# Improving machine learning with polymer physics
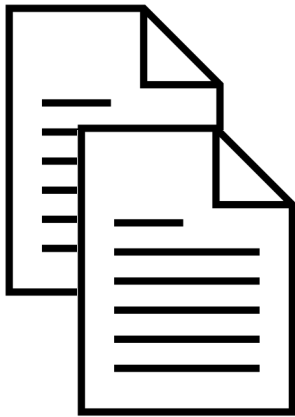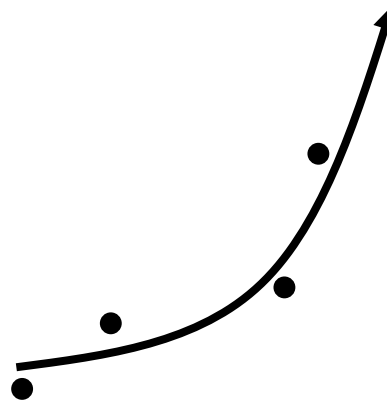
Debra J. Audus
2024 Artificial Intelligence for Materials Science (AIMS) Workshop
July 18, 2024
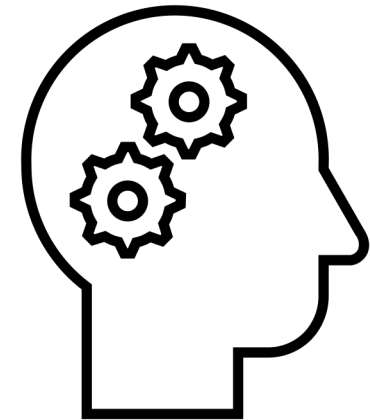
**NIST**
National Institute of
Standards and Technology
U.S. Department of Commerce

# Outstanding challenges

DATA
(fuel for ML/AI)

EXTRAPOLATION
(go beyond the dataset)

EXPLAINABILITY
(answer scientific questions)

# A path forward: use knowledge

## Theory-Guided Data Science: A New Paradigm for Scientific Discovery from Data

Anuj Karpatne, Gowtham Atluri, James H. Faghmous, Michael Steinbach, Arindam Banerjee, Auroop Ganguly, Shashi Shekhar, Nagiza Samatova, and Vipin Kumar

## Theory-Guided Machine Learning in Materials Science

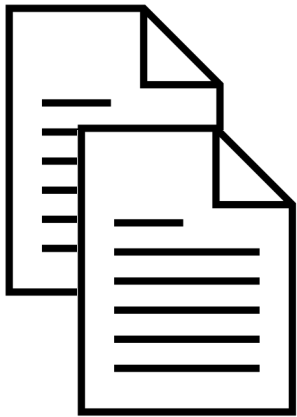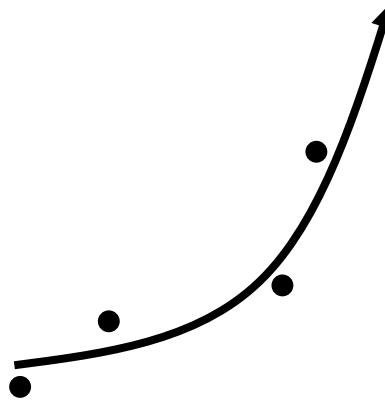Nicholas Wagner and James M. Rondinelli

## Embedding domain knowledge for machine learning of complex material systems
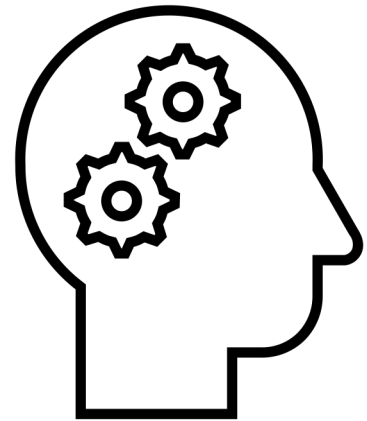
Christopher M. Childs and Newell R. Washburn

# Impact of theory



*Need less*
DATA
(fuel for ML/AI)

*Improves*
EXTRAPOLATION
(go beyond the dataset)

*May provide*
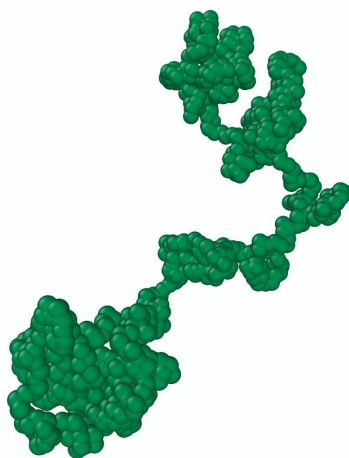EXPLAINABILITY
(answer scientific questions)

NIST

Bad
Solvent

Theta
solvent

Good
solvent

ML Input:
- $\ln N$
- $\alpha$ (solvent quality)

ML Output:
- $\ln R_g^2$

$R_g \sim N^{1/3}$

$R_g \sim N^{1/2}$

$R_g \sim N^{0.588}$

$R_g \sim N^{\nu}$

Austin McDannald, Brian DeCost (NIST)

# Benchmarks

**Direct**

$$x \equiv [\ln N, \alpha] \rightarrow \boxed{\text{GPR}} \rightarrow \ln R_g^2$$

