**Molecular science tasks**

biomolecule understanding

drug design

material/catalyst discovery

...

# Computational Methods for Molecular Science



**Molecular science tasks**

Macroscropic properties ← Thermodynamic ensemble ← Microscopic properties ← Electronic structure

- biomolecule understanding
- drug design
- material/catalyst discovery
- ...

- binding affinity
- stability
- activation energy
- conductivity

$E$

$\epsilon_\Delta$

$\mathbf{f}$

$\psi(\mathbf{r}_1, \cdots, \mathbf{r}_N)$

$\{\phi_i(\mathbf{r}_j)\}$

$\rho(\mathbf{r})$

# Computational Methods for Molecular Science

北京中关村学院 x 中关村人工智能研究院
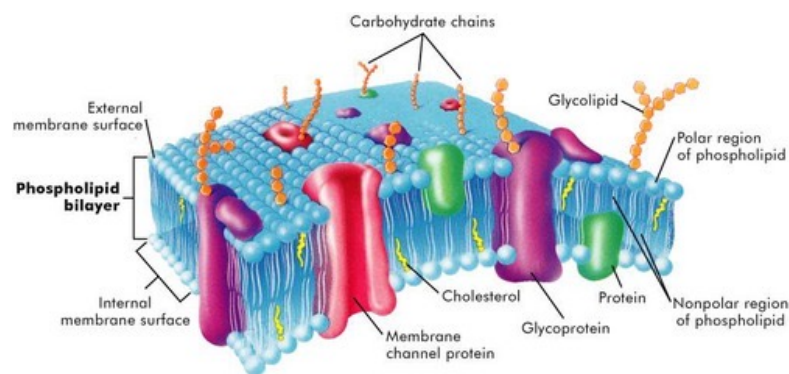Zhongguancun Academy · Zhongguancun Institute of Artificial Intelligence
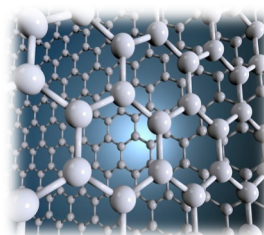
**Molecular science tasks**

