

# Explainable Artificial Intelligence for Materials Discovery

**DATAI Conference 2026**  
**Pamplona, Spain**

Valentin Vassilev Galindo

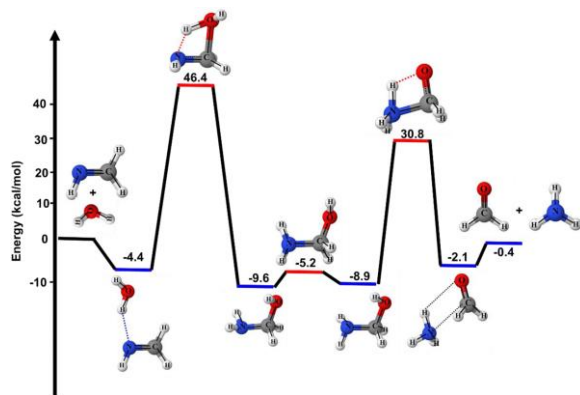
Department of Statistics, Computer Science and Mathematics  
Universidad Pública de Navarra, Pamplona, Spain

Navarra Artificial Intelligence Research Center (NAIR center)

5<sup>th</sup> May 2026

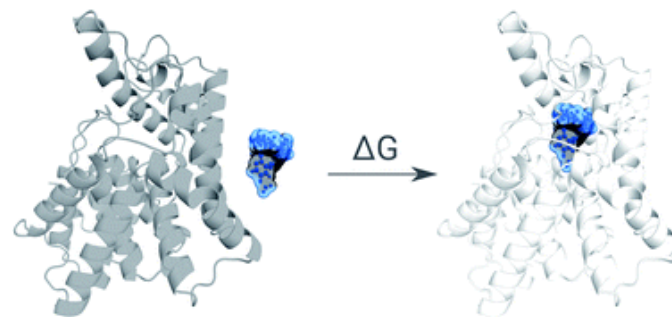
# Machine Learning for Computational Chemistry

## Reaction Mechanisms



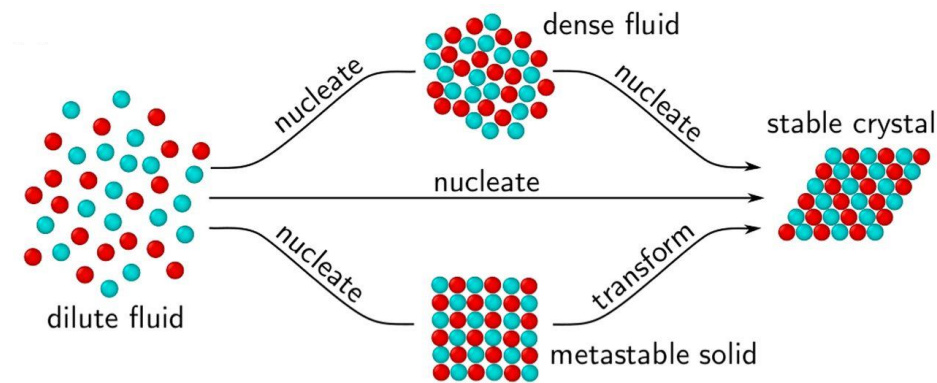
Ali, Y. *Sci. Rep.* **10**, 10995 (2020)

## Thermodynamic Binding Constants



Khalak, Y. *et al.*, *Chem. Sci.* **12**, 13958 (2021)

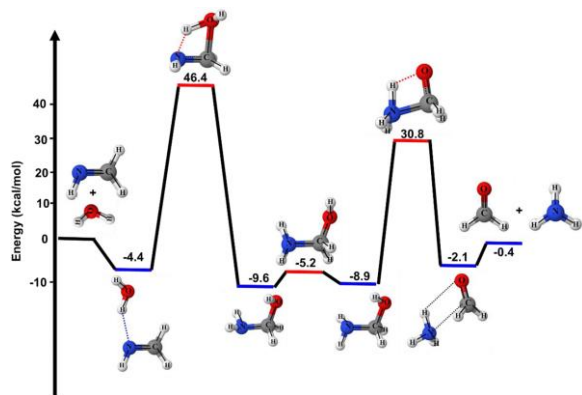
## Nucleation Events in Phase Transitions