**Theory**

$$t(x) = \ln N + \ln 1/6 = \ln R_g^2$$
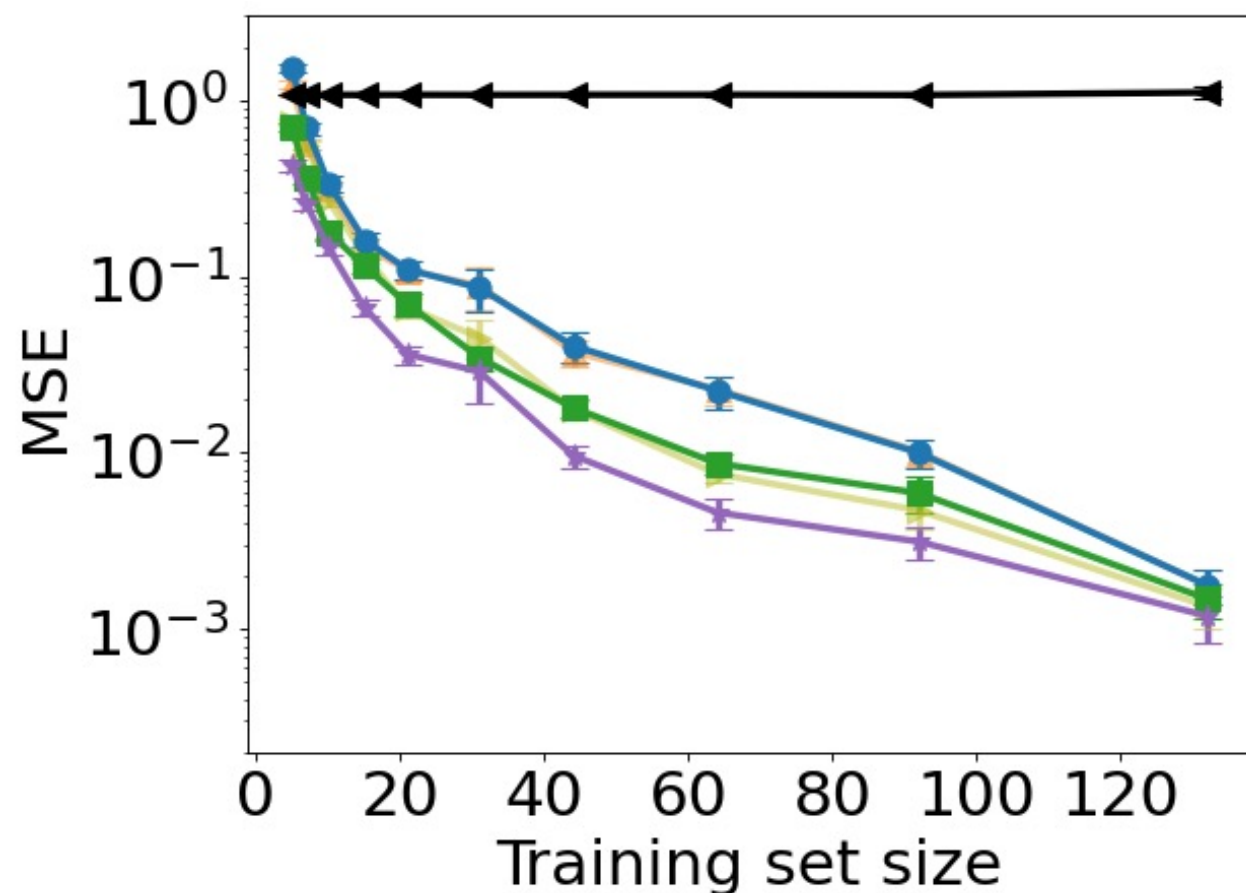
Purposely use an imperfect theory ($R_g \sim N^{1/2}$) because we often have theories that are only sometimes valid

# Parameterization: best!

**Direct**     $x \equiv [\ln N, \alpha] \rightarrow$ GPR $\rightarrow \ln R_g^2$

**Theory**     $t(x) = \ln N + \ln 1/6 = \ln R_g^2$

**Latent variable**     $x \rightarrow$   $t(x) \rightarrow$ GPR $\rightarrow \ln R_g^2$

**Difference**     $x \rightarrow$ GPR $\rightarrow \ln R_g^2 - t(x) \rightarrow$ sum $\rightarrow \ln R_g^2$   $t(x) \rightarrow$
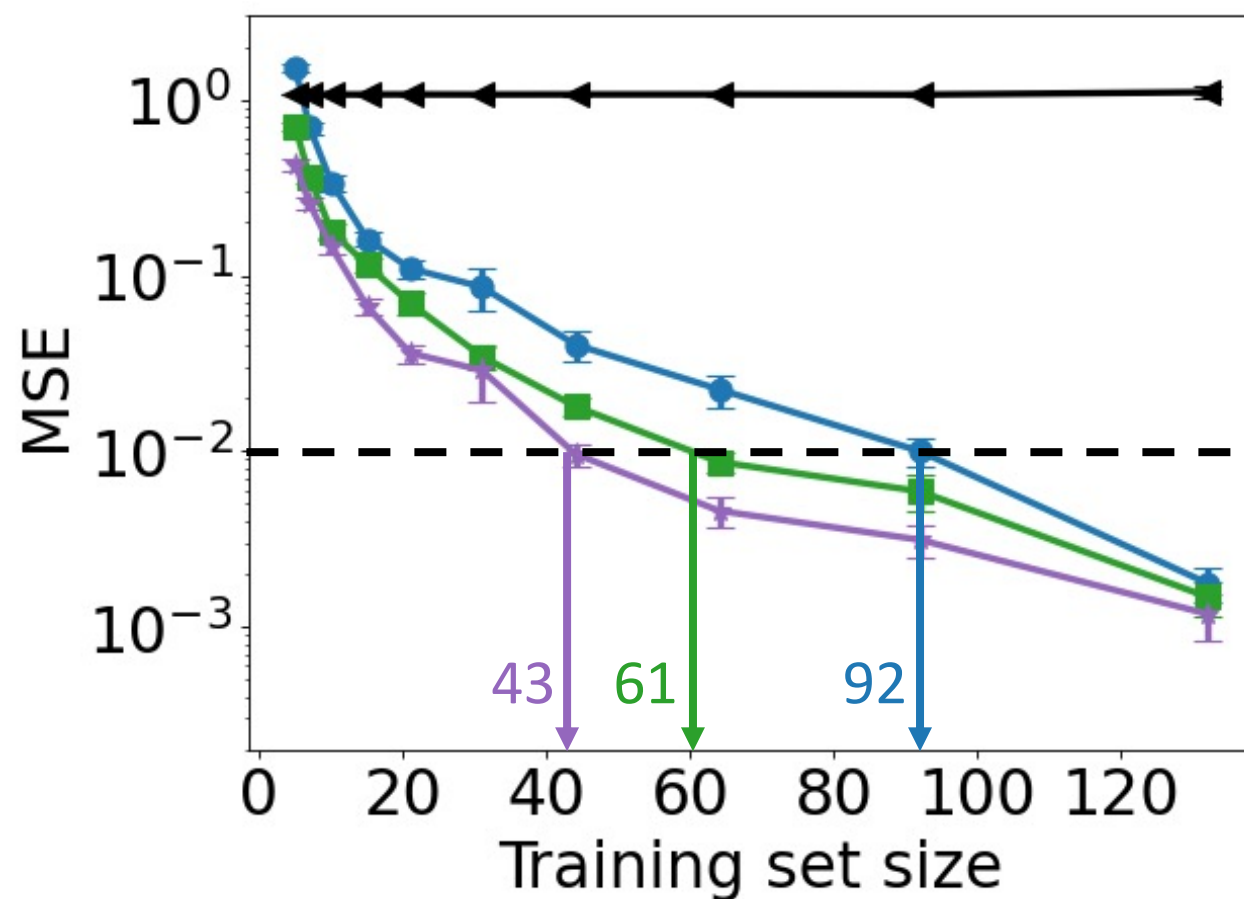
**Linear prior**     $\ln R_g^2 = \mathcal{GP}(\kappa + 2\nu \ln N + \lambda\alpha, k(x, x'))$

