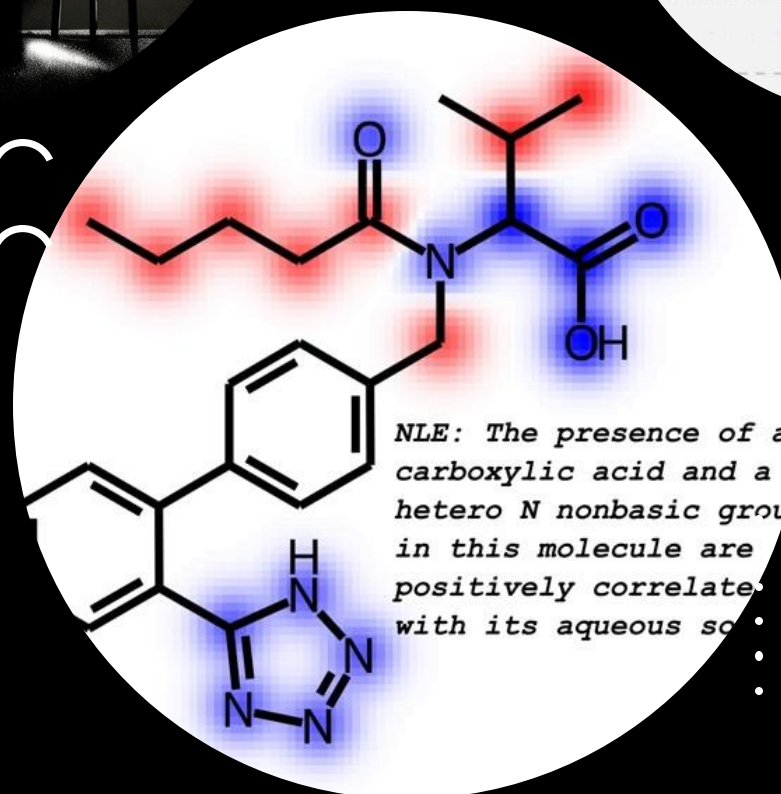
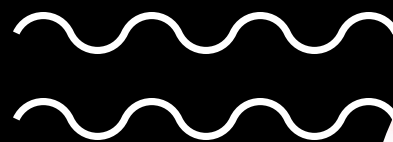
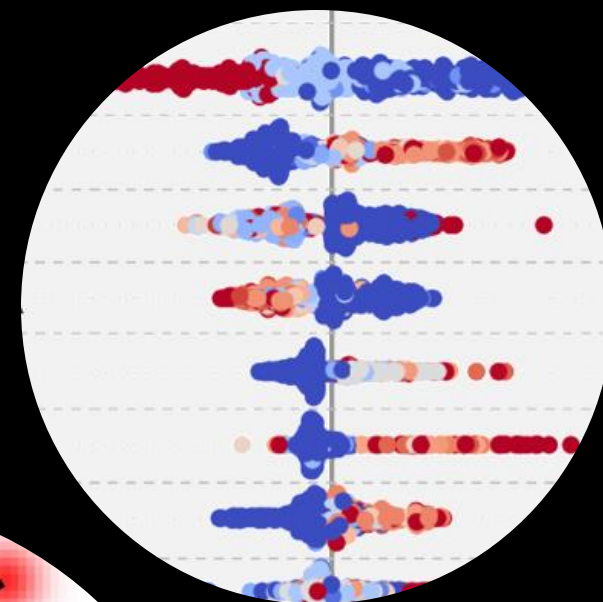


Explainable Artificial Intelligence

how to interrogate a black box?



The what and why's of explainable AI

Explainable AI (XAI) encompasses a broad set of topics that deal with the problem of trying to rationalize, understand, and/or distill meaning from ML models

This is important and under-leveraged:

- XAI makes ML models more practical and practicable
we are more likely to use models if we can understand why they work
- XAI is important from a regulatory or compliance perspective
dissemination and action depends on credibility and trust;
XAI is now becoming a part of guidelines and legislation
- XAI can help prevent misuse
if we understand why something works, we may anticipate when it won't

In physical science and engineering

- we are often tasked with *understanding*; it is part of our job
- there is bias against the use of black-box models and enhanced scrutiny
- we are probably leaving fruit on the tree; there's juice still to squeeze; etc.

Interpretation vs. Justification vs. Explanation

There are many useful routes in XAI; however, we should be clear on what it is that we are actually explaining!

- **Interpretation:** why has the model made the prediction it has made? Here, we are ascertaining an underlying cause for the decision. This explains the model but not the phenomena.
- **Justification:** why should the prediction be trusted? Here, we are considering whether there are valid reasons to support the model's decision.
- **Explanation:** how is the prediction actually linked to the input? Why are things the way that they are. This is more distilling cause and effect at the meta level. Such explanations usually afford an understanding that is actionable.

Explanations via post hoc analysis

ML models and their parameters can themselves be difficult to interpret; often, “explanations” are formulated via *post hoc* analysis

- **Training data importance:** what data is most important for formulating a given prediction? This provides context for a prediction.
- **Feature Importance:** what factors (usually of the input) were most influential in making the prediction? Analysis usually reveals important trends of the data.
- **Counterfactual analysis:** how can predictions/outcomes minimally differ and what changes? That’s the reason! right?

Training Data Importance

ML models are the way they are in response to the training data.
Therefore, we can posit that the prediction being the way it is
should be traceable to some subset of the training data.

