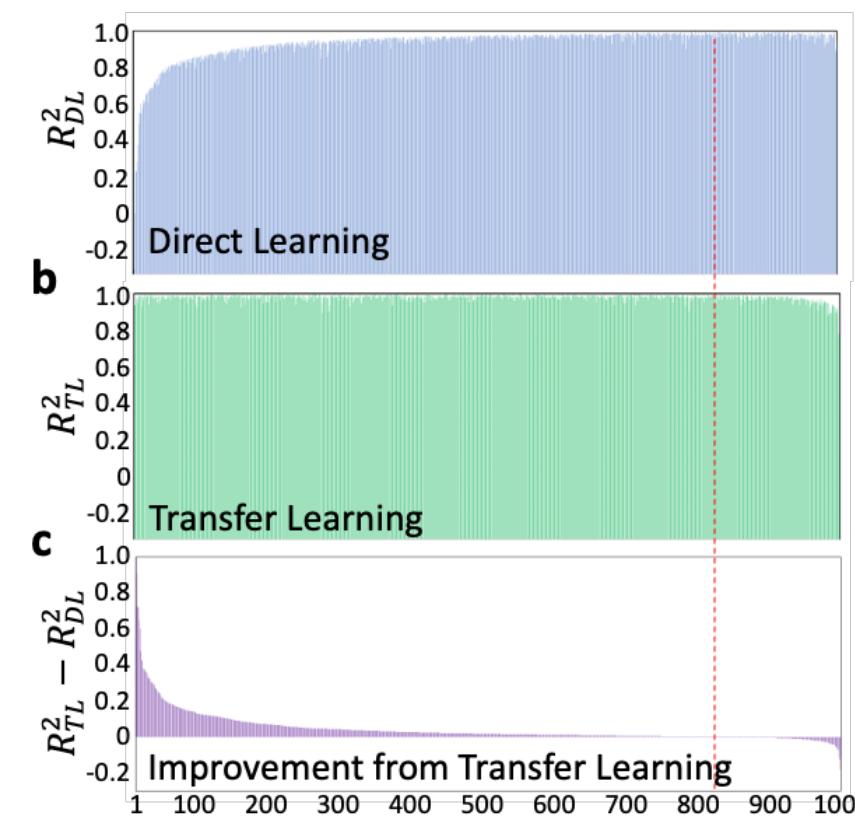
100 random data

Transfer learning

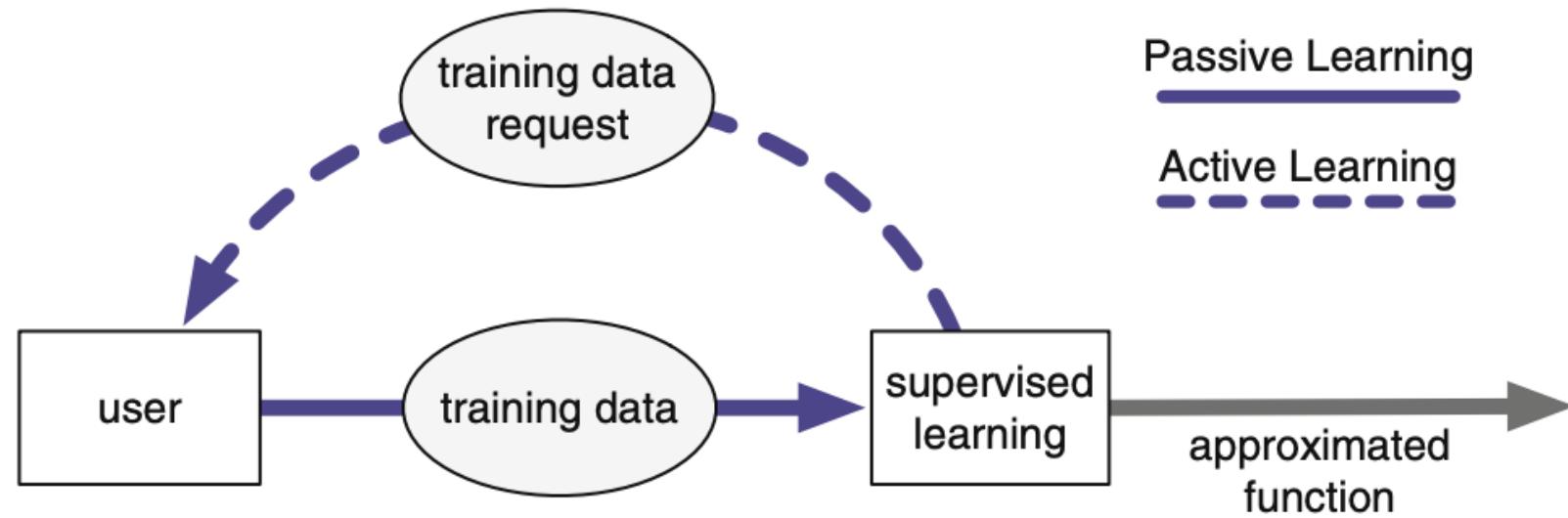


TL from hydrogen to methane

H₂ (100 bar, 243 K) to CH₄ (100 bar, 298 K)