biomolecule understanding

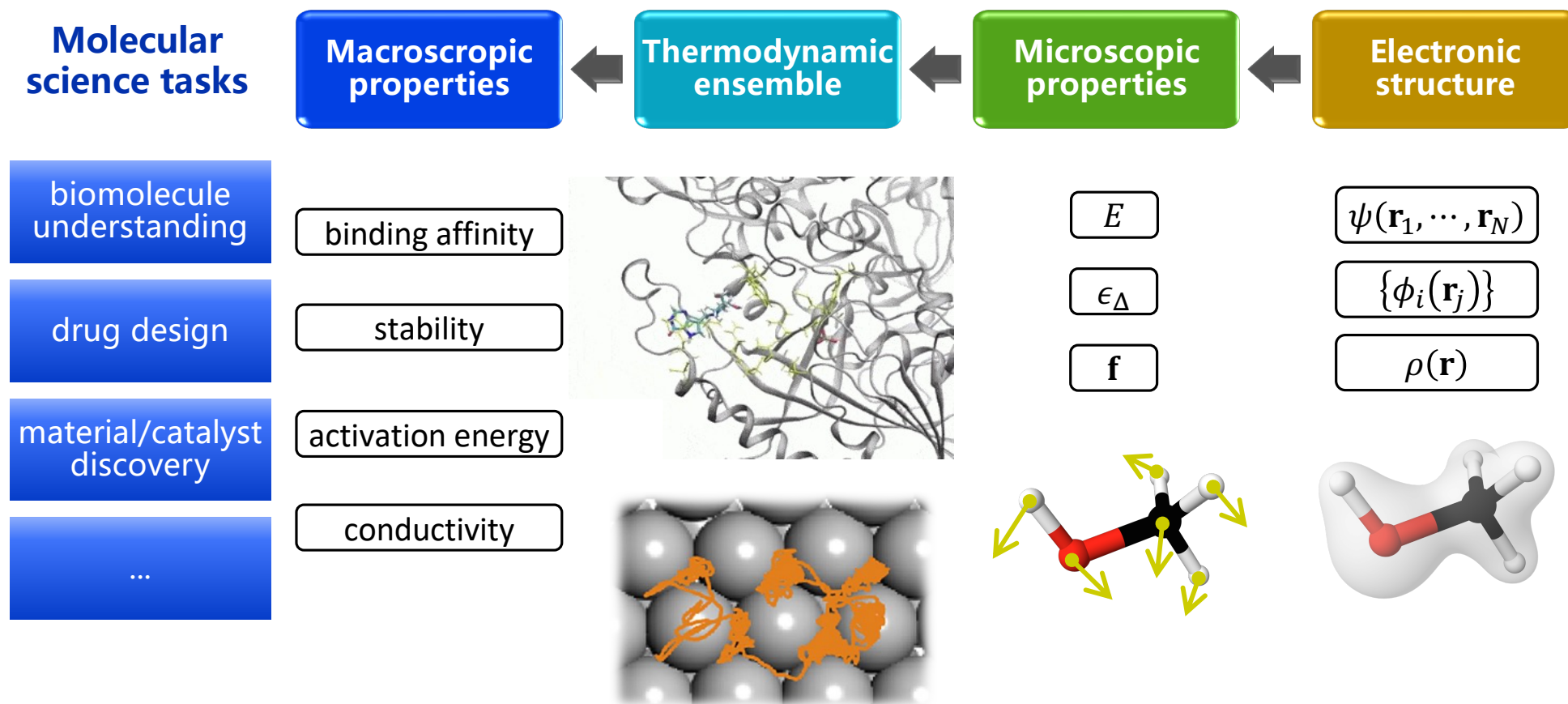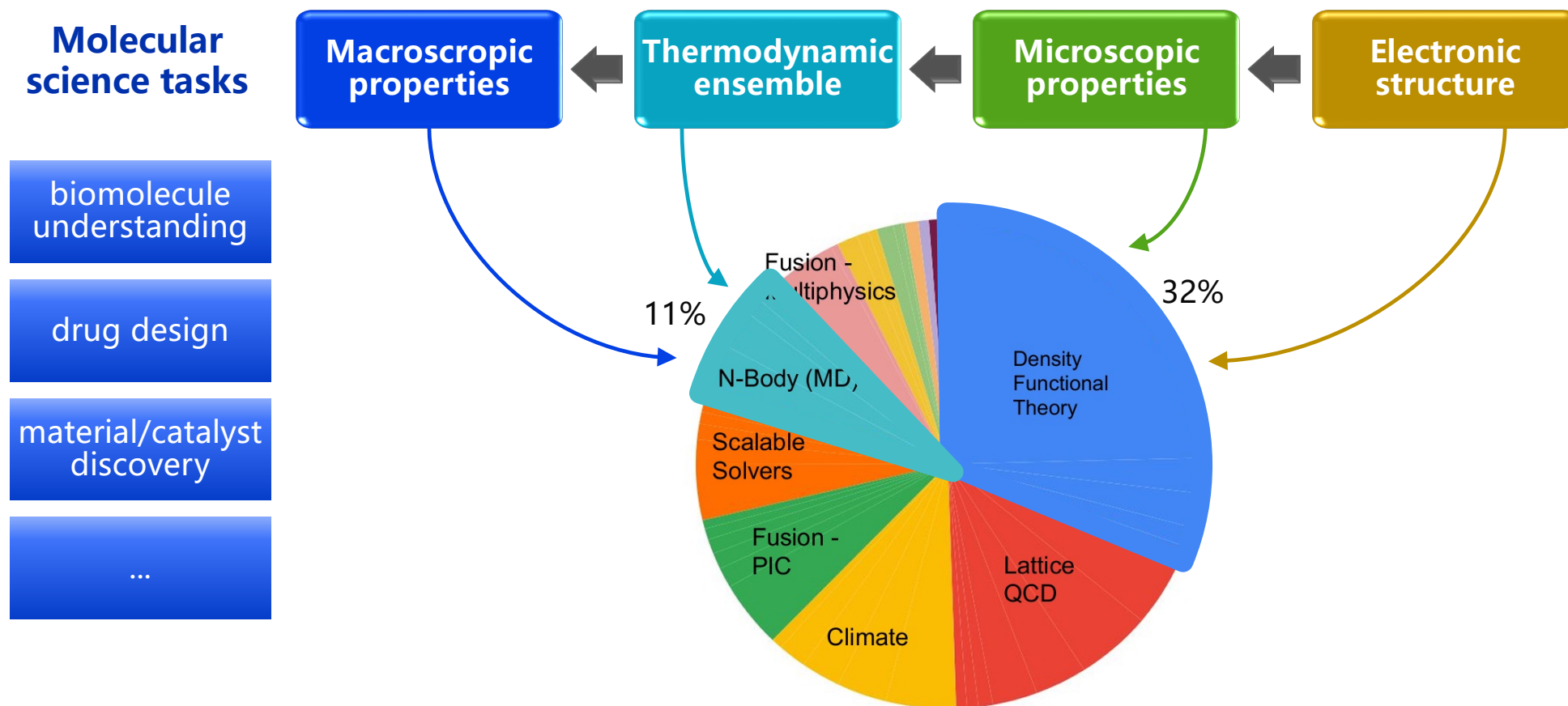drug design

material/catalyst discovery

...

| Macroscropic properties | ← | Thermodynamic ensemble | ← | Microscopic properties | ← | Electronic structure |



Fusion - Multiphysics

11%

N-Body (MD,

Scalable Solvers

Fusion - PIC

Climate

Density Functional Theory

32%

Lattice QCD

# AI for Scientific Computation

**Macroscropic properties**

**Thermodynamic ensemble**

**Microscopic properties**

**Electronic structure**

**Macroscopic property prediction**: scalability

**Structure prediction and sampling**: accuracy & scalability

**ML force field**: accuracy

**Hamiltonian prediction**: scalability

**Learned functional**: accuracy

# The Curse: Data Scarcity



system domain

level of accuracy

⚫ : available dataset

chemical space

property space (energy, structure, Hamiltonian, ...)

conformational space

**Even more axes**: system scale, temperature/pressure, reaction, solution environments, ...

# The Blessing: Physical Laws!

| Macroscropic properties | ← | Thermodynamic ensemble | ← | Microscopic properties | ← | Electronic structure |
|---|---|---|---|---|---|---|

Statistics

Statistical mechanics

Quantum mechanics

Schrödinger equation

$$\Delta G_{21} = \log \frac{p_1}{p_2}$$

$$p(\mathbf{R}) \propto \exp\left(-\frac{E(\mathbf{R})}{k\mathcal{T}}\right)$$

$$E(\mathbf{R}) = E[\psi_{\mathbf{R}}]$$

$$\widehat{H}_{\mathbf{R}}\psi = E\psi$$



Free energy landscape

open | closed

# Electronic Structure

Electron structure

?