Basic idea: evaluate predictions in the presence/absence of training with certain data

$$\mathcal{I}(\mathbf{x}_i, \mathbf{x}) = \hat{f}_{-\mathbf{x}_i}(\mathbf{x}) - \hat{f}(\mathbf{x})$$

“influence” of \mathbf{x}_i

*This possesses conceptual similarities to ablation analysis
(or is perhaps a specific manifestation thereof);
ablation analysis is usually used to identify important factors in performance*

Feature Importance Assessment

Feature importance is currently the most straightforward and common way of interpreting the outcomes of ML; there are both global and local variants

Basic formulations:

- *How responsive is the output to the inputs (gradient-esque evaluation)?*

consider a linear function: $\hat{y} = \boldsymbol{\theta}^T \mathbf{x} + b$

non-linear conceptual extrapolation: $\hat{f}(\mathbf{x}) \approx \hat{f}(\mathbf{x}') + \nabla \hat{f}(\mathbf{x}')(\mathbf{x} - \mathbf{x}')$

can use automatic differentiation methods or packages to implement; direct gradients may not be useful!

- *How is a prediction formulated in terms of relative contributions?*

Shapley Additive Explanations is a very common variant

$$\hat{f}(\mathbf{x}) = \sum_i \phi_i(\mathbf{x}) \quad \phi_i(\mathbf{x}) = \frac{1}{Z} \sum_{S \in N - \{x_i\}} v(S \cup x_i) - v(S)$$

SHAP analysis is model-agnostic, which is nice!

$$\phi_i(\mathbf{x}) \approx \frac{1}{M} \sum_{m=1}^M (\hat{f}(\mathbf{z}_{+i}(\sim \mathbf{x}_m)) - \hat{f}(\mathbf{z}_{-i}(\sim \mathbf{x}_m)))$$

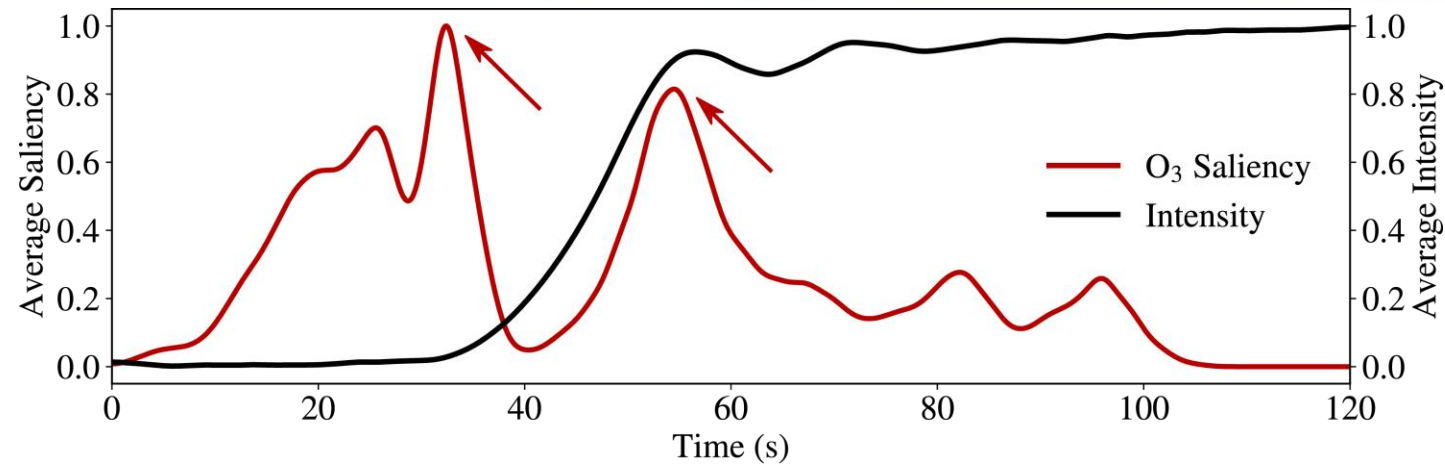
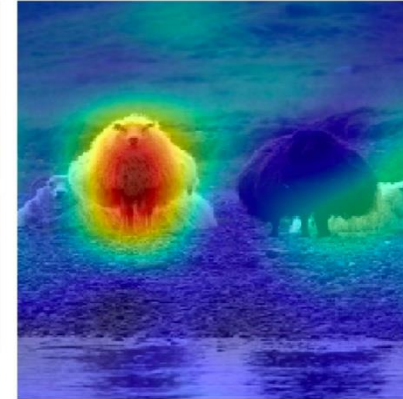
Saliency/Gradient Analysis: What's important?

$$S = \frac{\partial F(x)}{\partial x} = \frac{\partial \text{Model output}}{\partial \text{Model input}}$$

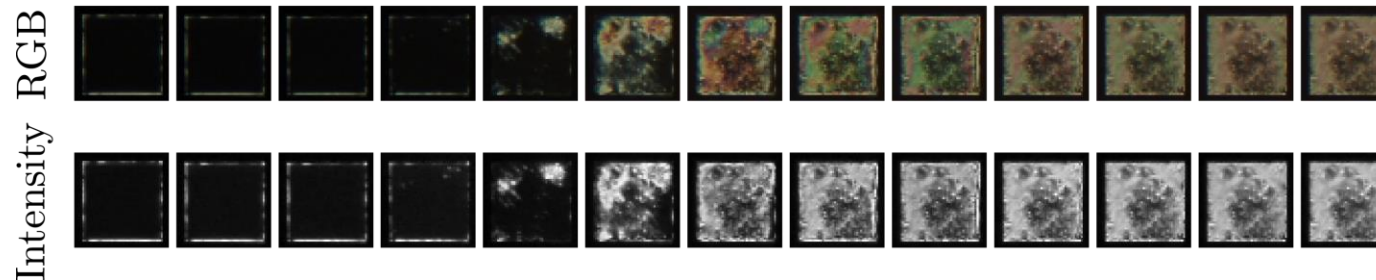
Model input

Prediction as "sheep"