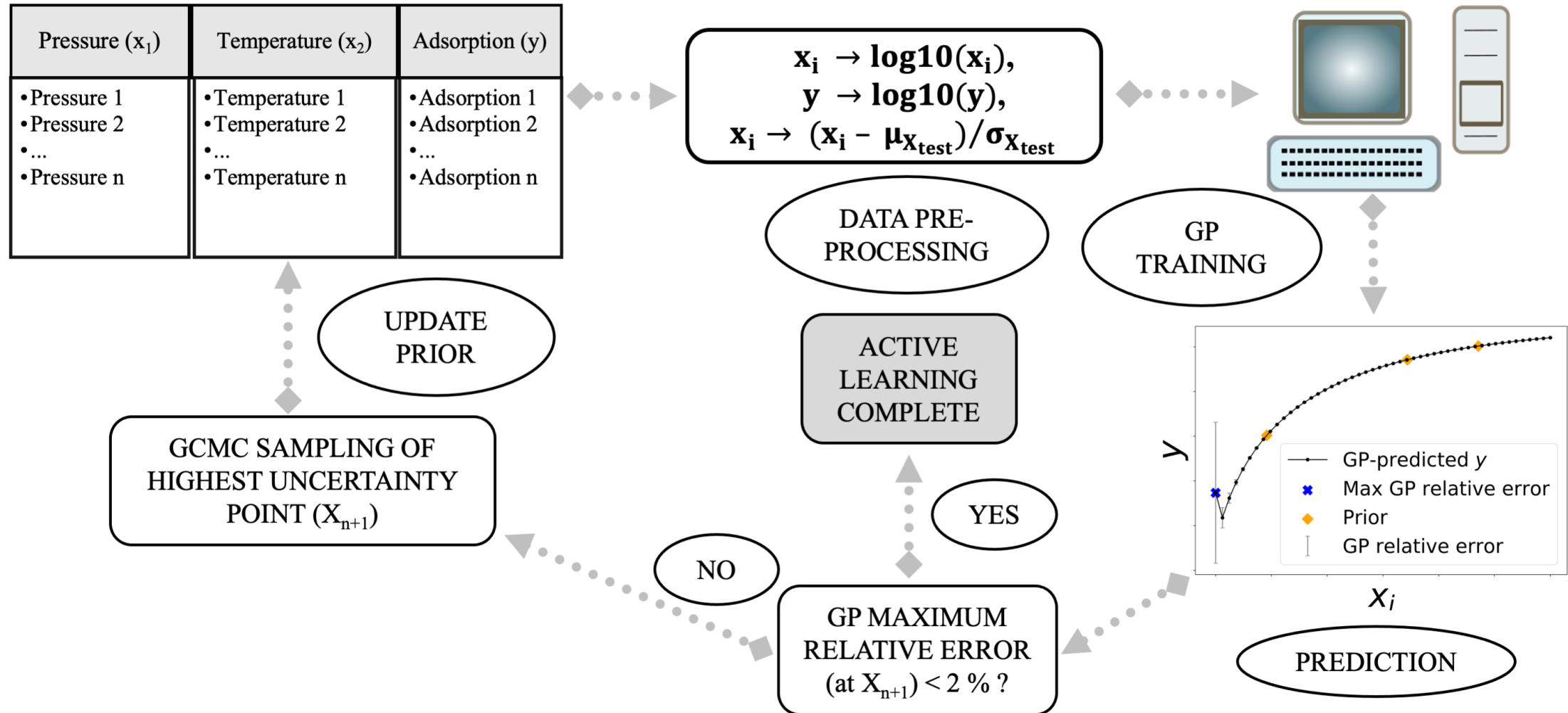


Active Learning

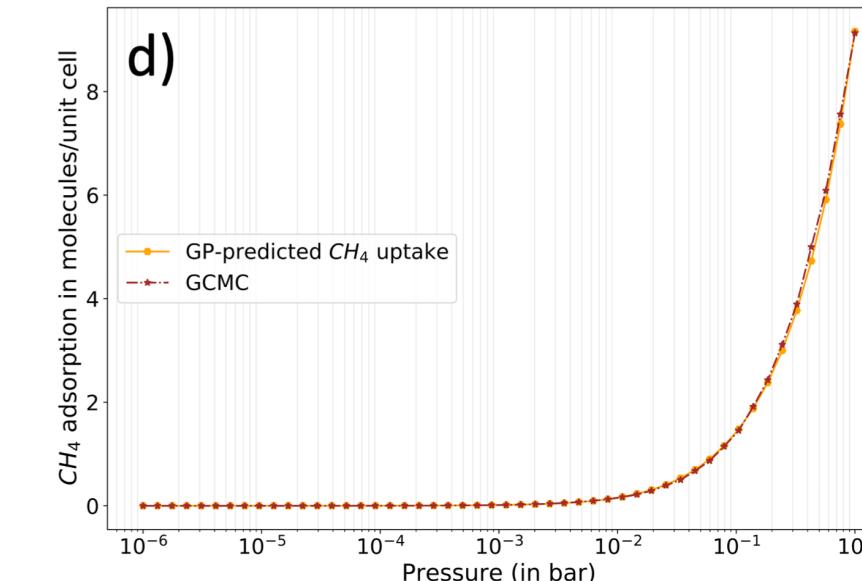
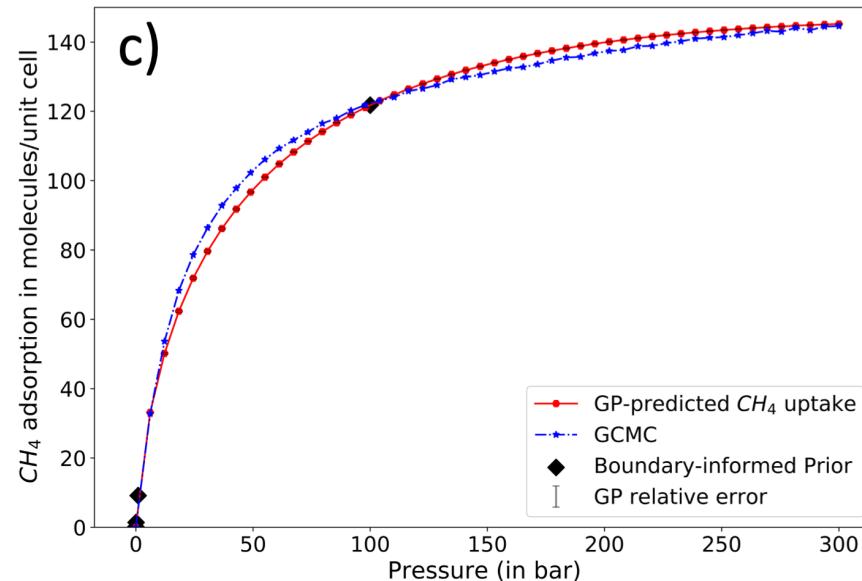
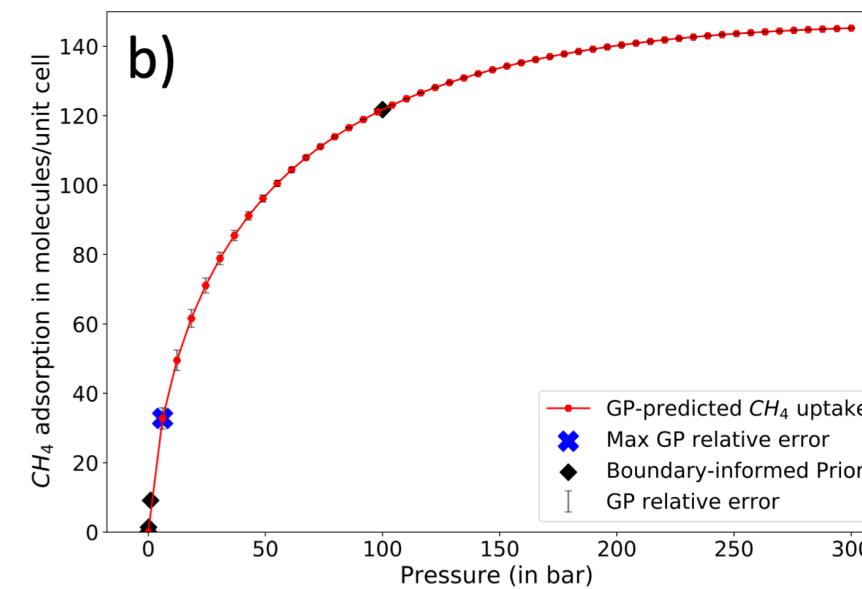
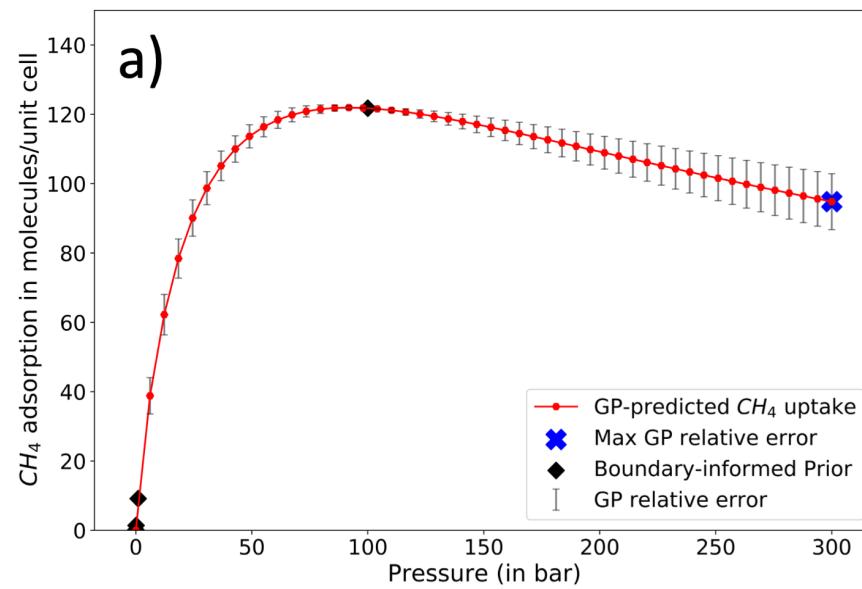


$$f \sim \mathcal{N} \left(\begin{bmatrix} \mu(x_1) \\ \mu(x_2) \\ \vdots \\ \mu(x_n) \end{bmatrix}, \begin{bmatrix} K(x_1, x_1) & \cdots & K(x_1, x_n) \\ K(x_2, x_1) & \ddots & \vdots \\ \vdots & \cdots & \vdots \\ K(x_n, x_1) & \cdots & K(x_n, x_n) \end{bmatrix} \right)$$