$e^-$
$\mathbf{r}_1$

$e^-$
$\mathbf{r}_2$

$Z_1$      $Z_2$

$\mathbf{R}_1$      $\mathbf{R}_2$

Molecular structure
$\mathbf{R}$

Schrödinger equation:
$$\widehat{H}_{\mathbf{R}}\psi(\mathbf{r}_1,\cdots,\mathbf{r}_N) = E\psi(\mathbf{r}_1,\cdots,\mathbf{r}_N)$$
$$\rightarrow E = \min_{\psi}\{E_{\mathbf{R}}[\psi] := \langle\psi|\widehat{H}|\psi\rangle\}$$

$\psi_{\mathbf{R}}^{\star}$

Density functional theory:
$$E = \min_{\mathbf{C}:\text{ orthonormal}} E_{\mathbf{R}}(\mathbf{C})$$
$$\rightarrow \underbrace{\mathbf{H}_{\mathbf{R}}(\mathbf{C})}\ \mathbf{C} = \underbrace{\mathbf{S}}\ \underbrace{\mathbf{C}}\ \underbrace{\boldsymbol{\epsilon}}$$

$\nabla_{\mathbf{C}}E_{\mathbf{R}}(\mathbf{C})$  $\nabla_{\mathbf{C}}$constraint  Lagrange multiplier

Self-Consistent Field Iteration
$$\mathbf{H}^{(1)}\mathbf{C} = \mathbf{SC}\boldsymbol{\epsilon} \qquad \mathbf{H}^{(2)}\mathbf{C} = \mathbf{SC}\boldsymbol{\epsilon}$$

$\mathbf{C}^{(0)}$   $\mathbf{H}^{(1)}$   $\mathbf{C}^{(1)}$   $\mathbf{H}^{(2)}$   ...

Calculate
$\mathbf{H}_{\mathbf{R}}(\mathbf{C}^{(0)})$

Calculate
$\mathbf{H}_{\mathbf{R}}(\mathbf{C}^{(1)})$

Hamiltonian solution
$\mathbf{H}_{\mathbf{R}}^{\star}$

$E$

$\epsilon_{\text{HOMO}}$

$\epsilon_{\text{LUMO}}$

$\epsilon_{\Delta}$

...

Physical-law Training!

**Hamiltonian Prediction**
$\widehat{\mathbf{H}}_{\theta}(\mathbf{R})$

Supervised Training

Zhang et al. Self-Consistency Training for Density-Functional-Theory Hamiltonian Prediction. *ICML*, 2024.

# Physical-Law Training for Electronic Structure

- Self-consistency training for Hamiltonian prediction:

$$L_{\text{self\_con}}(\theta) = \left\| \widehat{\mathbf{H}}_\theta(\mathbf{R}) - \mathbf{H_R}\left( \mathbf{C_R}\left( \widehat{\mathbf{H}}_\theta(\mathbf{R}) \right) \right) \right\|_{\text{F}}^2.$$

- **Label-free**: distinction from predicting other properties.
- Not just a regularization: it fully **determines** the solution.

· Generalization beyond available data

Largest labeled dataset (≤ 31 atoms) + unlabeled molecules → test on larger molecules

| Molecule | Setting | Hamiltonian | Derived properties | | | | |
|---|---|---|---|---|---|---|---|
| | | $\mathbf{H}\,[\mu E_h]\downarrow$ | $\epsilon\,[\mu E_h]\downarrow$ | $\mathbf{C}\,[\%]\uparrow$ | $\epsilon_{\text{HOMO}}\,[\mu E_h]\downarrow$ | $\epsilon_{\text{LUMO}}\,[\mu E_h]\downarrow$ | $\epsilon_{\Delta}\,[\mu E_h]\downarrow$ |
| ALA3 (42 atoms) | zero-shot | 237.71 | $6.54\times10^3$ | 52.24 | $6.90\times10^3$ | $9.51\times10^4$ | $9.79\times10^4$ |
| | self-con | **52.49** | $\mathbf{1.22\times10^3}$ | **94.46** | $\mathbf{2.07\times10^3}$ | $\mathbf{3.76\times10^3}$ | $\mathbf{2.69\times10^3}$ |
| DHA (56 atoms) | zero-shot | 397.87 | $1.84\times10^4$ | 20.15 | $1.11\times10^4$ | $1.90\times10^5$ | $1.85\times10^5$ |
| | self-con | **56.12** | $\mathbf{1.81\times10^3}$ | **83.51** | $\mathbf{1.99\times10^3}$ | $\mathbf{4.01\times10^3}$ | $\mathbf{2.34\times10^3}$ |

北京中关村学院 x 中关村人工智能研究院
Zhongguancun Academy   Zhongguancun Institute of Artificial Intelligence