Prediction as "cow"



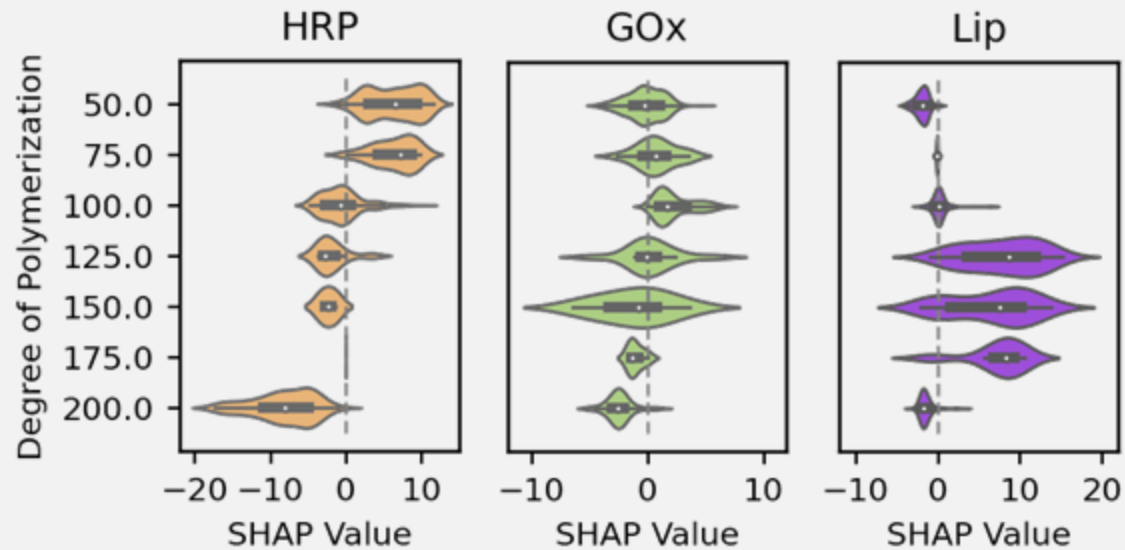
optical response
of LC array



SHAP Analysis Plots

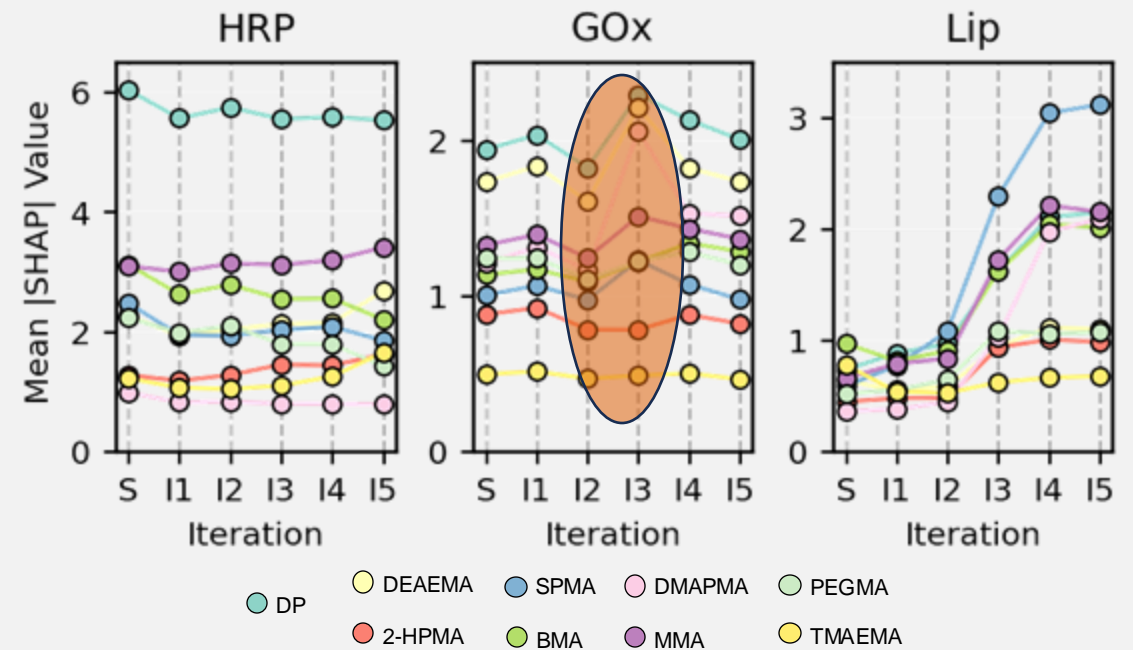
Feature Importance by SHapley Additive exPlanations

Relative Importance of Polymer Length



Each enzyme displays unique and intriguing preferences on the degree of polymerization