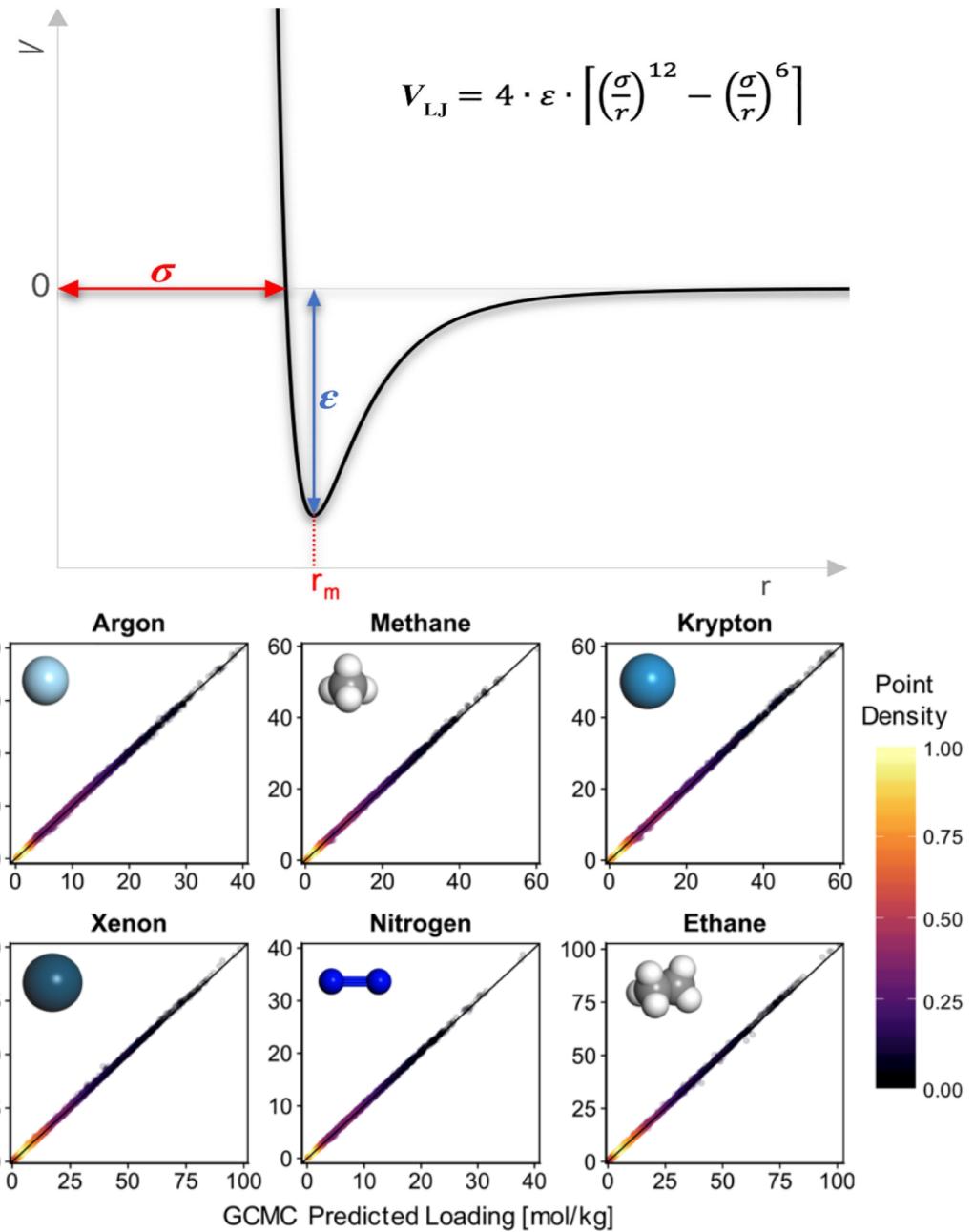
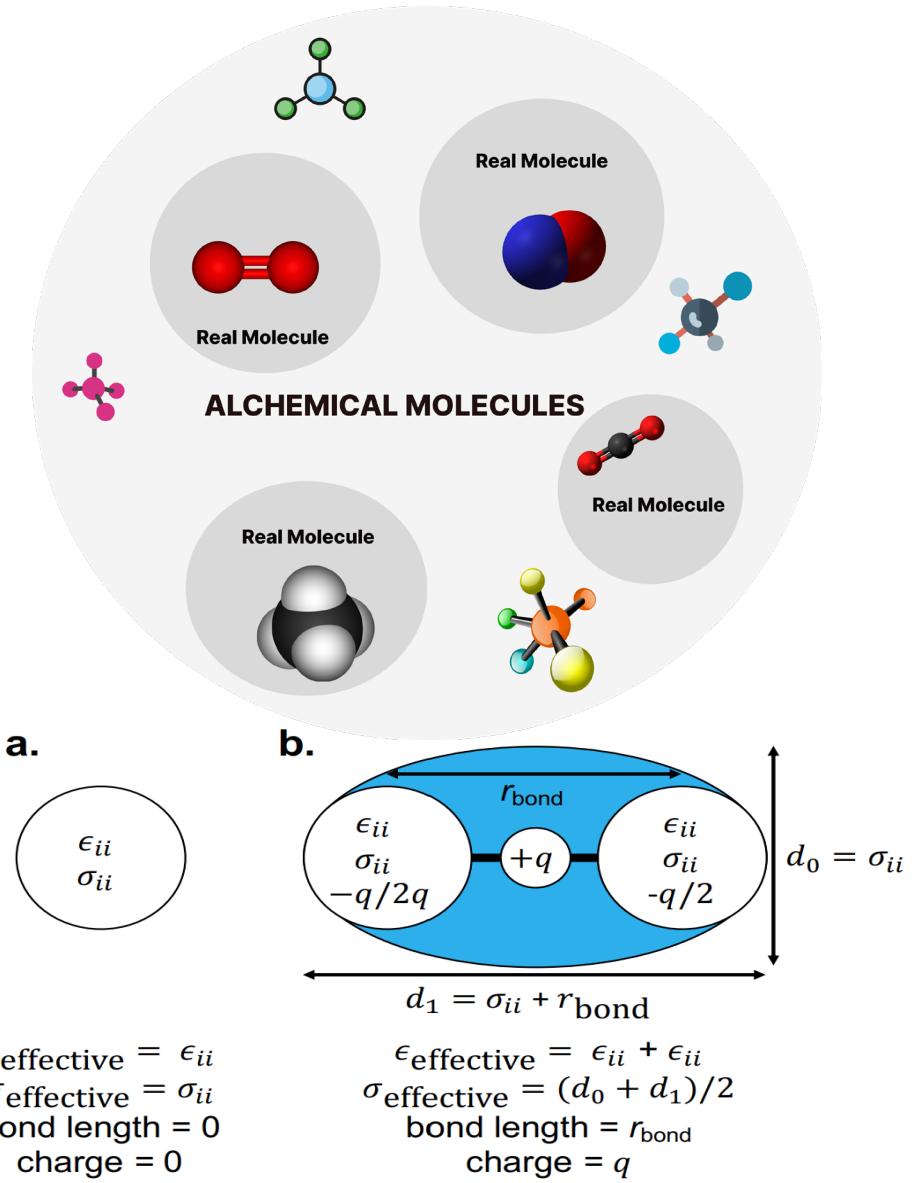
Sequential Design of Adsorption Simulations