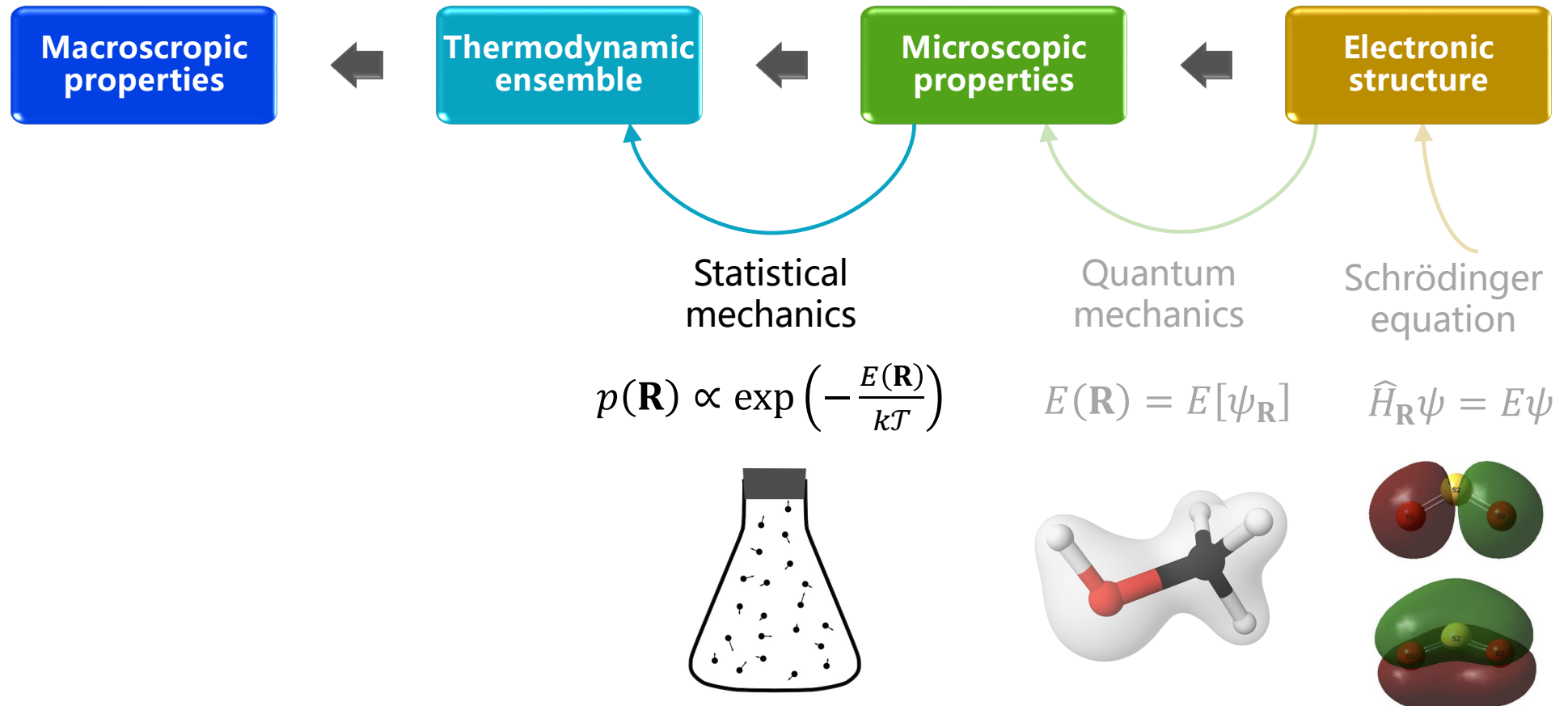· **Amortization effect**: more efficient than running DFT to generate labels.

- Passing physical-law information to microscopic properties

labeled dataset ($\leq$ 31 atoms) + unlabeled molecules $\rightarrow$ test on larger molecules

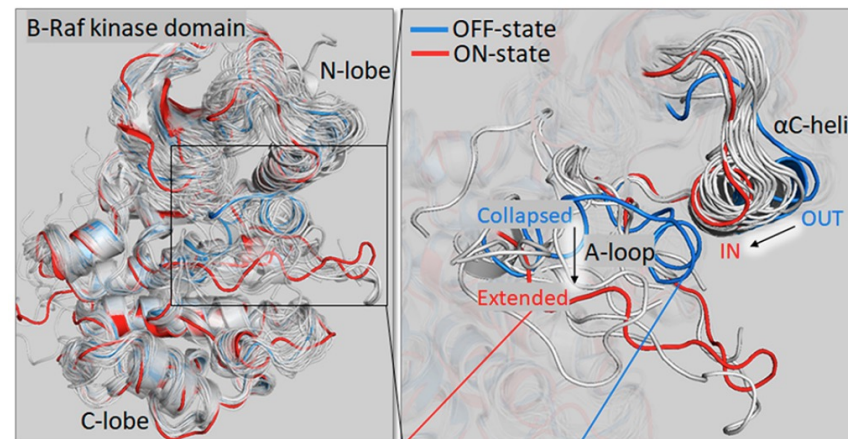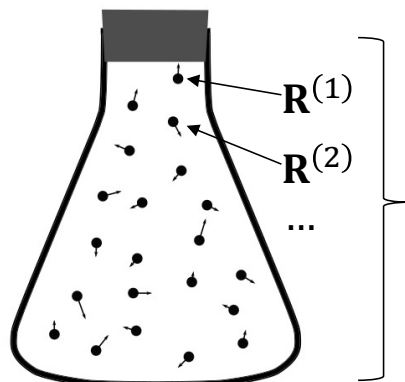| Molecule | Setting | $\epsilon_{HOMO}$ [$\mu E_h$] $\downarrow$ | $\epsilon_{LUMO}$ [$\mu E_h$] $\downarrow$ | $\epsilon_\Delta$ [$\mu E_h$] $\downarrow$ |
|---|---|---|---|---|
| ALA3 (42 atoms) | self-con | $\mathbf{2.07 \times 10^3}$ | $\mathbf{3.76 \times 10^3}$ | $\mathbf{2.69 \times 10^3}$ |
| | e2e (ET) | $1.74 \times 10^5$ | $7.72 \times 10^3$ | $2.38 \times 10^5$ |
| | e2e (Equiformer) | $2.38 \times 10^5$ | $1.16 \times 10^4$ | $2.27 \times 10^5$ |
| DHA (56 atoms) | self-con | $\mathbf{1.99 \times 10^3}$ | $\mathbf{4.01 \times 10^3}$ | $\mathbf{2.34 \times 10^3}$ |
| | e2e (ET) | $2.92 \times 10^5$ | $2.58 \times 10^4$ | $3.39 \times 10^5$ |
| | e2e (Equiformer) | $3.76 \times 10^5$ | $2.31 \times 10^4$ | $4.17 \times 10^5$ |