Fang, H. *et al.*, *PNAS* **117**, 27927 (2020)

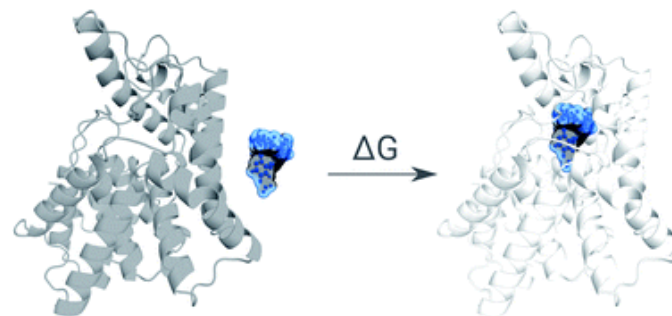
# Machine Learning for Computational Chemistry

## Reaction Mechanisms



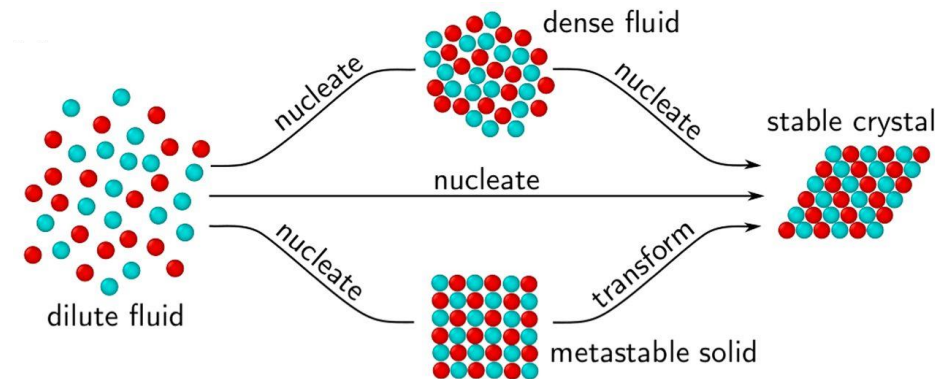
Ali, Y. *Sci. Rep.* **10**, 10995 (2020)

## Thermodynamic Binding Constants

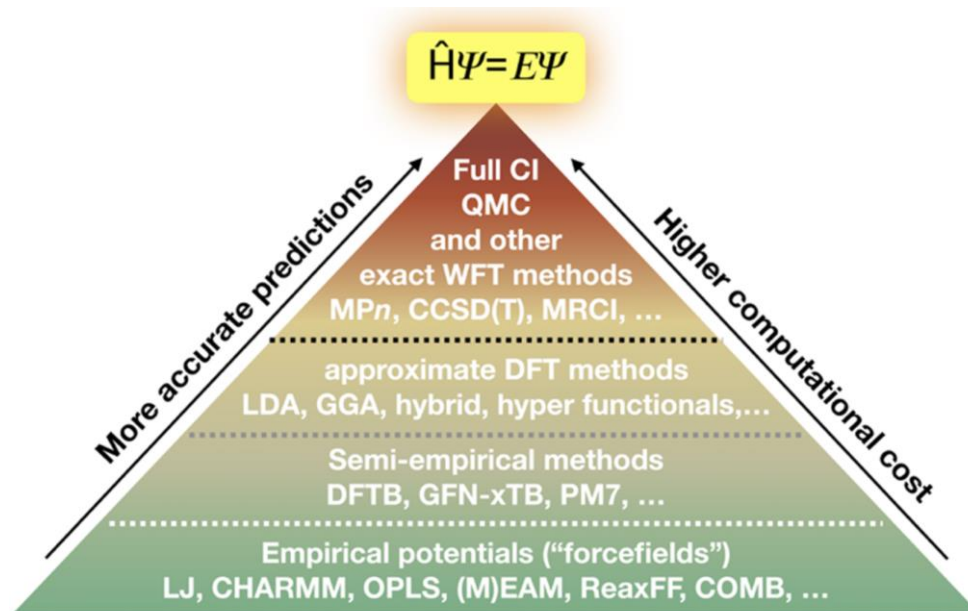


Khalak, Y. *et al.*, *Chem. Sci.* **12**, 13958 (2021)