Methane Isotherm

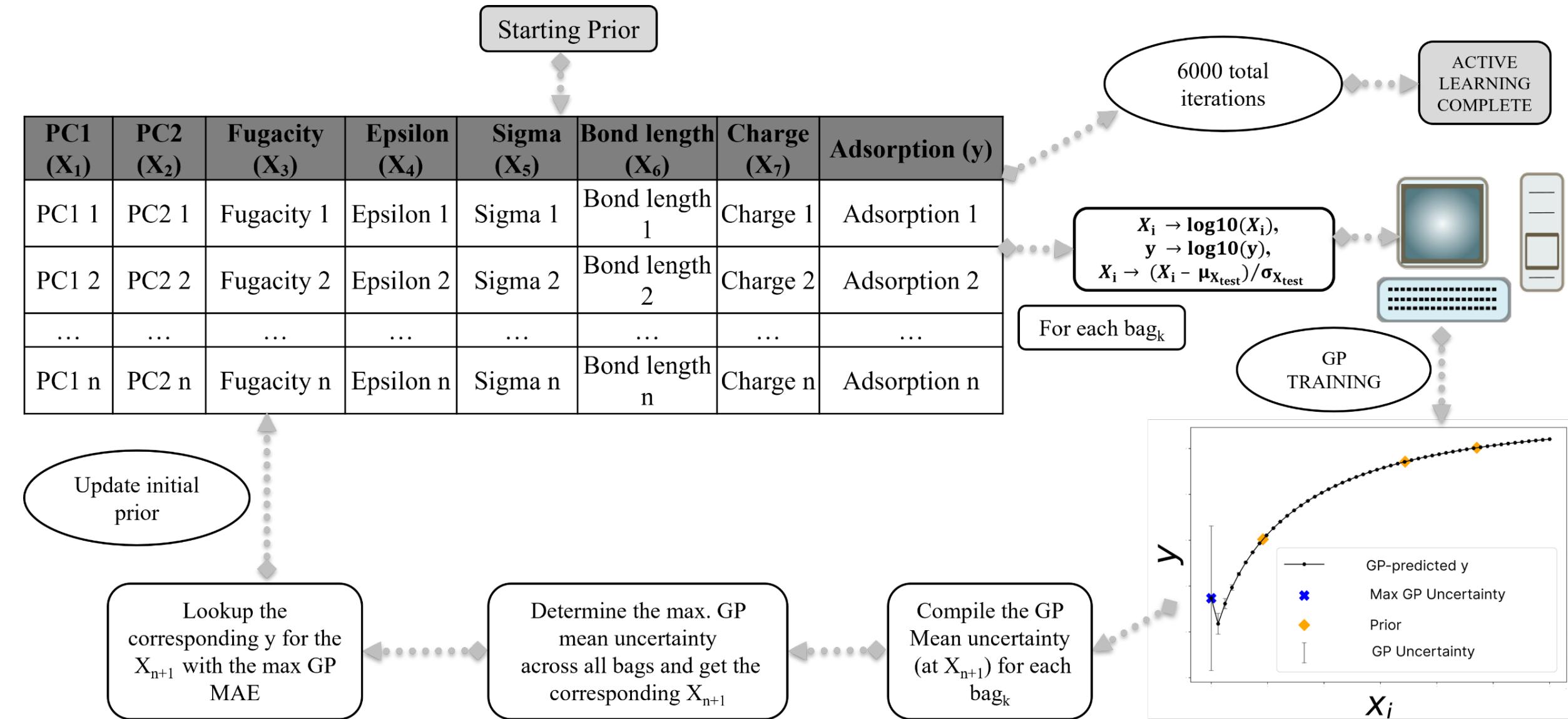


Universal Adsorption through Alchemical Species



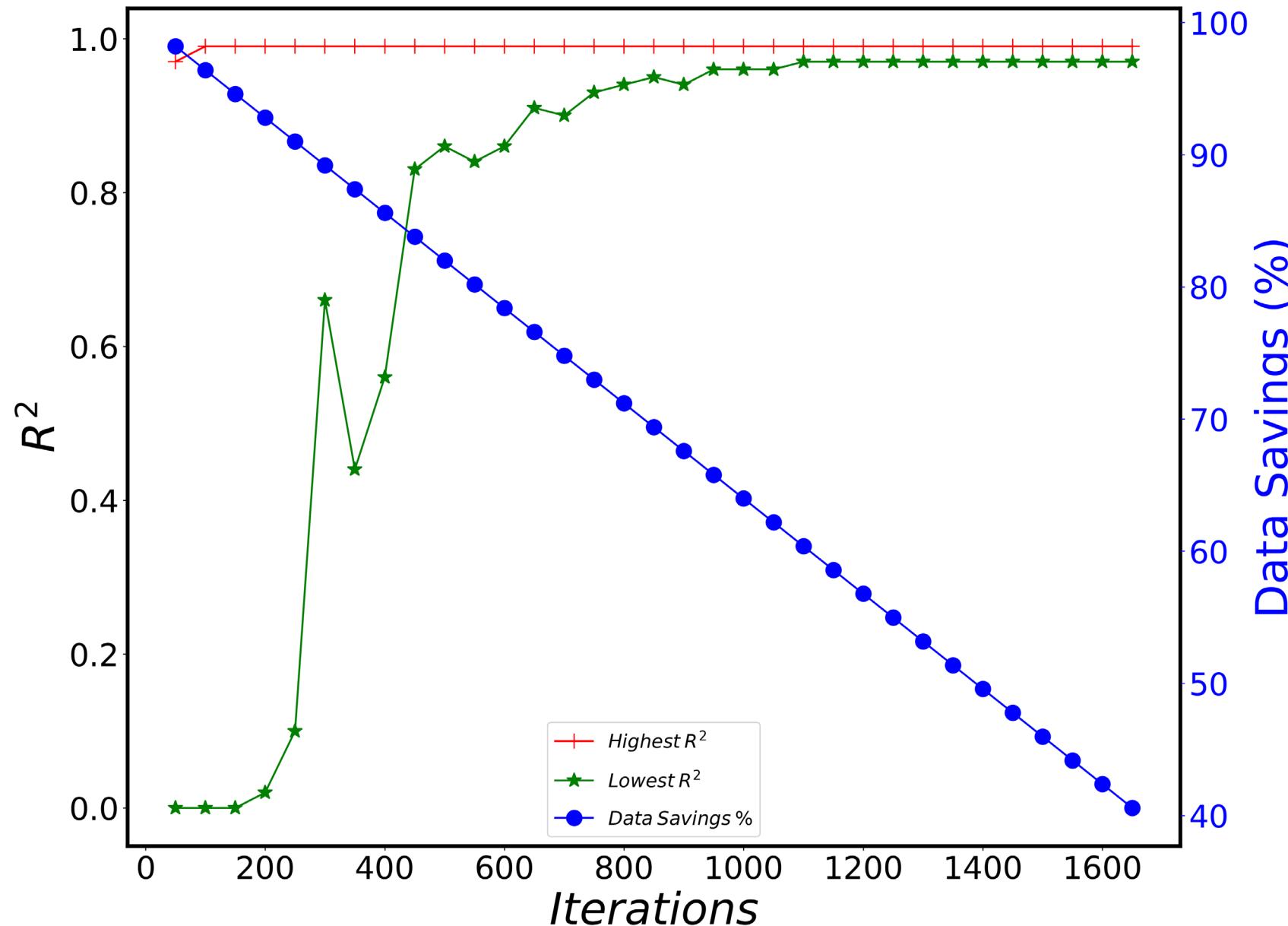
AL in alchemical space

$$N \sim f(F, \varepsilon, \sigma, l, q)$$