# The Blessing: Physical Laws!

| Macroscropic properties | ← | Thermodynamic ensemble | ← | Microscopic properties | ← | Electronic structure |
|---|---|---|---|---|---|---|

Statistical mechanics

Quantum mechanics

Schrödinger equation

$$p(\mathbf{R}) \propto \exp\left(-\frac{E(\mathbf{R})}{k\mathcal{T}}\right)$$

$$E(\mathbf{R}) = E[\psi_{\mathbf{R}}]$$
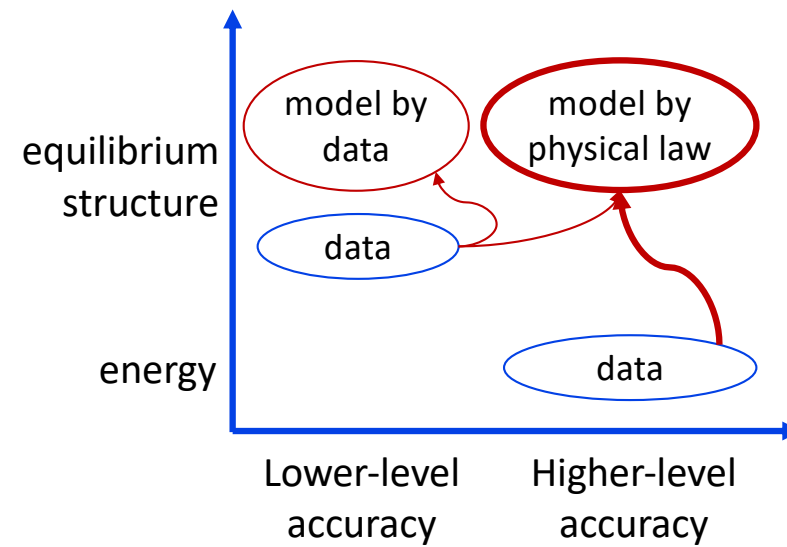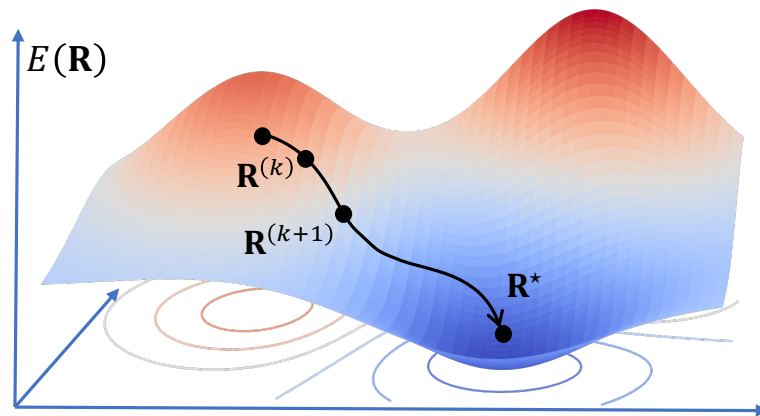
$$\hat{H}_{\mathbf{R}}\psi = E\psi$$

# Thermodynamic Ensemble

- A molecule exists in real world with structures **R** following a distribution.

  - More detailed description/knowledge



$\mathbf{R}^{(1)}$

$\mathbf{R}^{(2)}$

…

https://en.wikipedia.org/wiki/
Ensemble_(mathematical_phy
sics)

北京中关村学院 Zhongguancun Academy x 中关村人工智能研究院 Zhongguancun Institute of Artificial Intelligence



$E(\mathbf{R})$

$\mathbf{R}^{(k)}$

$\mathbf{R}^{(k+1)}$

$\mathbf{R}^{\star}$



equilibrium structure

model by data

model by physical law

data

energy

data

Lower-level accuracy

Higher-level accuracy

# Connection between Structure and Energy

- Optimality consistency:

$$\mathbf{R}^\star = \underset{\mathbf{R}}{\arg\min}\, E(\mathbf{R}) \qquad \Rightarrow \qquad \min_\theta \mathbb{E}\, \max\{0,\; E(\mathbf{R}^\star_\theta) - E(\mathbf{R}^\star_\theta + \boldsymbol{\eta})\}.$$

$$\mathbf{D}_\theta(\boldsymbol{\epsilon}, t) \text{ for large } t \approx T.$$

- Score consistency:

$$\mathbf{R}^\star \sim \exp\left(-\frac{E(\mathbf{R})}{k\mathcal{T}}\right) \text{ for small } \mathcal{T} \quad \Rightarrow \quad \min_\theta \mathbb{E}_{\mathbf{R}} \left\| \nabla \log p_\theta(\mathbf{R}) + \frac{\nabla E(\mathbf{R})}{k\mathcal{T}} \right\|^2.$$

$$\frac{\alpha_t\, \mathbf{D}_\theta(\mathbf{R},t) - \mathbf{R}}{\sigma_t^2} \text{ for small } t \approx 0$$

Ren et al. Physical Consistency Bridges Heterogeneous Data in Molecular Multi-Task Learning. *NeurIPS*, 2024.