**Parameterization**     $\ln R_g^2 = 2\nu[\alpha] \ln N + \kappa[\ln N, \alpha]$

Respects a scaling exponent dependent on solvent quality

**Best performance when full functional form of theory is used**

# Need substantially less data with theory

**Direct** $\quad x \equiv [\ln N, \alpha] \rightarrow \boxed{\text{GPR}} \rightarrow \ln R_g^2$

**Theory** $\quad t(x) = \ln N + \ln 1/6 = \ln R_g^2$

**Latent variable**
$$\begin{array}{c} x \\ t(x) \end{array} \rightarrow \boxed{\text{GPR}} \rightarrow \ln R_g^2$$
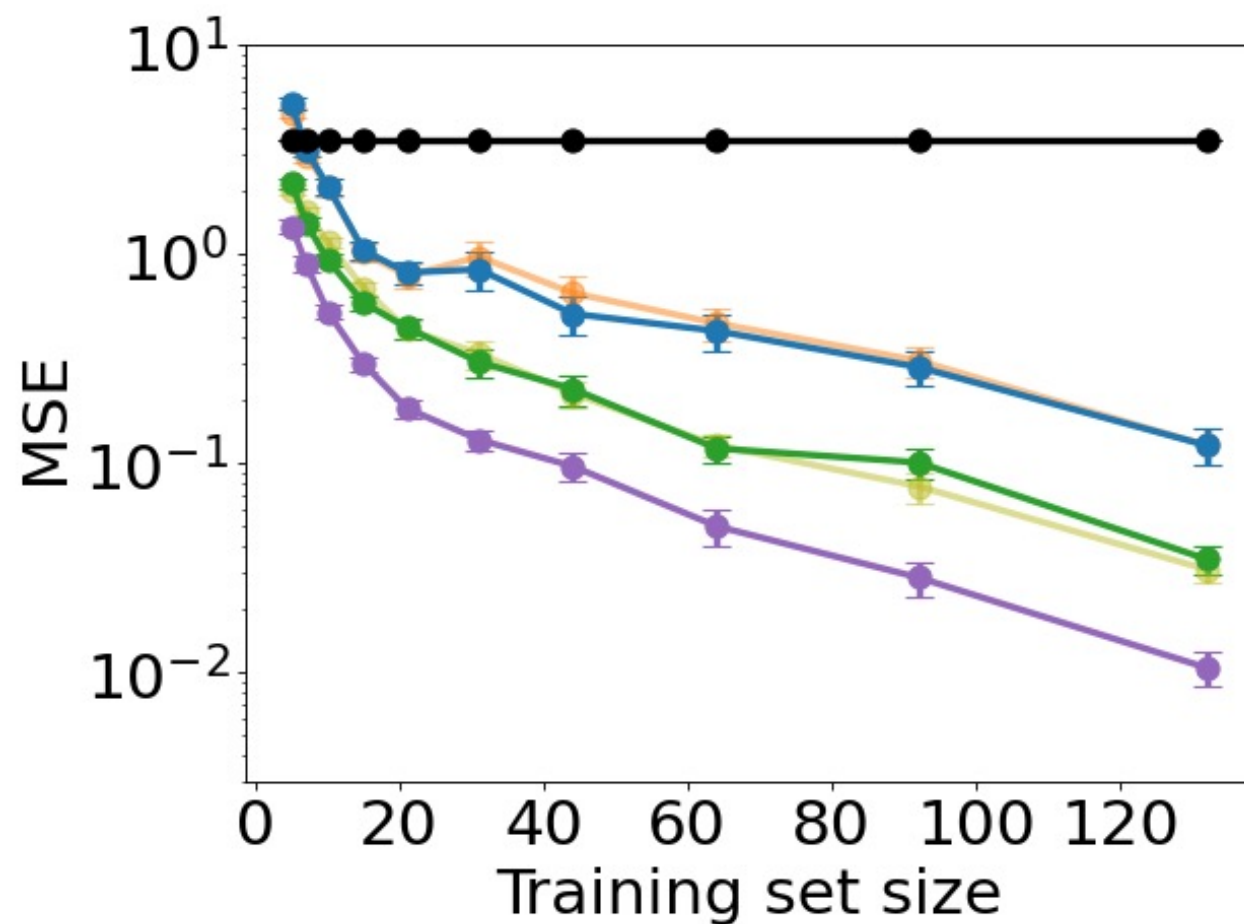
**Difference** $\quad x \rightarrow \boxed{\text{GPR}} \rightarrow \ln R_g^2 - t(x) \quad t(x) \rightarrow \boxed{\text{sum}} \rightarrow \ln R_g^2$