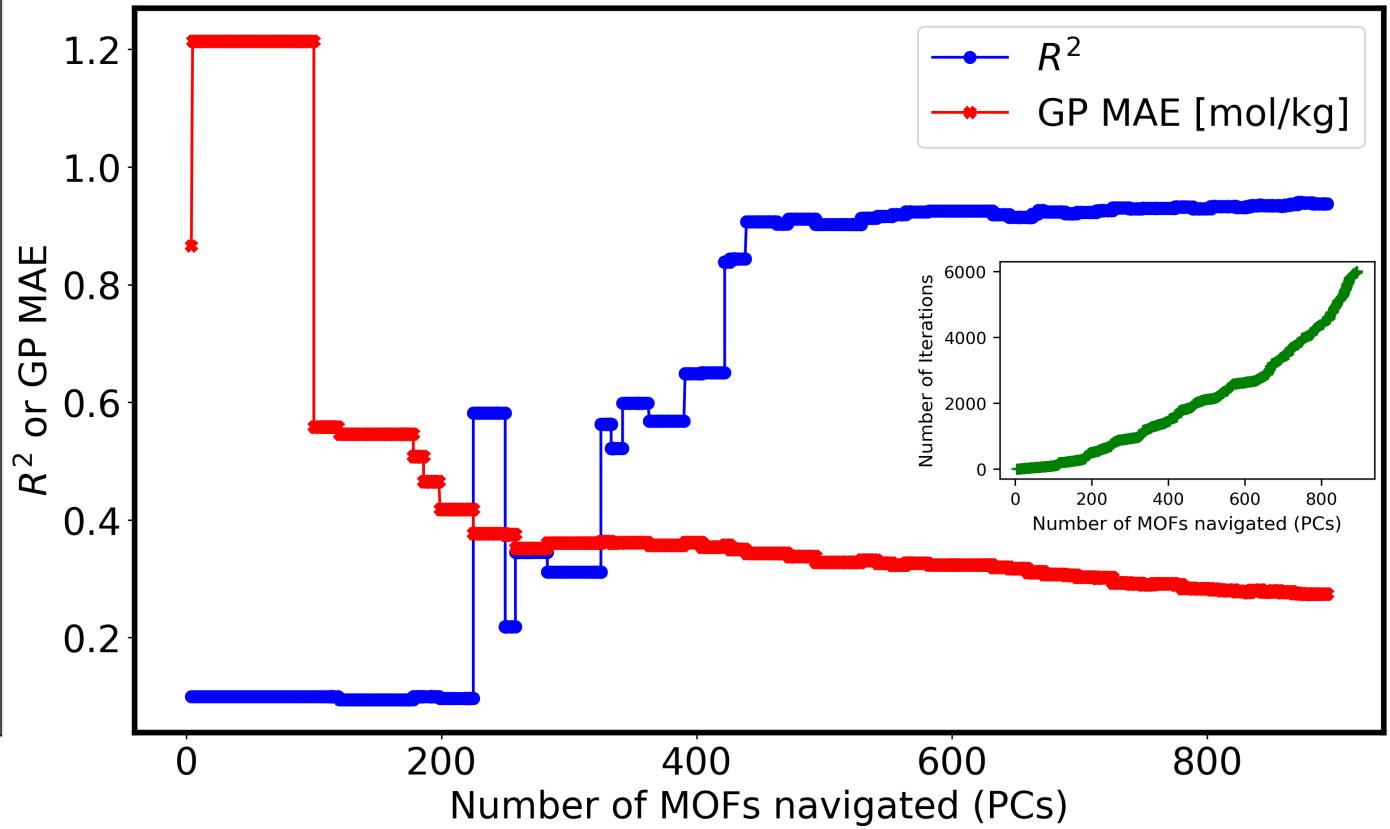
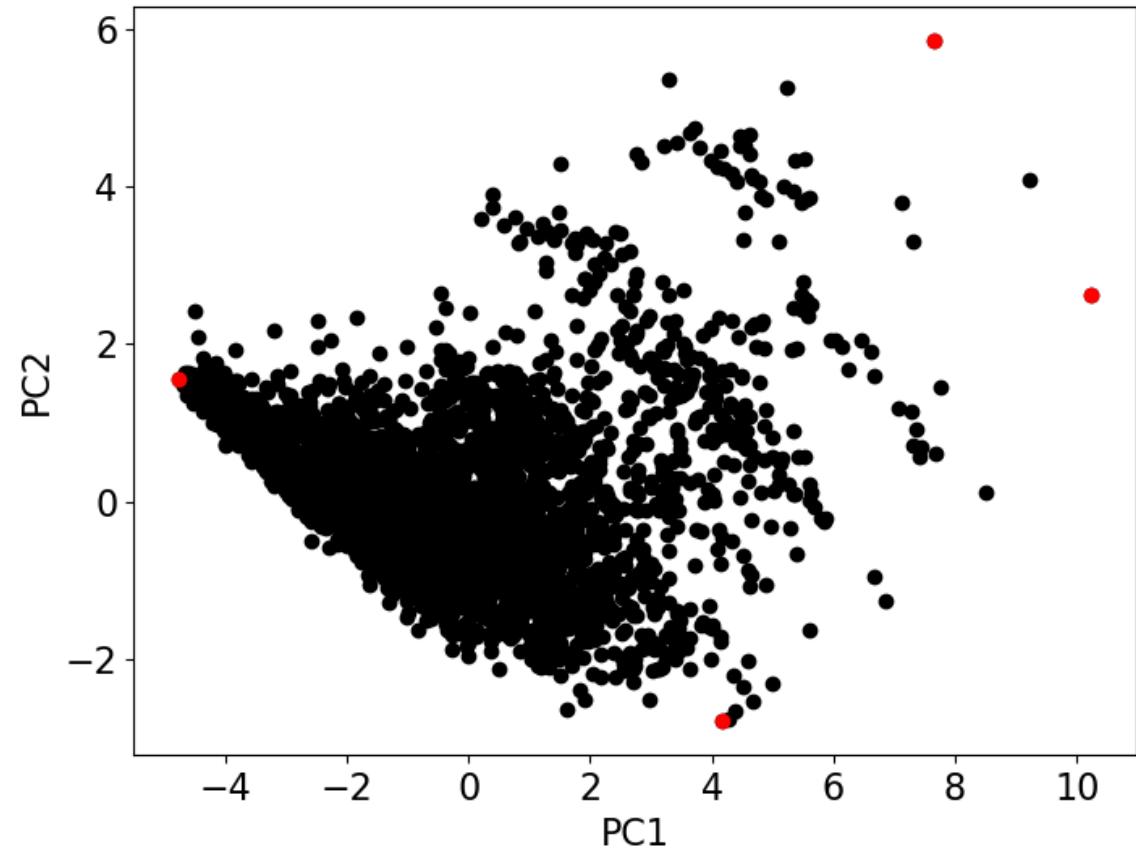
Significant Data Savings!

$$N \sim f(F, \varepsilon, \sigma, l, q)$$



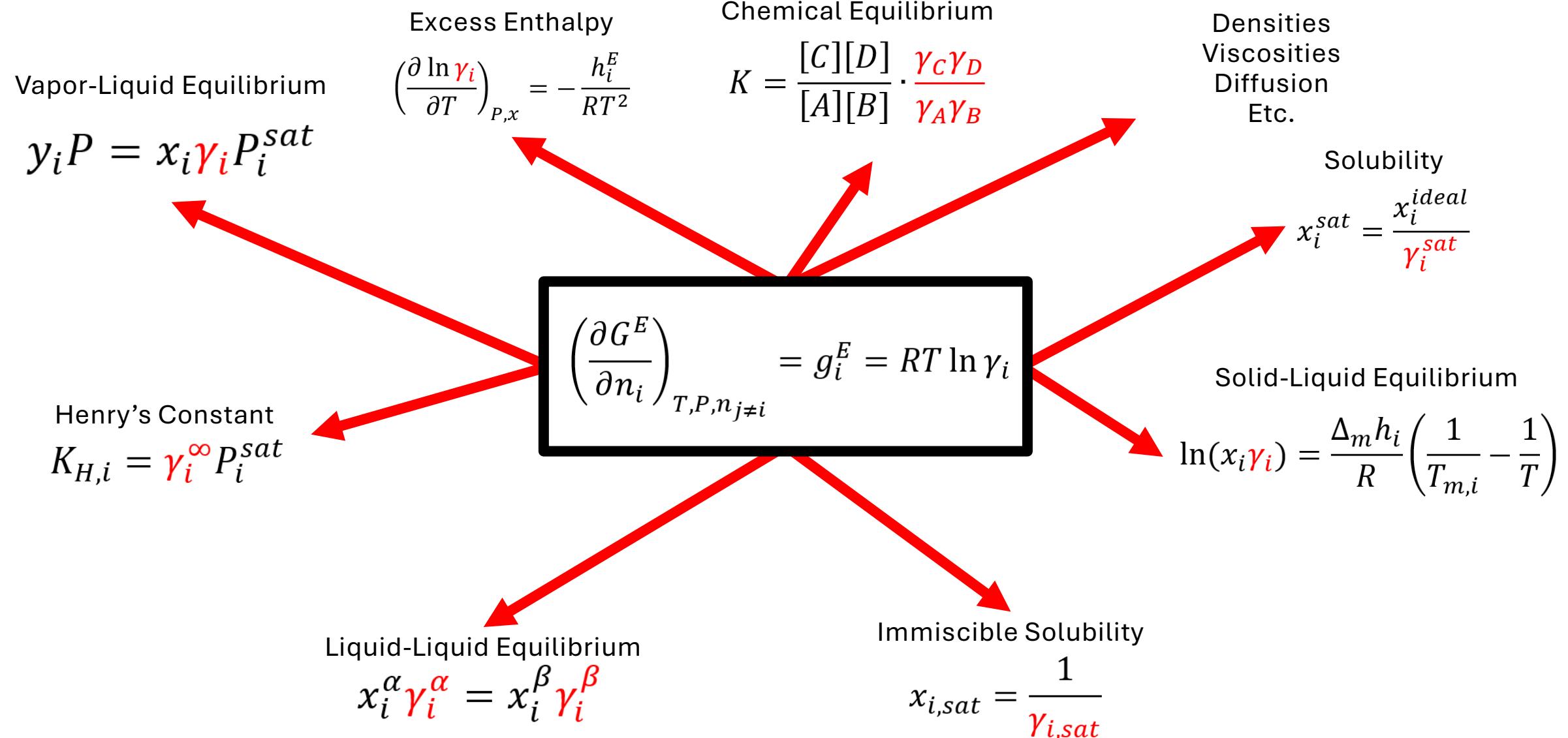
Navigating Alchemical and Material Spaces