Evolution of Importance during Active Learning



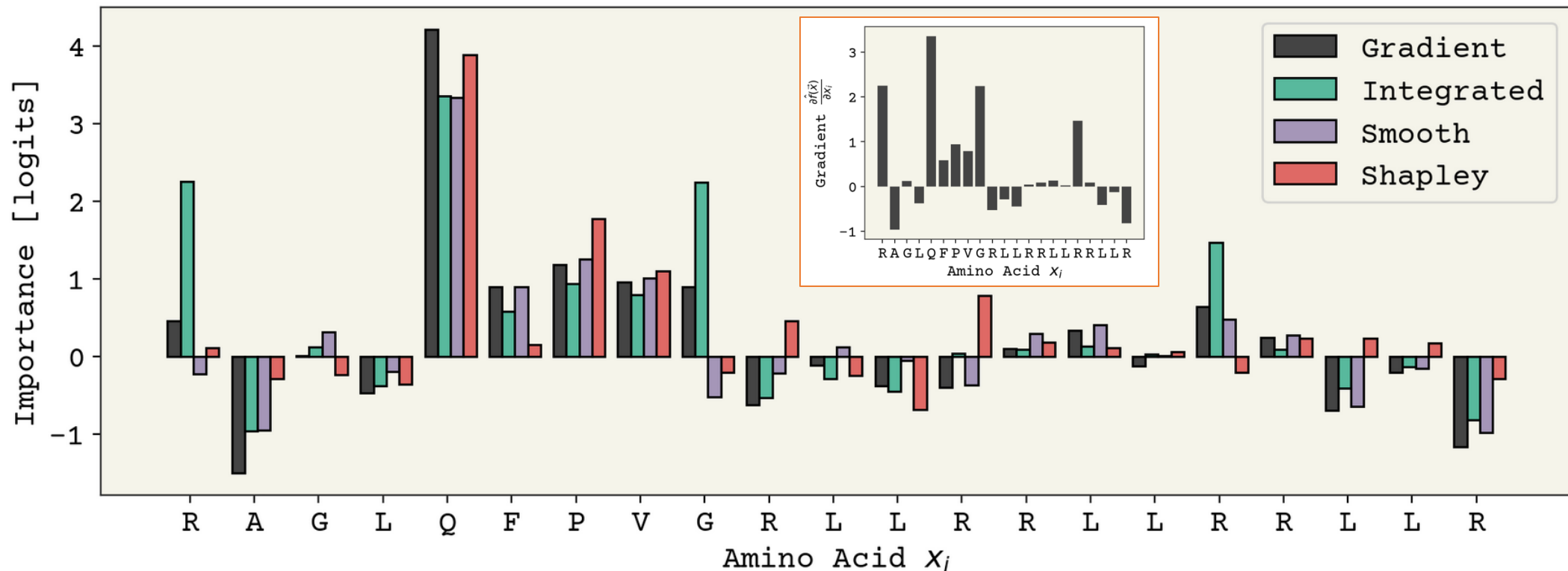
Discovery of high-performing systems can be a steady march or be a challenging excursion

Comparison of Feature Importance Methods

Task: is a peptide homolytic (i.e., will it kill red blood cells) ?



deep learning for
molecules & materials



Counterfactual Analysis

Counterfactual analysis aims to provide an explanation by comparison. Given a particular example/behavior, the idea is to find something that is close in all respects except the behavior to be explained.

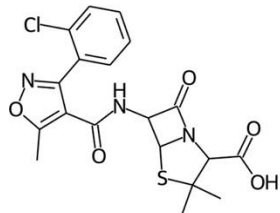
$$\min d(x, x')$$

such that

$$\hat{f}(x) \neq \hat{f}(x') \text{ or } \hat{f}(x) - \hat{f}(x') > \delta$$

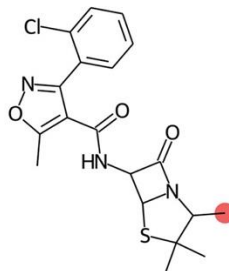
*permeate blood-
brain barrier?*

Base
 $f(x) = 0.000$



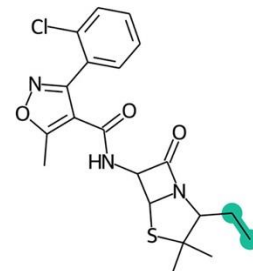
no

Similarity = 0.80
Counterfactual 1
 $f(x) = 1.000$



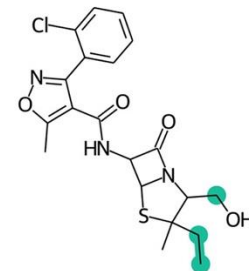
yes

Similarity = 0.77
Counterfactual 2
 $f(x) = 1.000$



yes

Similarity = 0.71
Counterfactual 3
 $f(x) = 1.000$



yes

Counterfactual Analysis

Counterfactual analysis aims to provide an explanation by comparison. Given a particular example/behavior, the idea is to find something that is close in all respects except the behavior to be explained.

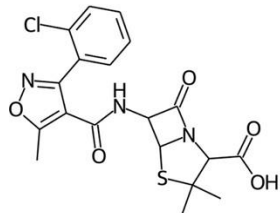
$$\min d(x, x')$$

such that

$$\hat{f}(x) \neq \hat{f}(x') \text{ or } \hat{f}(x) - \hat{f}(x') > \delta$$

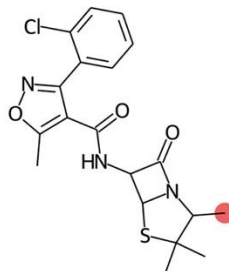
*permeate blood-
brain barrier?*

Base
 $f(x) = 0.000$



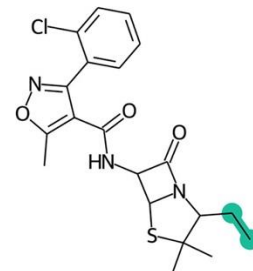
no

Similarity = 0.80
Counterfactual 1
 $f(x) = 1.000$



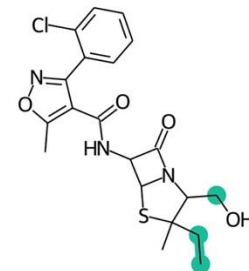
yes

Similarity = 0.77
Counterfactual 2
 $f(x) = 1.000$



yes

Similarity = 0.71
Counterfactual 3
 $f(x) = 1.000$



yes

Counterfactual Analysis

Counterfactual analysis aims to provide an explanation by comparison. Given a particular example/behavior, the idea is to find something that is close in all respects except the behavior to be explained.