**Evaluation**:
error w.r.t high-accuracy structure

| Test Set | PCQ | | QM9 | |
|---|---|---|---|---|
| Struct. Stat. | Mean | Min | Mean | Min |
| low-accuracy structure data | 1.189 | 0.655 | 0.928 | 0.545 |
| +high-accuracy energy data | **1.158** | **0.645** | **0.848** | **0.490** |

**Training**: low-accuracy structure data +high-accuracy energy data

Ren et al. Physical Consistency Bridges Heterogeneous Data in Molecular Multi-Task Learning. *NeurIPS*, 2024.

- Data are biased from finite-length simulation

# Connection between Distribution and Energy

- Energy function defines the target: $p_0(\mathbf{R}_0) \propto \exp\left(-\frac{E(\mathbf{R}_0)}{k\mathcal{T}}\right)$.

- Diffusion-model learning target:

$$\mathbf{s}_\theta(\mathbf{R}_t, t) \to \underset{\theta}{\arg\min}\, \mathbb{E}_{\underbrace{p_0(\mathbf{R}_0)}p(\mathbf{R}_t|\mathbf{R}_0)} \left\| \mathbf{s}_\theta(\mathbf{R}_t, t) - \frac{1}{\alpha_t}\nabla\log\underbrace{p_0(\mathbf{R}_0)}\right\|^2.$$

$$\mathbf{s}_{\theta_{\text{debiased}}}(\mathbf{R}_t, t) \to \qquad \underset{\approx\, q(\mathbf{R}_0)}{}\qquad\qquad \underset{-\frac{\nabla E(\mathbf{R}_0)}{k\mathcal{T}}}{}$$

$\approx q(\mathbf{R}_0)$
data distribution

$-\frac{\nabla E(\mathbf{R}_0)}{k\mathcal{T}}$

- Energy function corrects data bias:

$$\left\| \mathbf{s}_{\theta_{\text{debiased}}}(\mathbf{R}_t, t) - \nabla\log p_t(\mathbf{R}_t)\right\|^2 \leq \left\| \mathbf{s}_{\theta_{\text{data}}}(\mathbf{R}_t, t) - \nabla\log p_t(\mathbf{R}_t)\right\|^2.$$

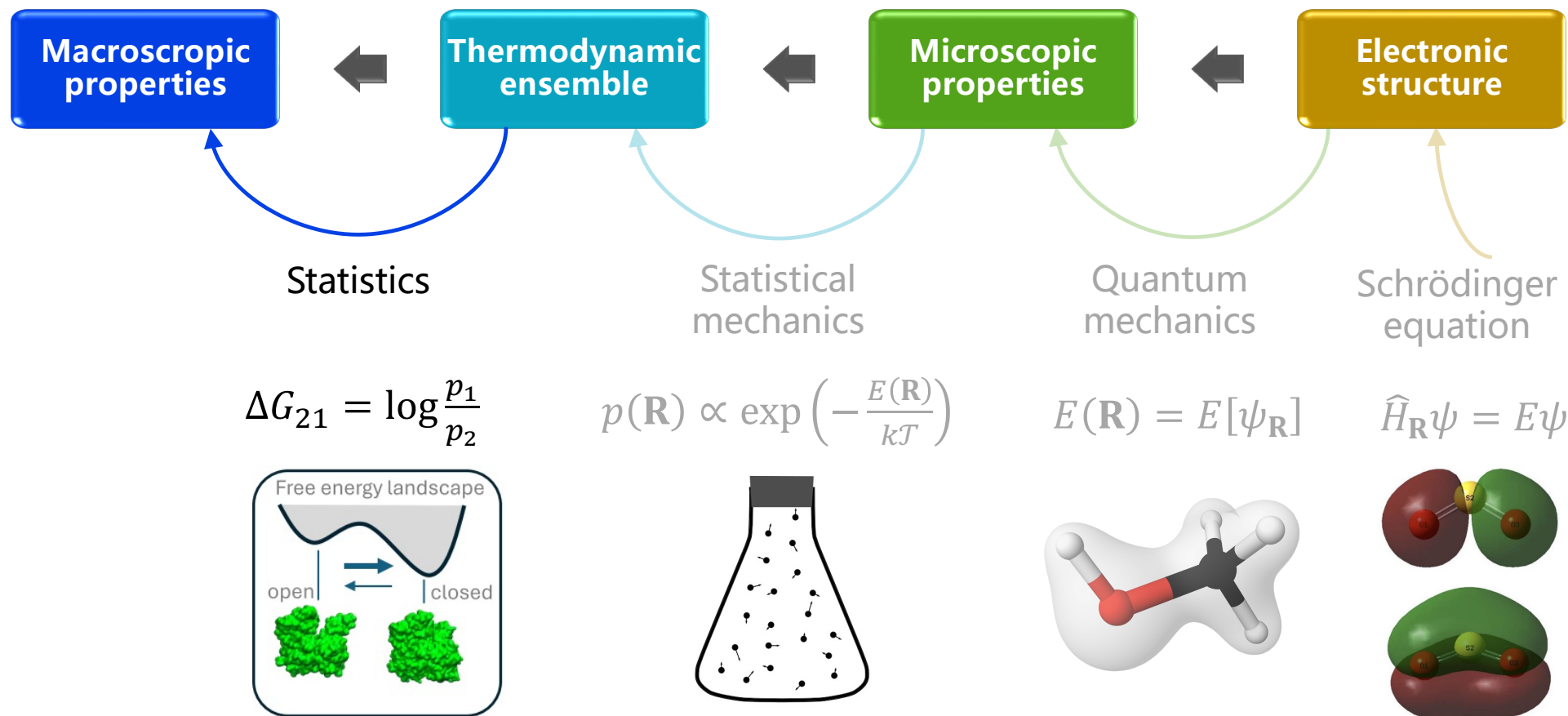Guo et al. Potential Score Matching: Debiasing Molecular Structure Sampling with Potential Energy Guidance. *TMLR*, 2025.