## Nucleation Events in Phase Transitions

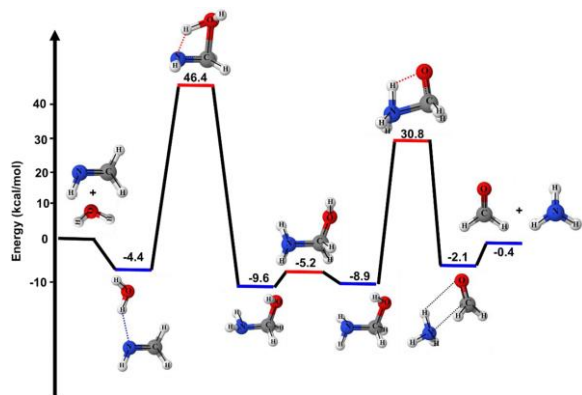


Fang, H. *et al.*, *PNAS* **117**, 27927 (2020)



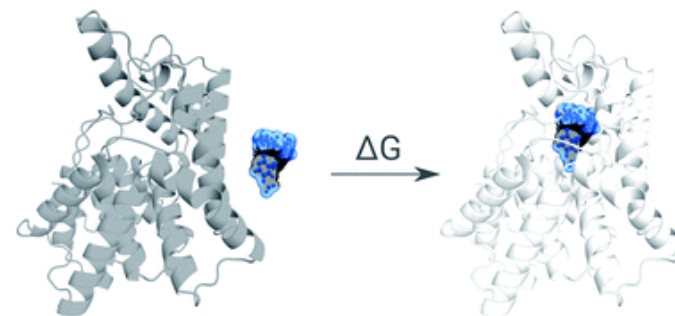
# Machine Learning for Computational Chemistry

## Reaction Mechanisms



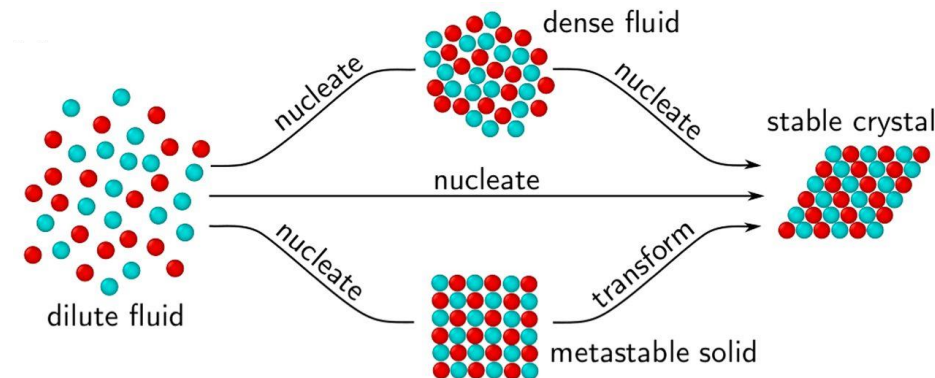
Ali, Y. *Sci. Rep.* **10**, 10995 (2020)