**Linear prior** $\quad \ln R_g^2 = \mathcal{GP}(\kappa + 2\nu \ln N + \lambda\alpha, k(x, x'))$

**Parameterization** $\quad \ln R_g^2 = 2\nu[\alpha]\ln N + \kappa[\ln N, \alpha]$

# Testing extrapolation (larger $N$)

**Direct**     $x \equiv [\ln N, \alpha] \rightarrow$ [GPR] $\rightarrow \ln R_g^2$

**Theory**     $t(x) = \ln N + \ln 1/6 = \ln R_g^2$

**Latent variable**     $\begin{matrix} x \\ t(x) \end{matrix} \rightarrow$ [GPR] $\rightarrow \ln R_g^2$

**Difference**     $x \rightarrow$ [GPR] $\rightarrow \ln R_g^2 - t(x) \rightarrow$ [sum] $\rightarrow \ln R_g^2$,   $t(x)$

**Linear prior**     $\ln R_g^2 = \mathcal{GP}(\kappa + 2\nu \ln N + \lambda\alpha, k(x, x'))$

**Parameterization**     $\ln R_g^2 = 2\nu[\alpha] \ln N + \kappa[\ln N, \alpha]$

Incorporating theory improves extrapolation

Audus *et. al.*, *ACS Macro Letters* **2023**, 11, 1117-1122; https://github.com/usnistgov/taml

# Interpretability for parameterization

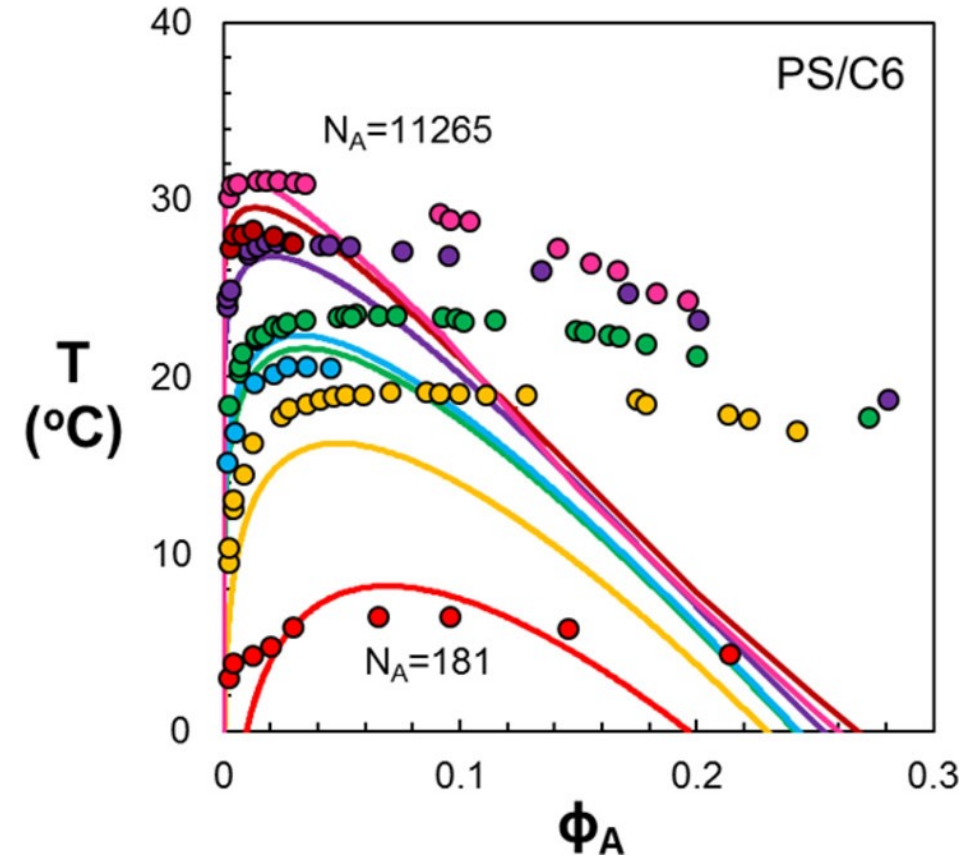$$\ln R_g^2 = 2\nu\,[\alpha]\ln N + \kappa[\ln N, \alpha]$$

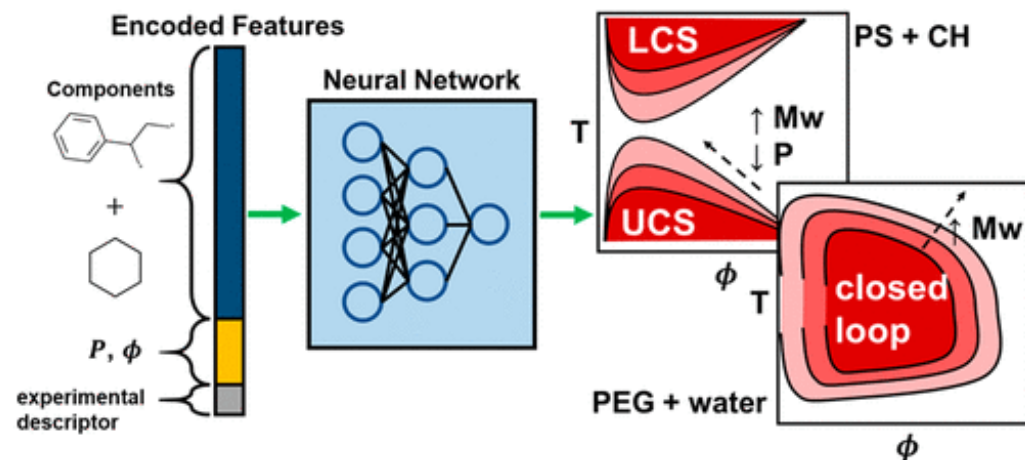Get a prediction for the scaling exponent that includes known limits

Solid lines = prediction; Shaded regions = 95% confidence intervals

$$\frac{\Delta G}{n_T k_B T} = \frac{\phi_A}{N_A} \log \phi_A + \frac{\phi_B}{N_B} \log \phi_B + \chi_{AB} \phi_A \phi_B$$