# Accuracy beyond Training Data

Guo et al. Potential Score Matching: Debiasing Molecular Structure Sampling with Potential Energy Guidance. *TMLR*, 2025.

# The Blessing: Physical Laws!

**Macroscropic properties** ← **Thermodynamic ensemble** ← **Microscopic properties** ← **Electronic structure**

Statistics

Statistical mechanics

Quantum mechanics

Schrödinger equation

$$\Delta G_{21} = \log \frac{p_1}{p_2}$$

$$p(\mathbf{R}) \propto \exp\left(-\frac{E(\mathbf{R})}{k\mathcal{T}}\right)$$

$$E(\mathbf{R}) = E[\psi_{\mathbf{R}}]$$

$$\hat{H}_{\mathbf{R}}\psi = E\psi$$



Free energy landscape

open | closed

# Stability Calculation from Ensemble Model

Thermodynamic ensemble $p_\theta(\mathbf{R})$ → $\Delta G$ calculation

Foldedness $x = f(\mathbf{R})$

Regularize/ calibrate

$p_{\text{folded}}$, $p_{\text{unfolded}}$



Train — Error 0.73 kcal/mol, Correlation 0.67

Test — Error 0.76 kcal/mol, Correlation 0.66

**outperforms black-box methods** that directly predict $\Delta G$ from sequences.

Stability $\log \dfrac{p_{\text{folded}}}{p_{\text{unfolded}}}$ ← Folding free energy $\Delta G = \log \dfrac{p_{\text{folded}}}{p_{\text{unfolded}}}$
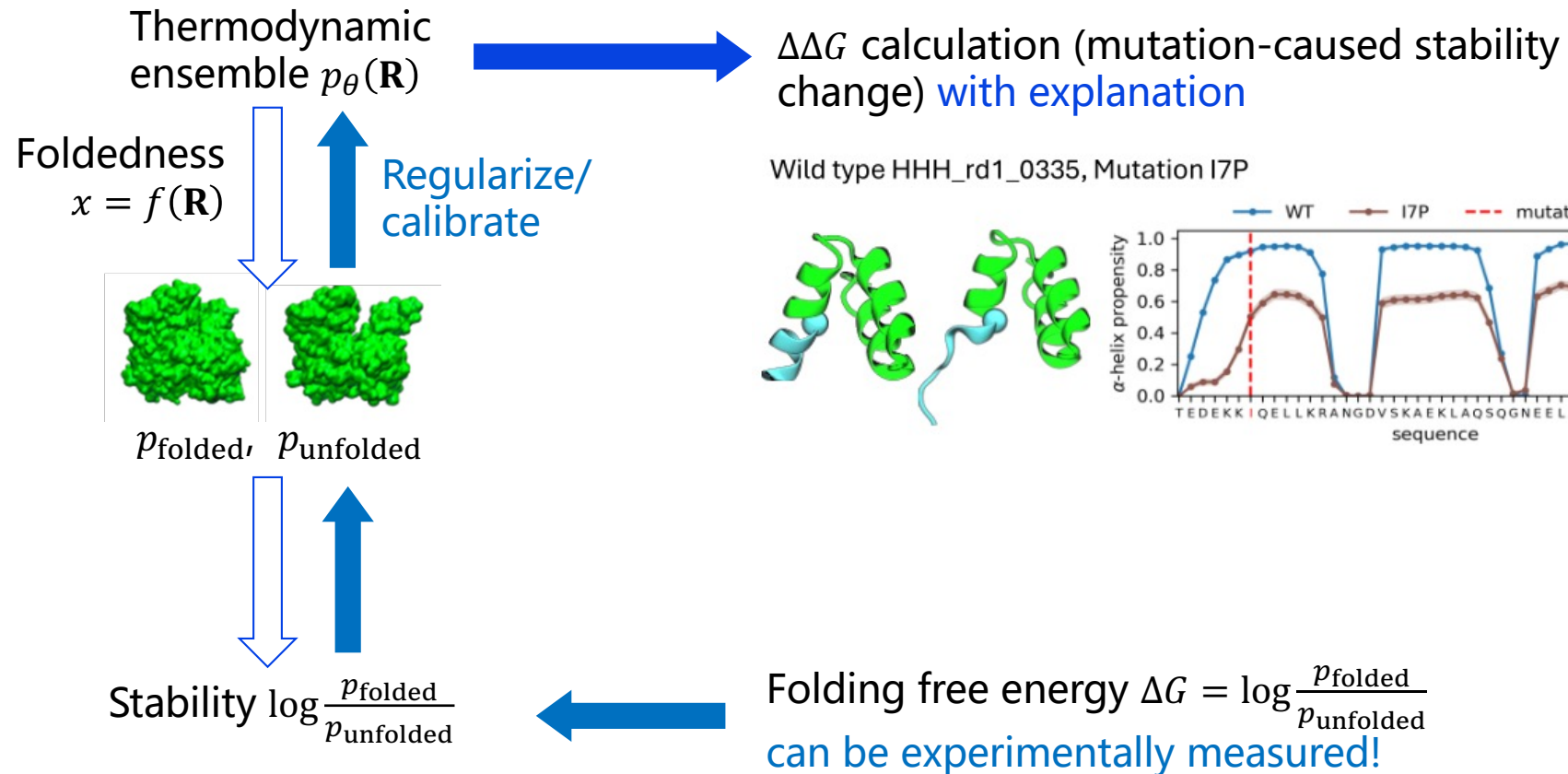
**can be experimentally measured!**