## Thermodynamic Binding Constants

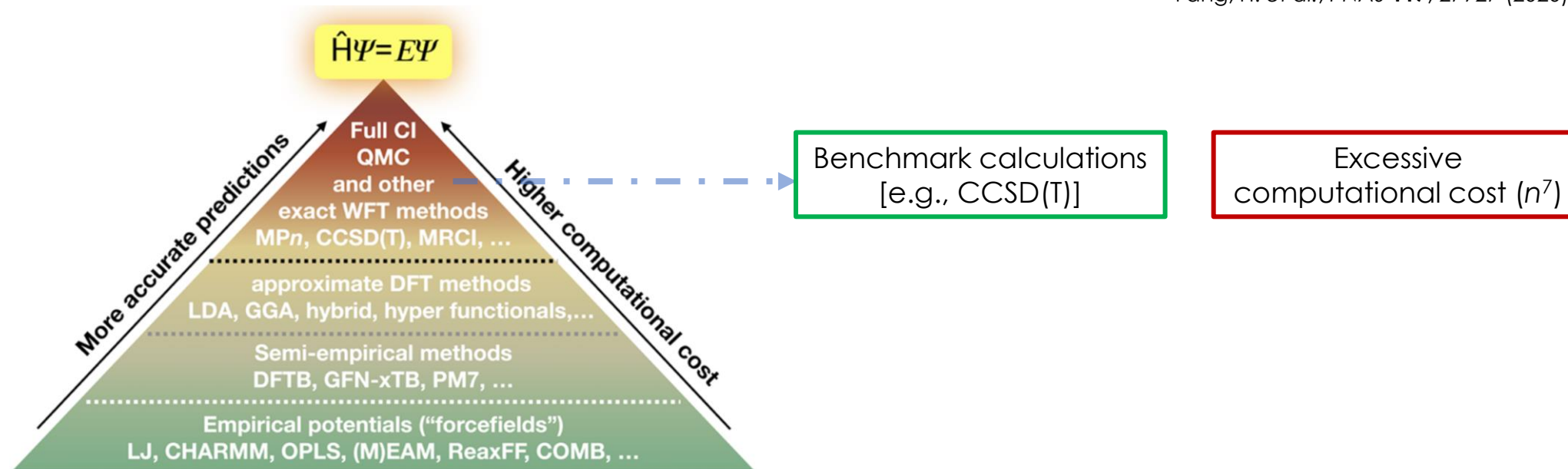


Khalak, Y. *et al.*, *Chem. Sci.* **12**, 13958 (2021)

## Nucleation Events in Phase Transitions

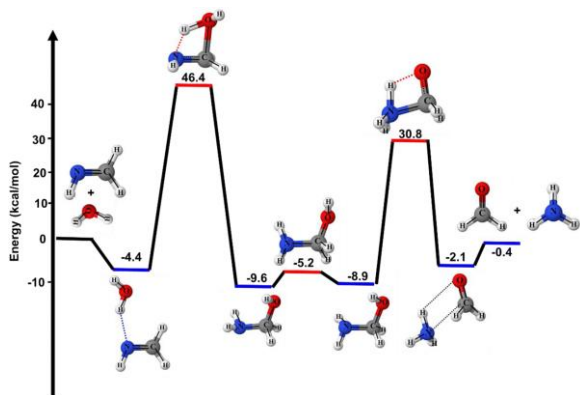


Fang, H. *et al.*, *PNAS* **117**, 27927 (2020)



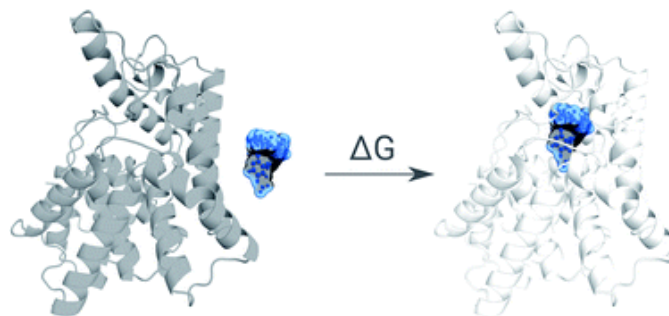
# Machine Learning for Computational Chemistry