$$N \sim f(PC1, PC2, F, \varepsilon, \sigma, l, q)$$

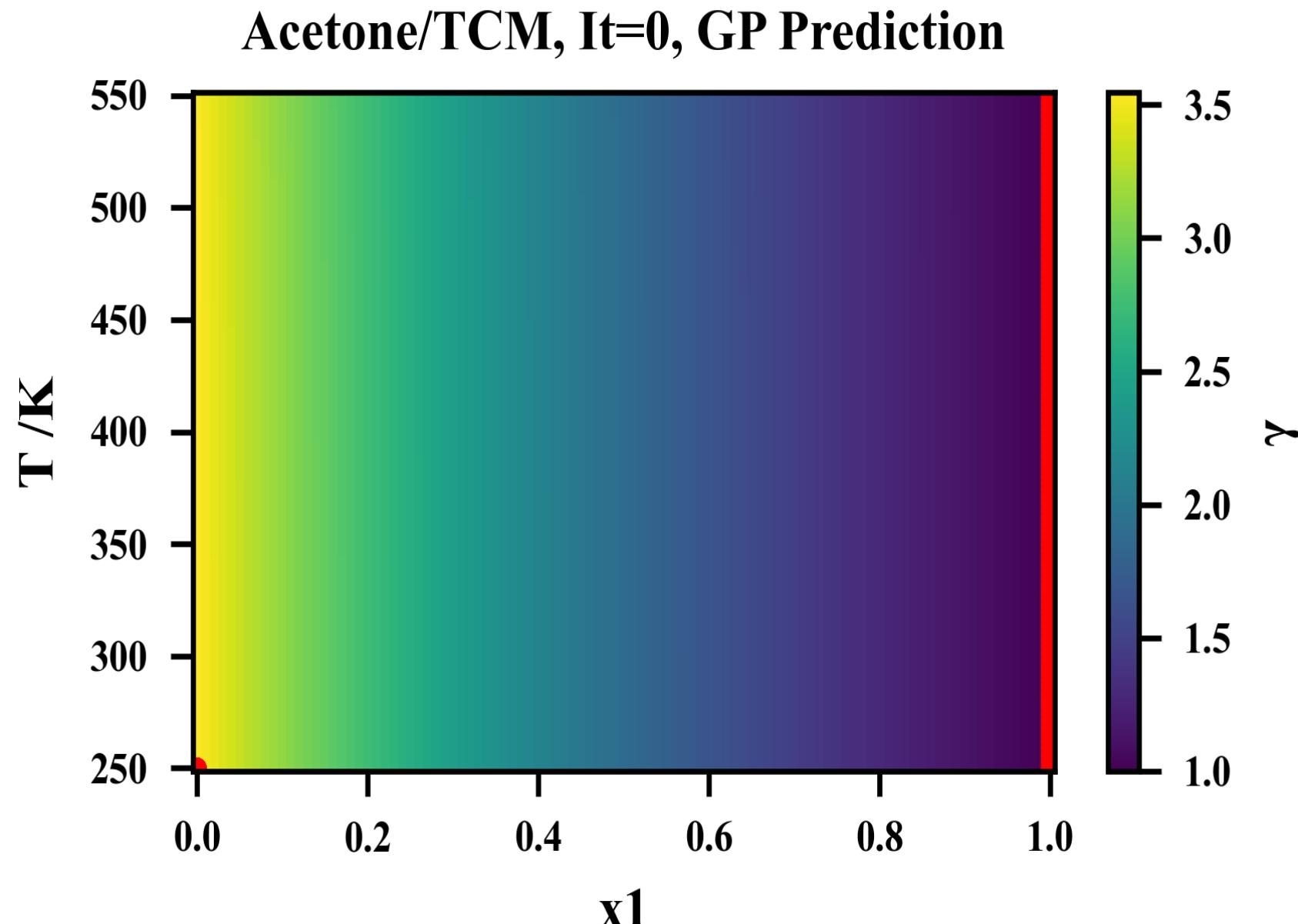


6,000 vs. 5,000,000 data

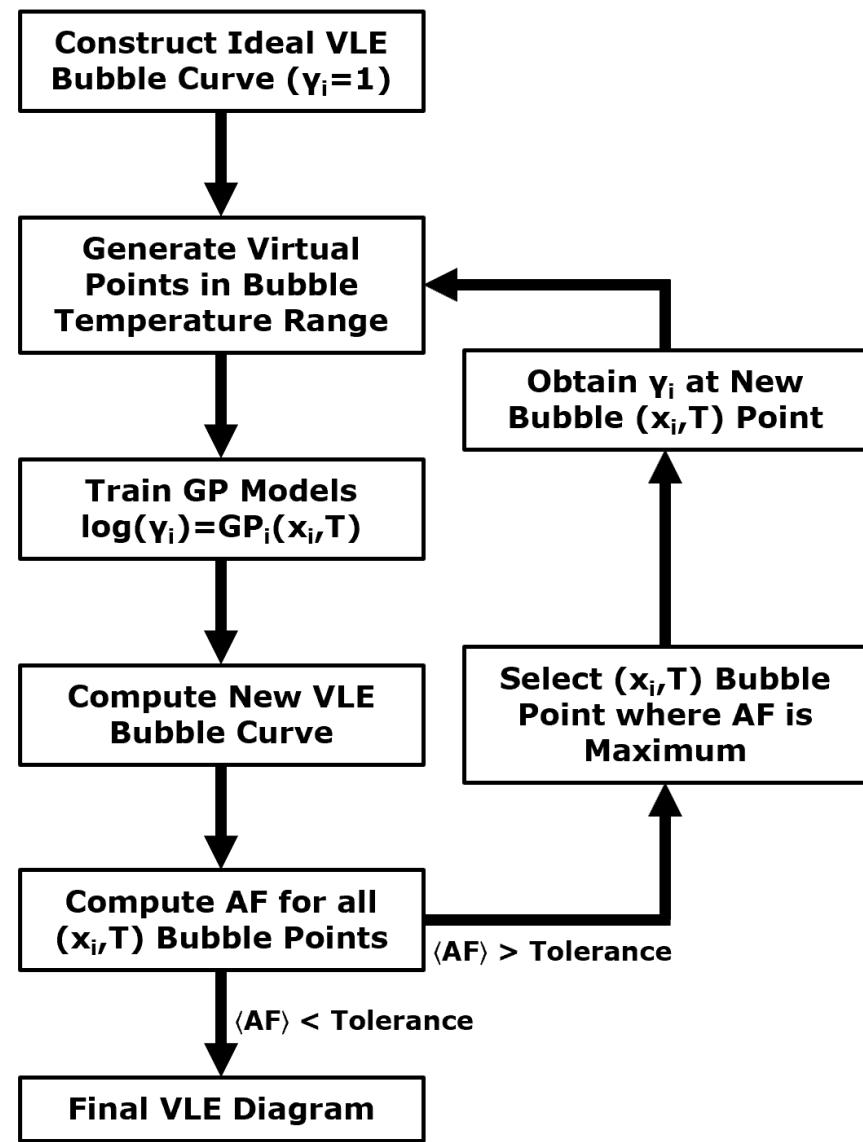
Active Learning on Activity Coefficients