Configurational Entropy
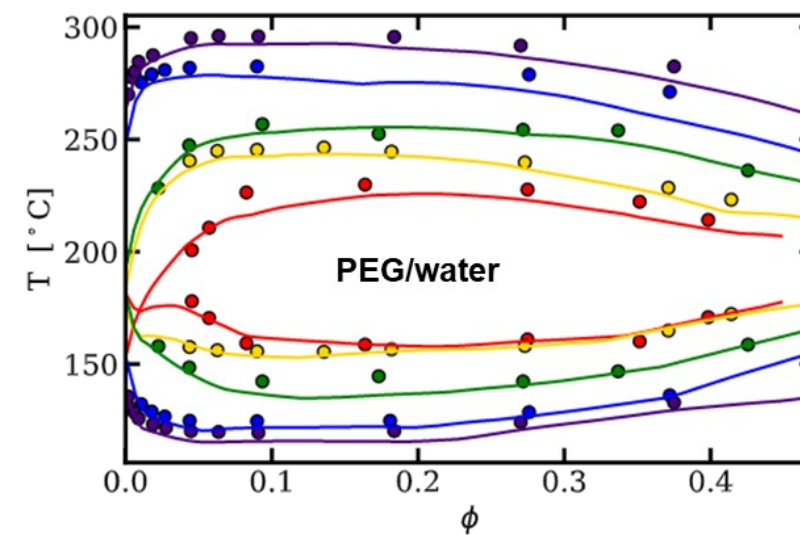
Pairwise Interactions (Enthalpy)

Current predictive models – poor global agreement, requires empirical parameters (e.g. $\chi$(T, p))

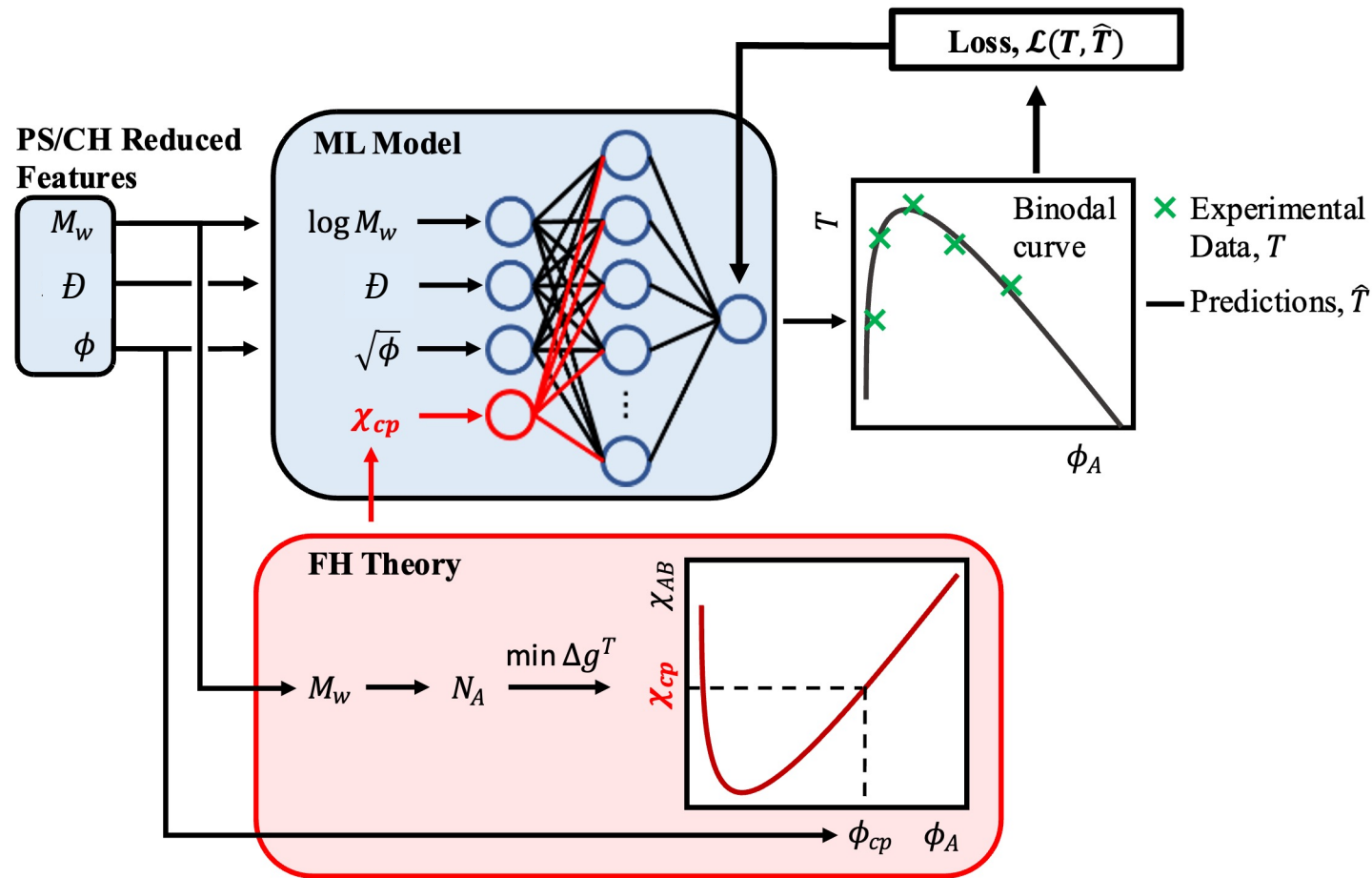Jeffrey Ethier, Devin Ryan, Richard Vaia (AFRL)

# Data driven predictions



Can we improve predictions with theory?