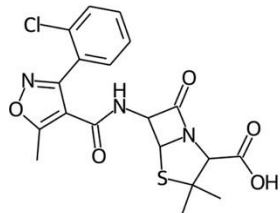
$$\min d(x, x')$$

such that

$$\hat{f}(x) \neq \hat{f}(x') \text{ or } \hat{f}(x) - \hat{f}(x') > \delta$$

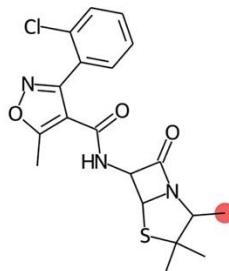
*permeate blood-
brain barrier?*

Base
 $f(x) = 0.000$



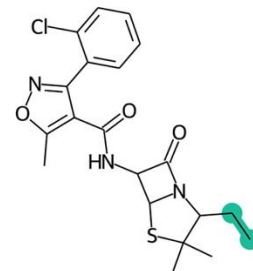
no

Similarity = 0.80
Counterfactual 1
 $f(x) = 1.000$



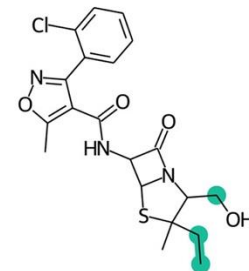
yes

Similarity = 0.77
Counterfactual 2
 $f(x) = 1.000$



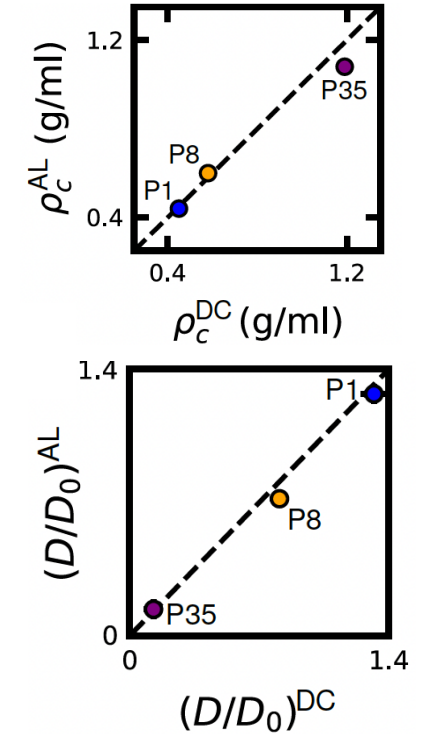
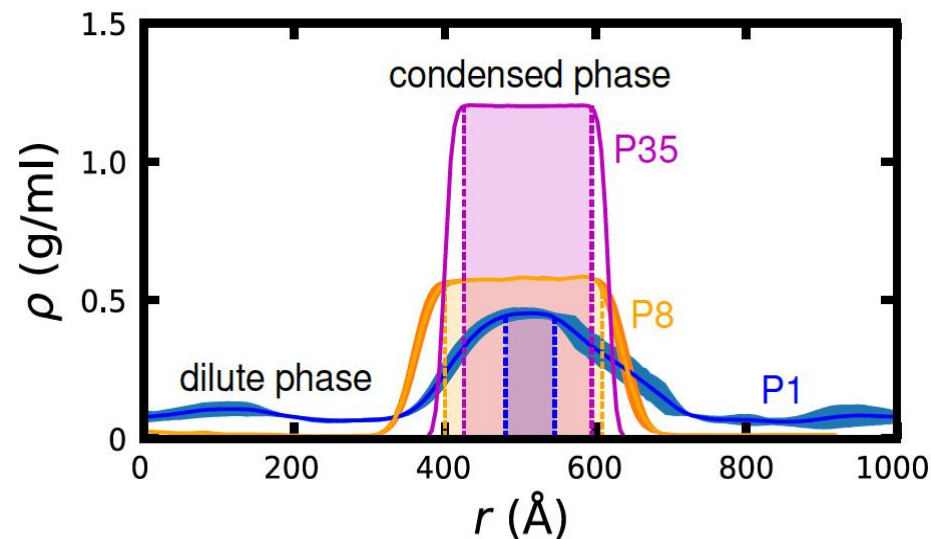
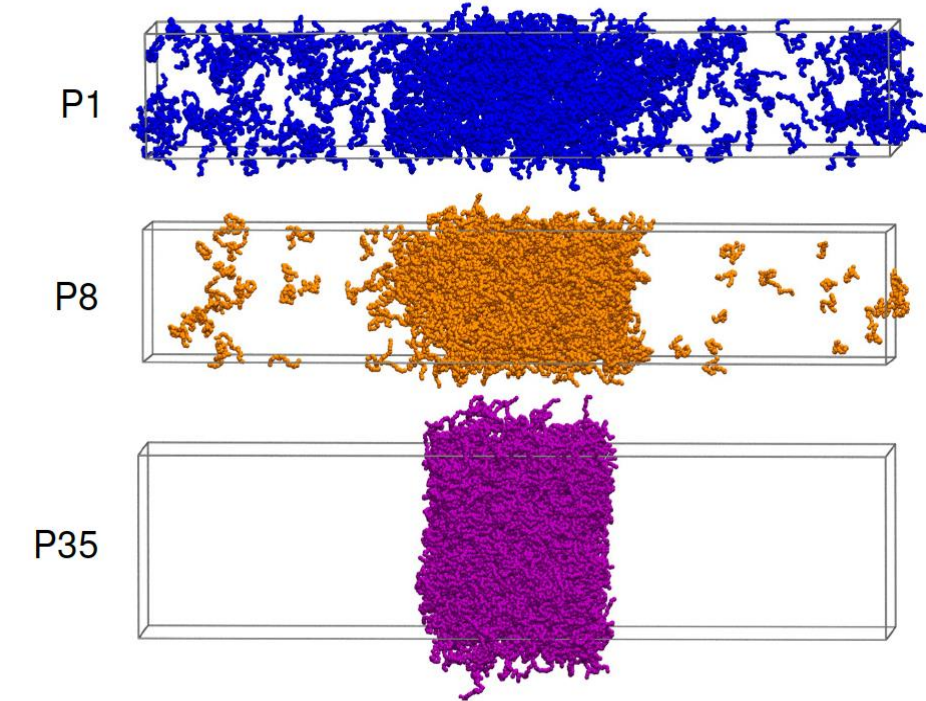
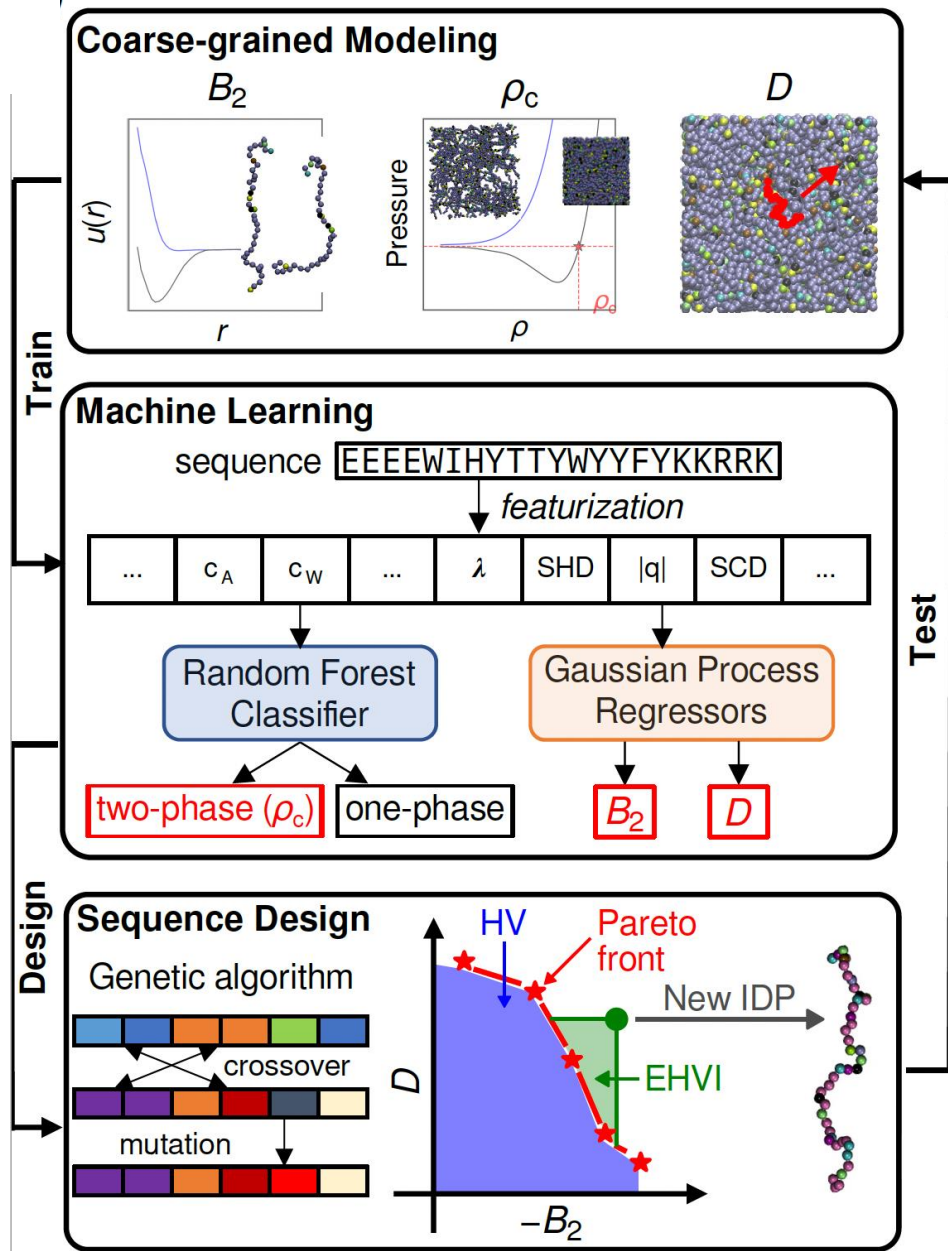
yes