NRTL as ground truth; GPs for activity coefficients



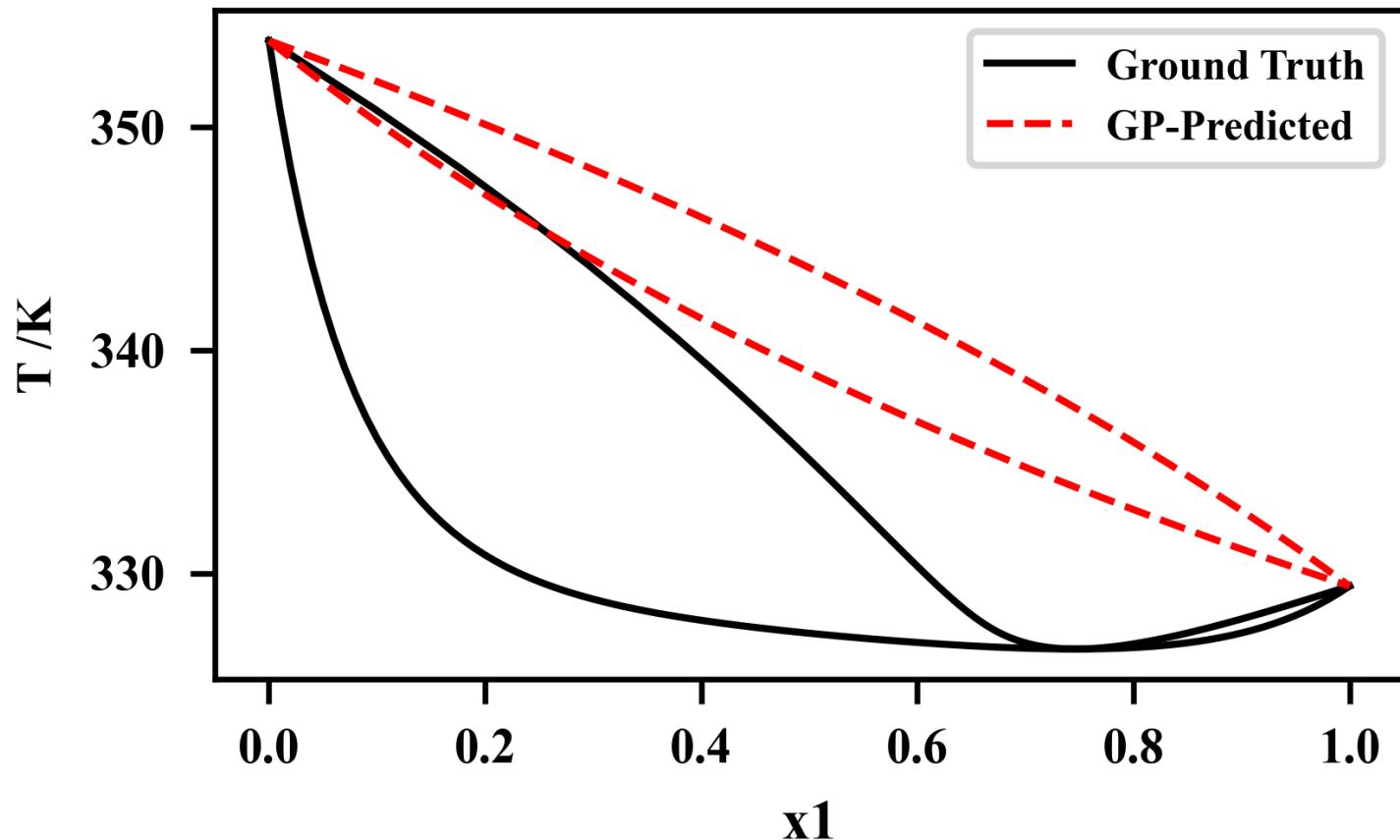
Thermodynamics-Informed AL



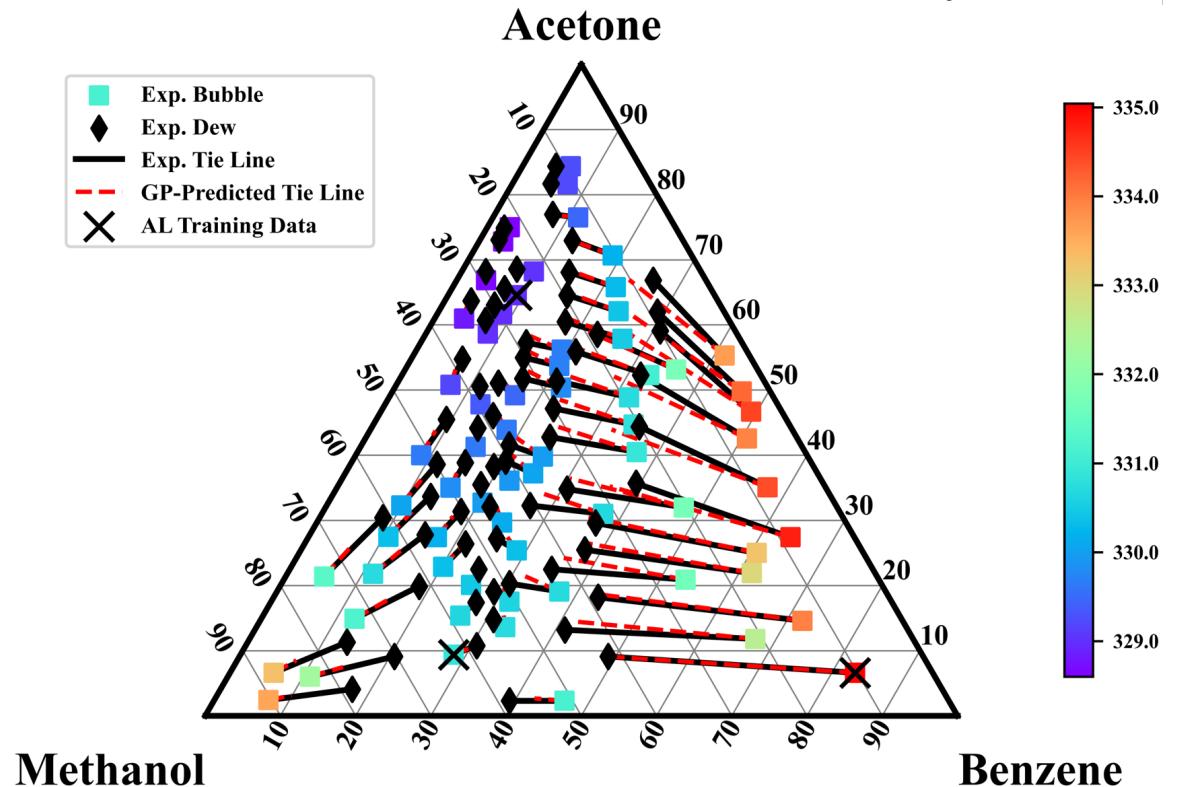
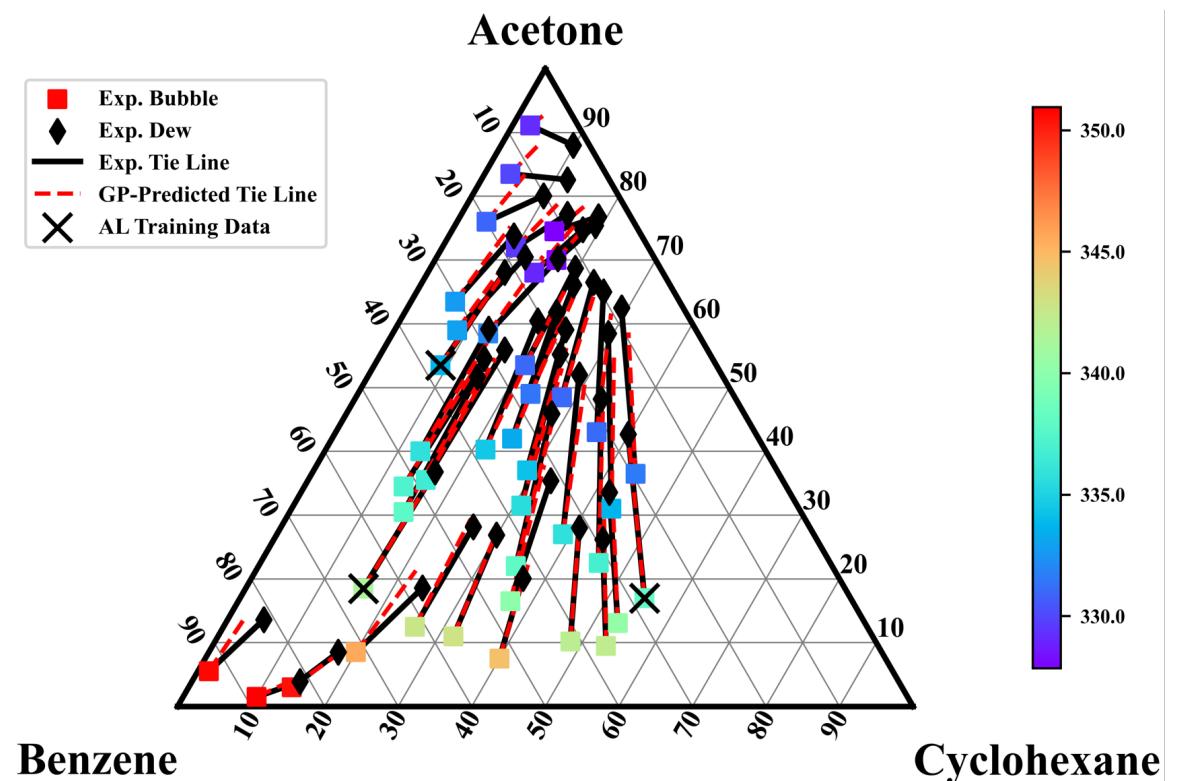
$$y_i P = x_i \gamma_i P_i^{sat}$$