# Theory informed models



**Baseline Model (Prior Knowledge)**

Inputs: $\log M_w, Đ, \sqrt{\phi}$
Outputs: $\hat{T}$

**$\chi$-Informed Model**

Inputs: $\chi_{cp}, \log M_w, Đ, \sqrt{\phi}$
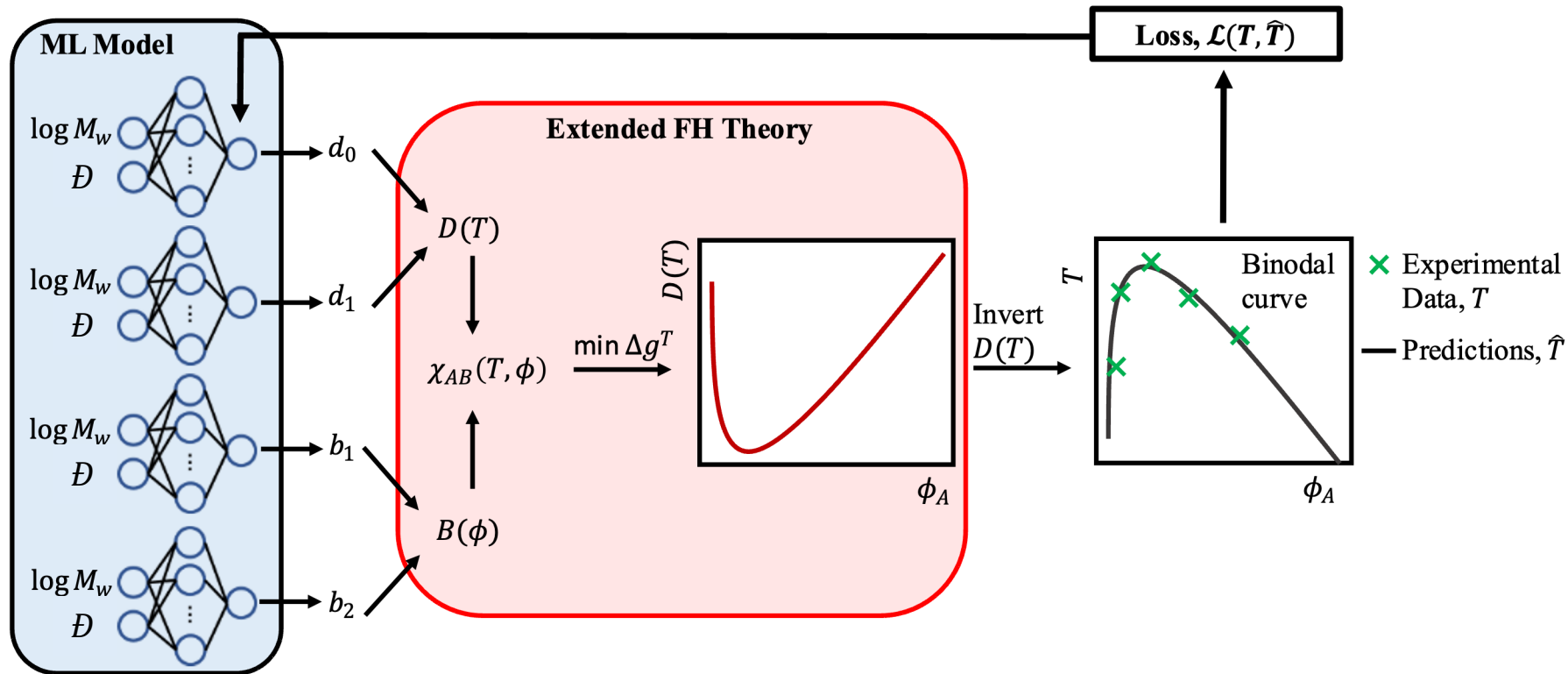Outputs: $\hat{T}$

**Chi2T Model**

Inputs: $\chi_{cp}$
Outputs: $\hat{T}$

# Theory constrained model



$$\frac{\Delta G}{n_T k_B T} = \frac{\phi_A}{N_A} \ln \phi_A + \frac{(1 - \phi_A)}{N_B} \ln(1 - \phi_A) + \phi_A \int_{\phi_A}^{1} \chi_{AB}(T, \phi) d\phi$$
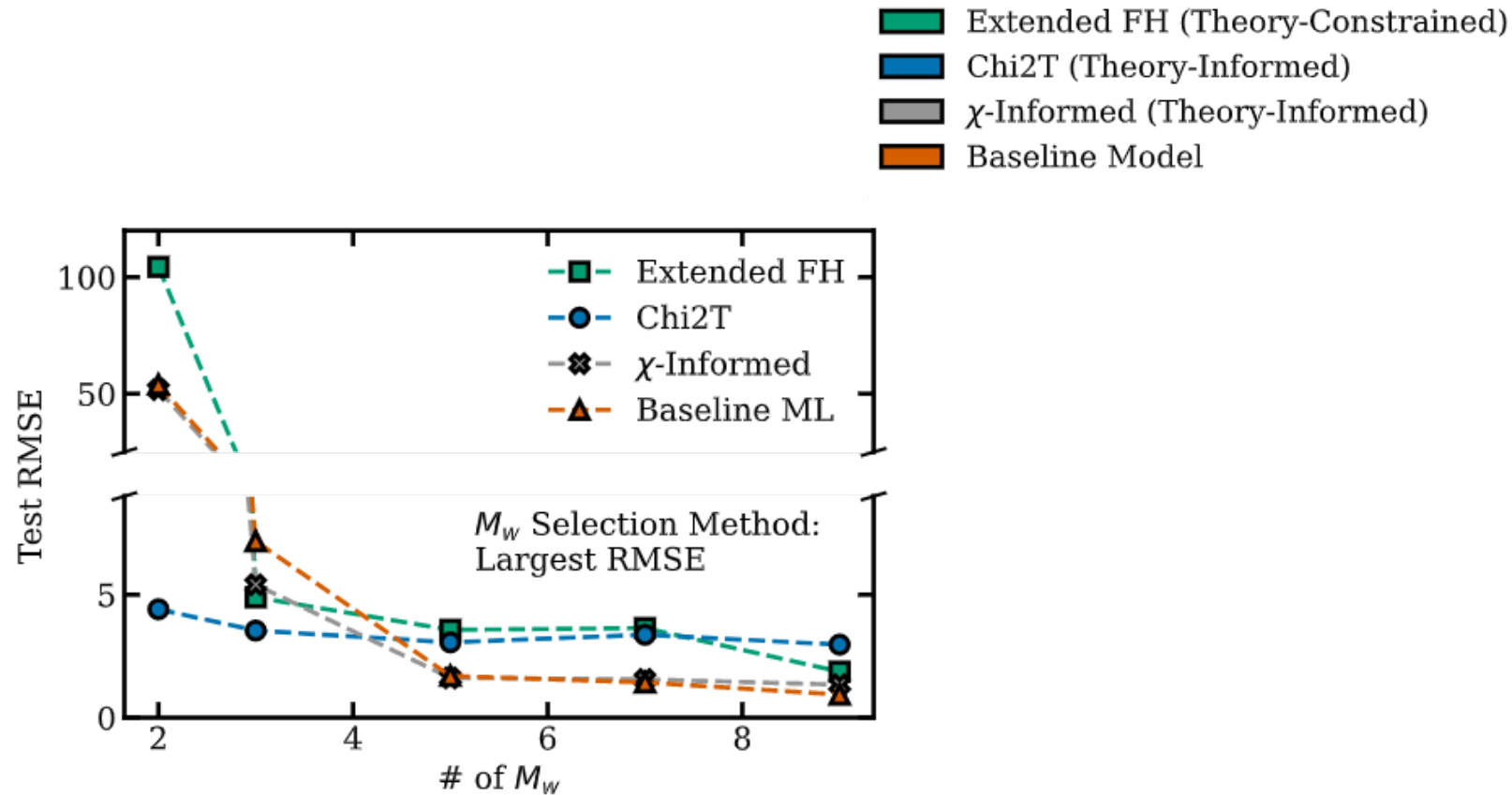
$$D(T)B(\phi_A)$$