Similarity = 0.71
Counterfactual 3
 $f(x) = 1.000$



yes

Counterfactual Analysis



brute-force

ML-assisted

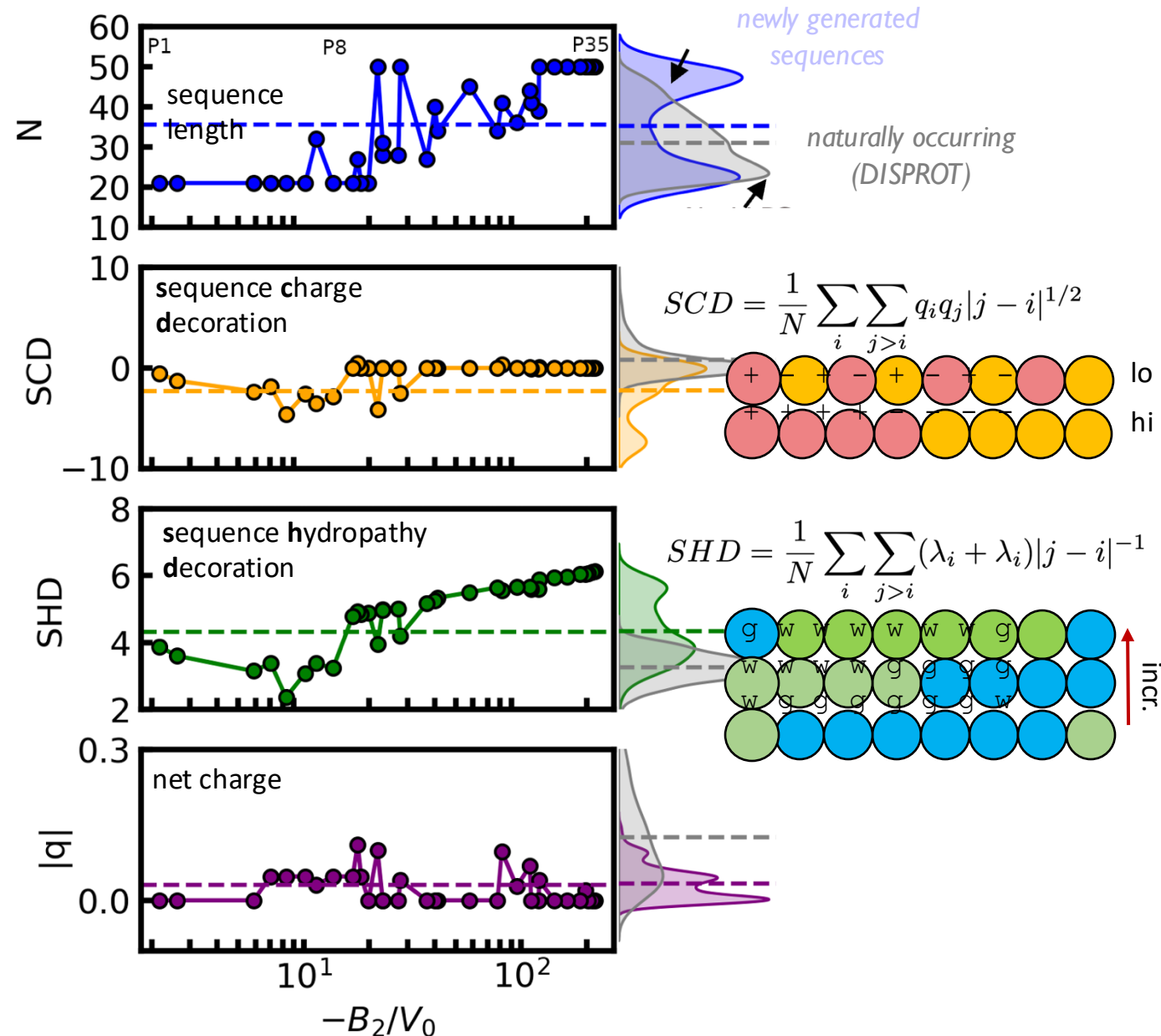
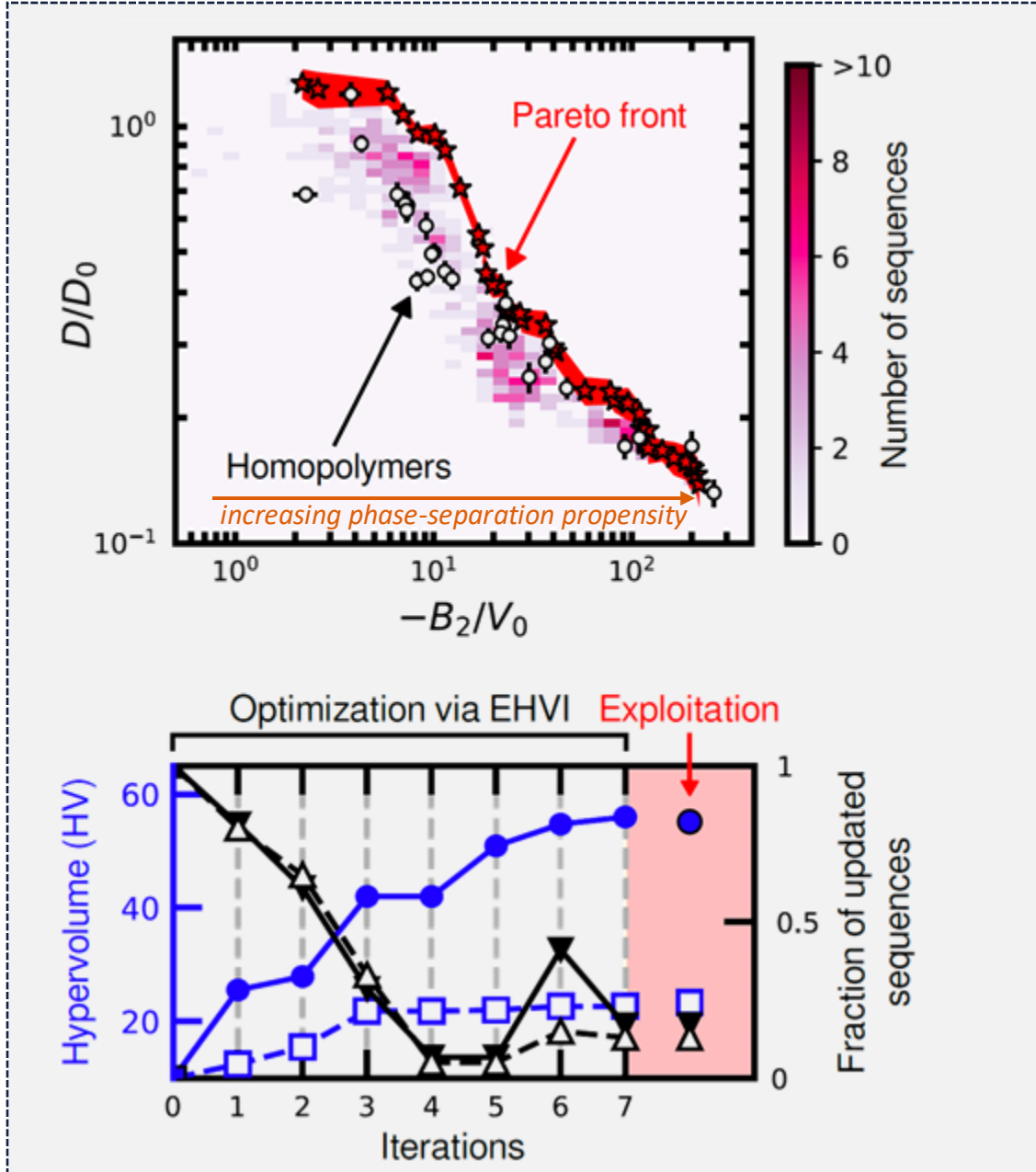
3 replicas
1000 chains
~10 μ s simulation