$$S_{\gamma_i} = \left| \frac{\partial}{\partial \gamma_i} \left(\frac{x_i \gamma_i P_i^{sat}}{P} \right) \right| S_{\gamma_i} = \frac{x_i P_i^{sat}}{P} S_{\gamma_i}$$

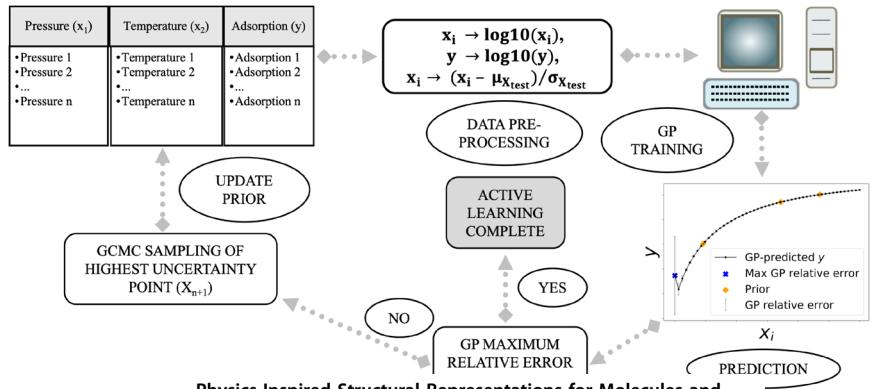
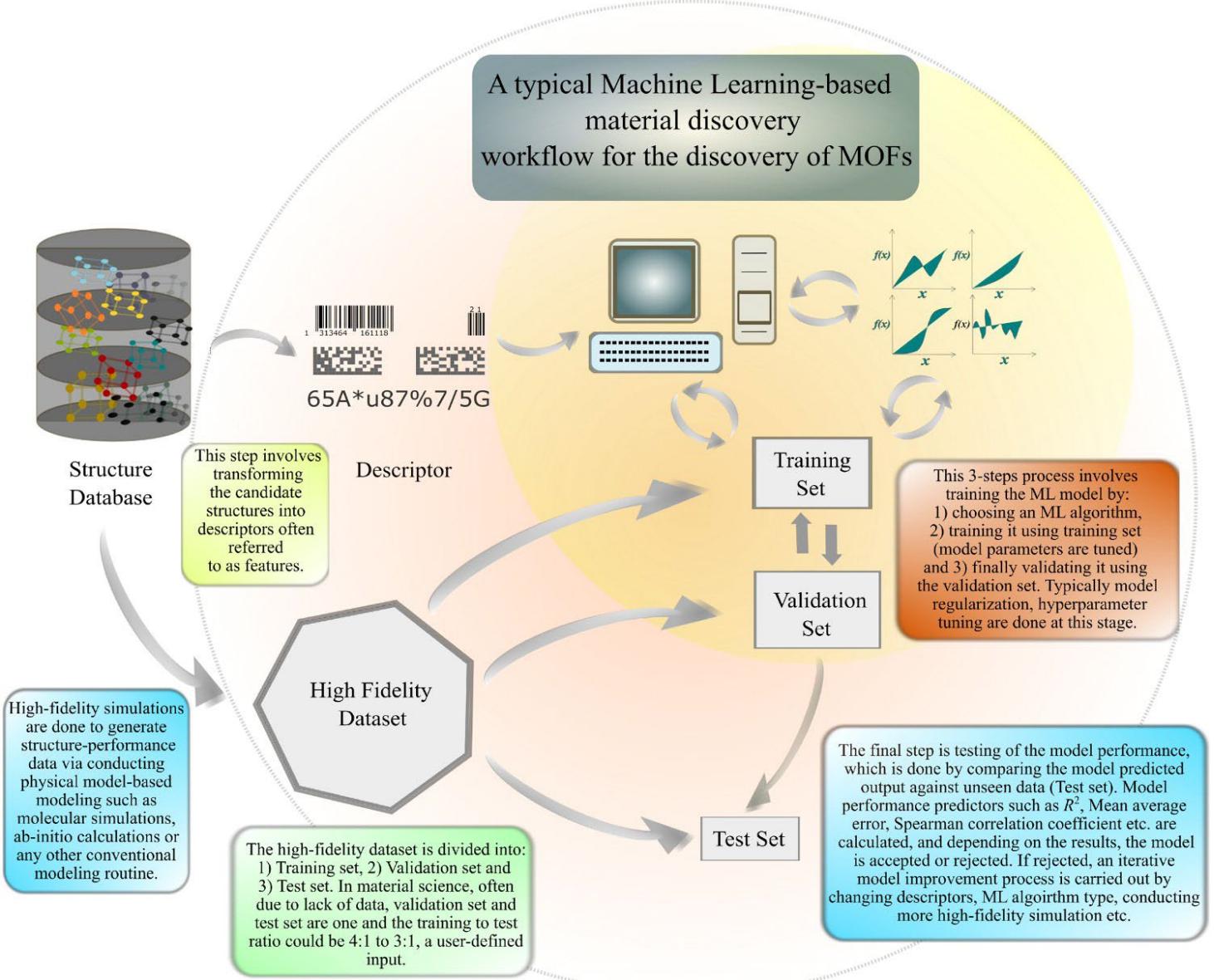
$$AF(x_i, T) = s_y^{\Sigma} = \sum_{i=1}^N s_{\gamma_i}(x_i, T)$$



Multicomponent VLE with AL



Summary



Felix Musil, Andrea Grisafi, Albert P. Bartók, Christoph Ortner, Gábor Csányi, and Michele Ceriotti*