$$d_0 + \frac{d_1}{T}$$

$$1 + b_1 \phi_A + b_2 \phi_A^2$$

**Extended FH Model**

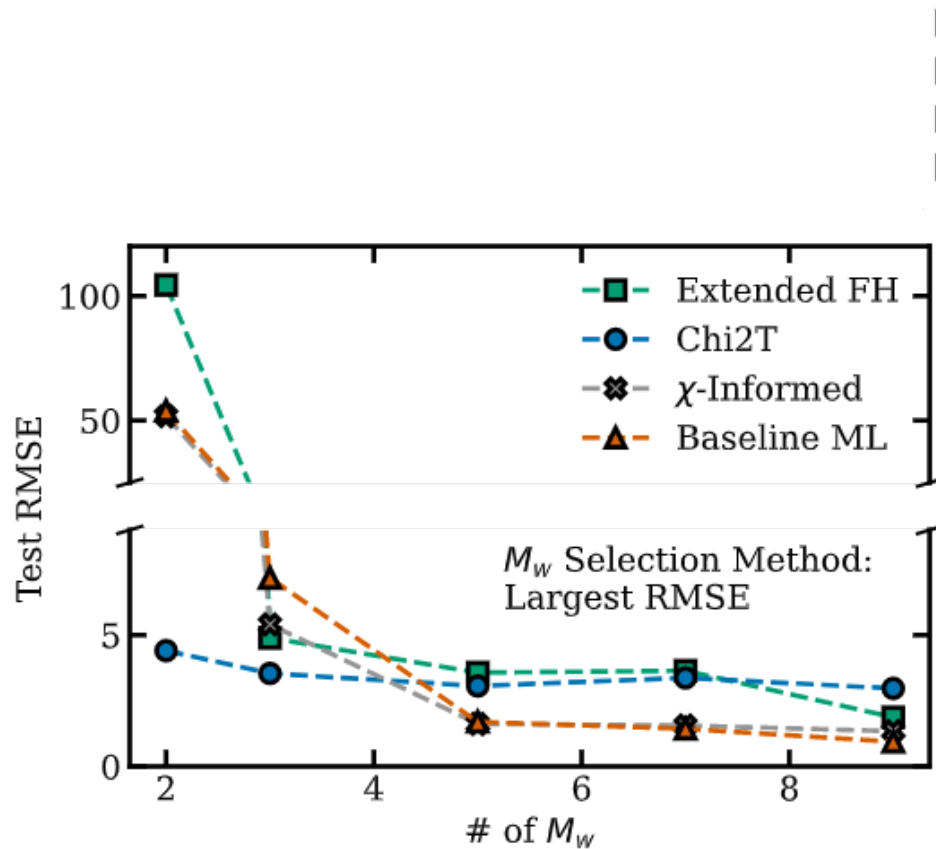Inputs: $\log M_w$, Đ
Output: $\hat{T}$

Chi2T works best for extreme data scarcity
$\chi$-informed model works for intermediate
and then saturates

χ-informed model is best for extrapolation

# Interpretability: Critical Points

**Extended FH Model**

$\chi$**-Informed Model\***



9 M$_W$s
Đ < 1.2
0.0041
$\theta = 305.8\ K$

$\theta = 306.3\ K$
0.0041
all data
Đ < 1.2

$T_c \sim \dfrac{1}{\sqrt{N_A}} + \dfrac{1}{2N_A}$

$\theta_{exp} = 307.2$ (Shultz & Flory)

Able to predict $T_c$

Ethier, Audus, et al. *Giant* **2023**, 15, 100171

*estimated at largest *T* in predicted cloud point curves

# Interpretability: Critical Points

**Extended FH Model**

$\chi$-**Informed Model***



9 $M_W$s
$Đ < 1.2$
$\theta = 305.8\ K$
0.0041

$\theta = 306.3\ K$
all data
$Đ < 1.2$
0.0041

$$T_c \sim \frac{1}{\sqrt{N_A}} + \frac{1}{2N_A}$$

$\theta_{exp} = 307.2$ (Shultz & Flory)

Able to predict $T_c$

---- -0.39log $N$ - 0.12 (log $N$ > 3.5)

---- -0.36log $N$ - 0.15 (log $N$ > 3.5)

$\phi_c \sim N^{-0.38}$  for large $N$

Able to capture scaling

Ethier, Audus, et al. *Giant* **2023**, 15, 100171

*estimated at largest $T$ in predicted cloud point curves

**NIST**

- Theory can reduce the data burden, improve extrapolation and provide explainability
- Many methods exist for embedding theory
- Need to be careful not to make theory overly complex
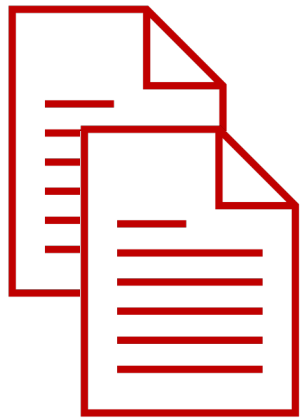
Next steps:

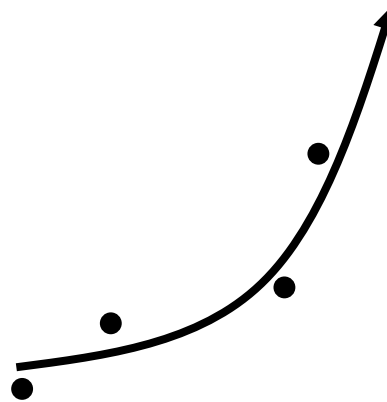- Explore extrapolation with other chemistries and polymer topologies





$\theta = 306.3\ K$

Audus *et. al.*, *ACS Macro Letters* **2023**, 11, 1117-1122;
https://github.com/usnistgov/taml

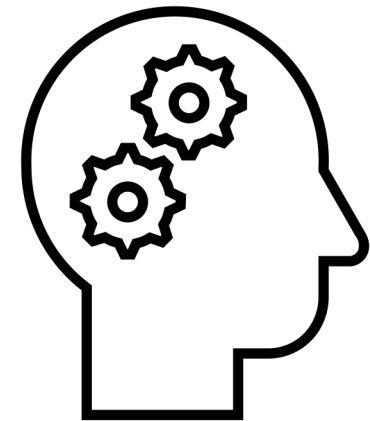Ethier, Audus, et al. *Giant* **2023**, 15, 100171
https://pppdb.uchicago.edu

DATA
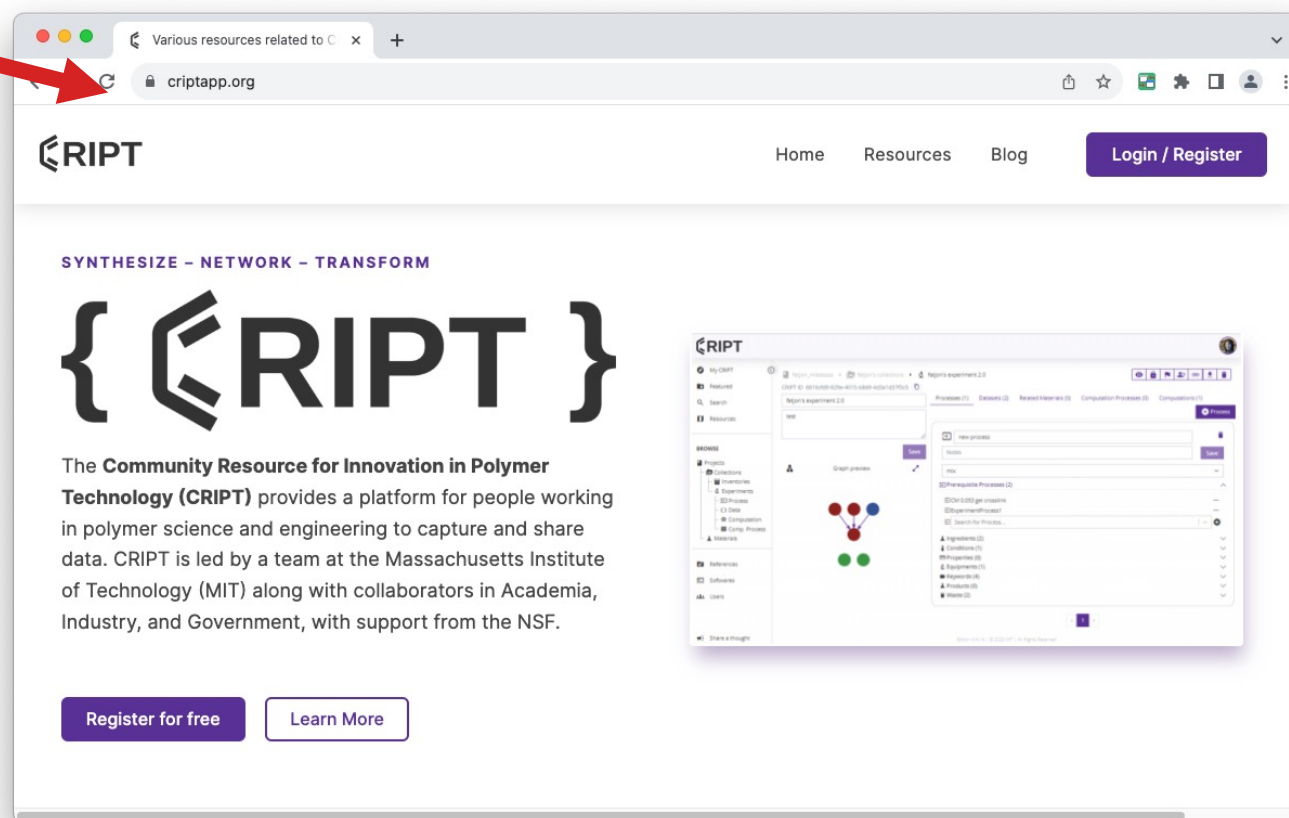(fuel for ML/AI)

EXTRAPOLATION
(go beyond the
dataset)

EXPLAINABILITY
(answer scientific
questions)

# Community Resource for Innovation in Polymer Technology (CRIPT)

https://criptapp.org

Goal: enable accelerated discovery

Bradley Olsen,* Dylan Walsh, Weizhong Zou, Nathan Rebello, Jiale Shi, Michael Deagan, Bruno Leao, Tzyy-Shyang Lin (MIT),

Kaoru Aou (Dow), Ken Kroenlein (Citrine Informatics), Juan de Pablo, Ludwig Schneider, Joshua Mysona (U. Chicago)

Debra Audus (NIST), Ardiana Osmani (CRIPT project leader), and the CRIPT Development Team

# CRIPT publications

Data model: Walsh,…, Audus, *et al. ACS Cent Sci* **2023**, 9, 3, 330-338

CRIPT overview: Deagen,…, Audus *et al. Cell Reports Physical Science* **2022**, 3, 101126

BigSMILES: Lin *et al. ACS Cent Sci* **2019** 5, 9, 1523-1531

Scheme generation: Deagen *et al. Macromolecules* **2023** 57, 1, 42-53

Search: Rebello et al. *J. Chem. Inf. Model.* 2023 63, 6555-6568

Similarity: Shi, …, Audus, Olsen, *Macromolecules* **2023** 56, 18, 7344-7357