Lewis et al., Scalable emulation of protein equilibrium ensembles with generative deep learning, *Science*, 2025.
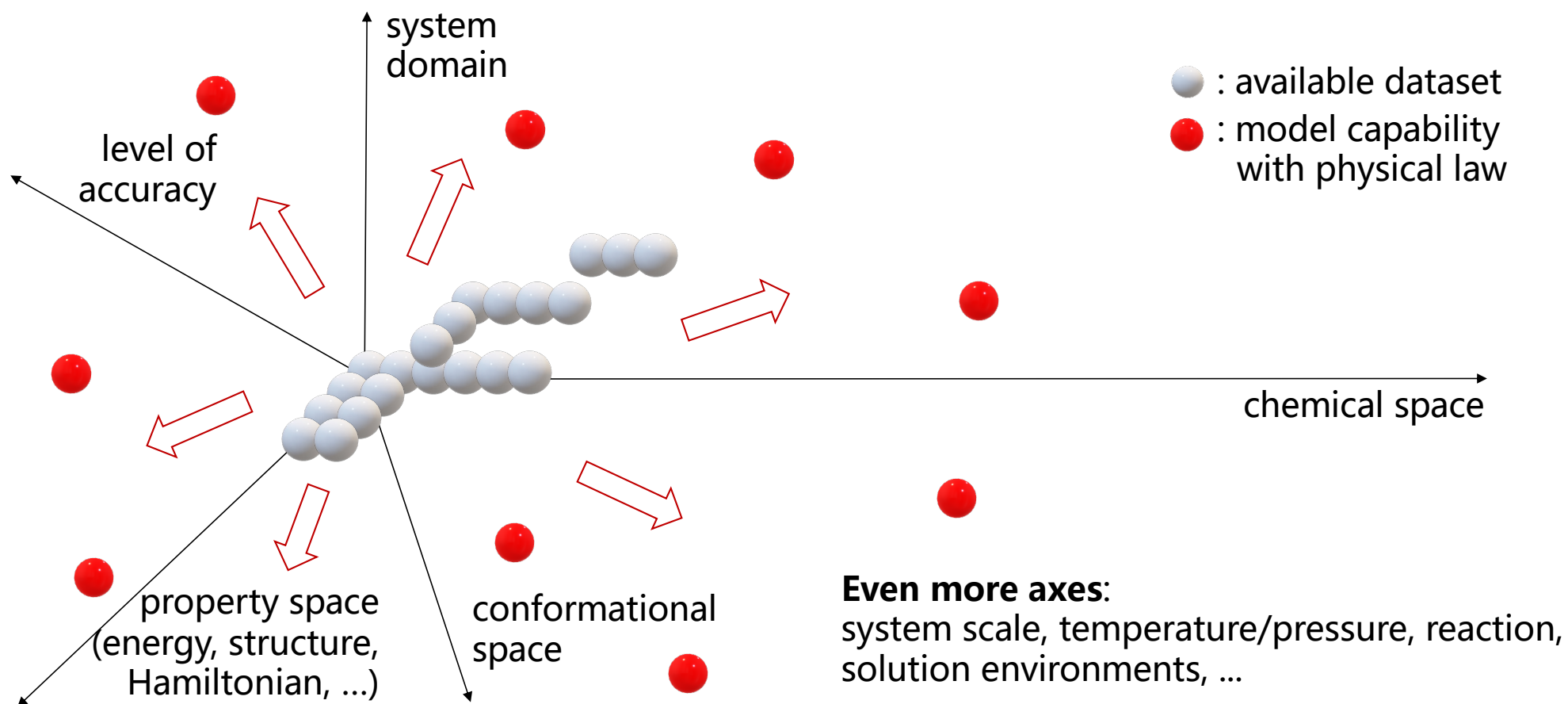
# Stability Calculation from Ensemble Model

Thermodynamic ensemble $p_\theta(\mathbf{R})$ → $\Delta\Delta G$ calculation (mutation-caused stability change) with explanation

Foldedness $x = f(\mathbf{R})$

Regularize/ calibrate

Wild type HHH_rd1_0335, Mutation I7P



$p_{\text{folded}}$, $p_{\text{unfolded}}$

Stability $\log \dfrac{p_{\text{folded}}}{p_{\text{unfolded}}}$ ← Folding free energy $\Delta G = \log \dfrac{p_{\text{folded}}}{p_{\text{unfolded}}}$ can be experimentally measured!

Lewis et al., Scalable emulation of protein equilibrium ensembles with generative deep learning, *Science*, 2025.