vs.

30 replicas
100 chains
100 ns simulation

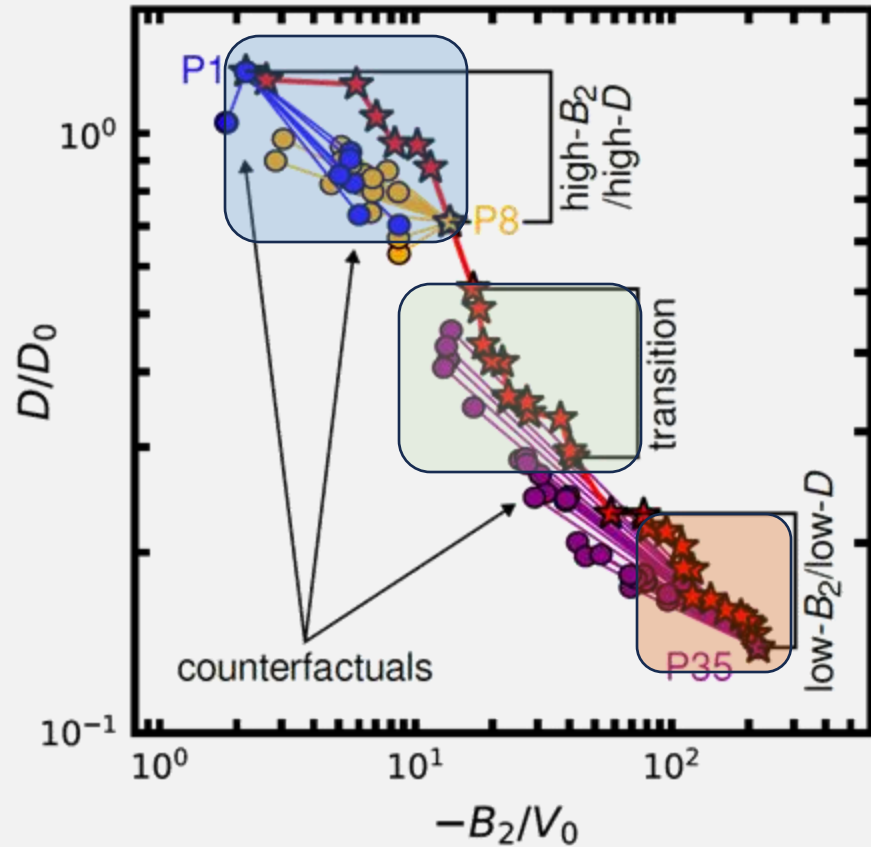
150x faster

Counterfactual Analysis



Counterfactual Analysis

The notion of a **counterfactual** can provide some insight as to why any particular sequence is “on the Front” compared to a “close” sequence that is not



- **high-high**: polyampholytes with subtle sequence variations
- **transition**: highly degenerate region with many options
- **low-low**: neutral sequences with strong hydrophobic patterning

$$\langle \Delta x_k \rangle \equiv \frac{1}{n_{C_{ij}}} \sum_{j=1}^{n_{C_{ij}}} (x_k^{(P_i)} - x_k^{(C_{ij})})$$

Average differences (Pareto - counterfactual)