## Reaction Mechanisms



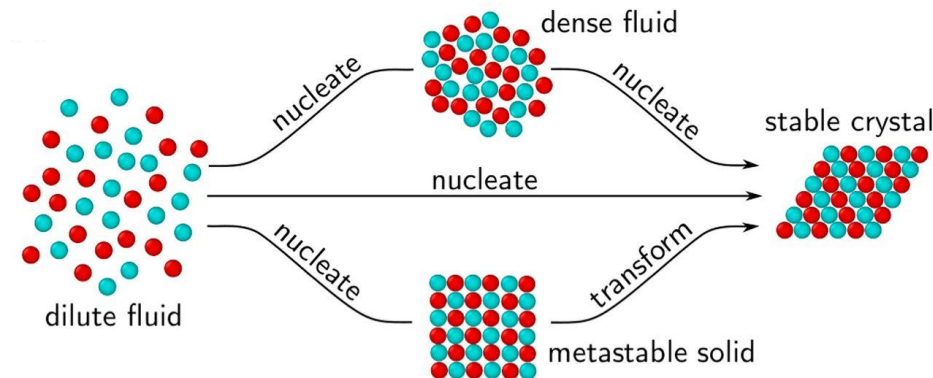
Ali, Y. *Sci. Rep.* **10**, 10995 (2020)

## Thermodynamic Binding Constants

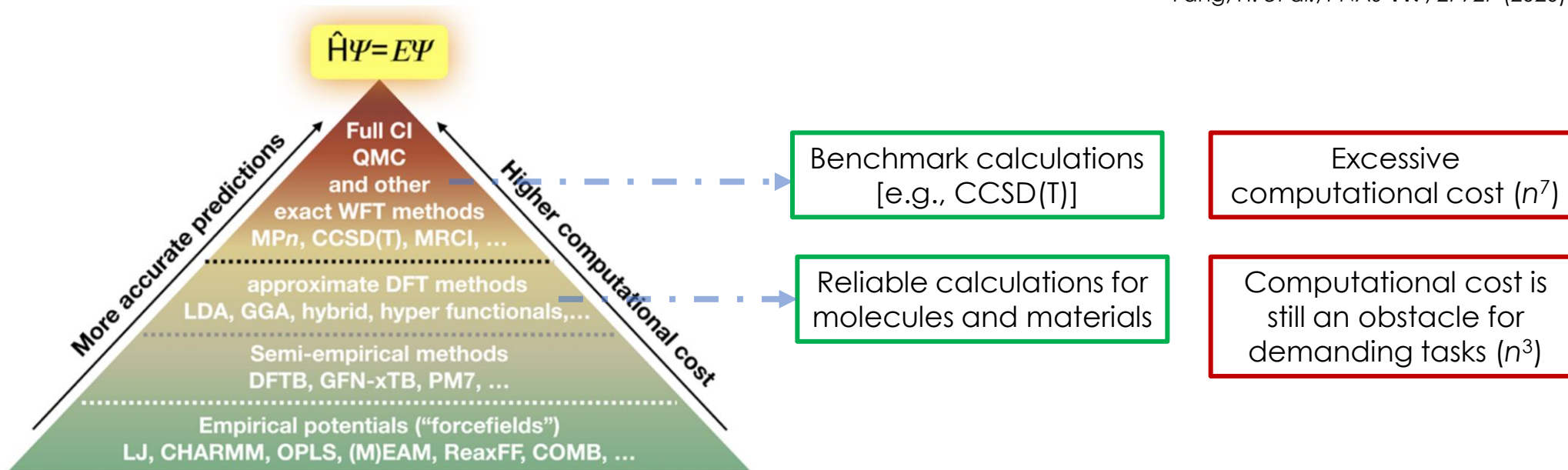


Khalak, Y. *et al.*, *Chem. Sci.* **12**, 13958 (2021)

## Nucleation Events in Phase Transitions

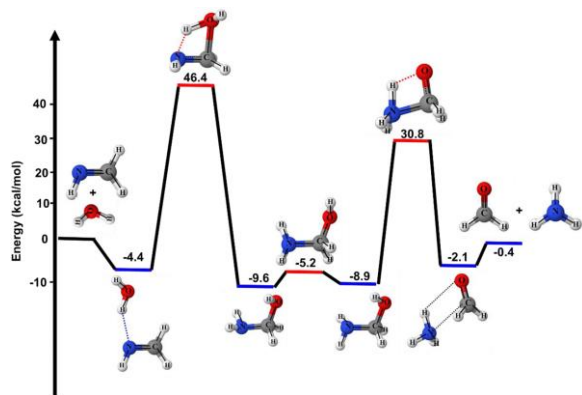


Fang, H. *et al.*, *PNAS* **117**, 27927 (2020)



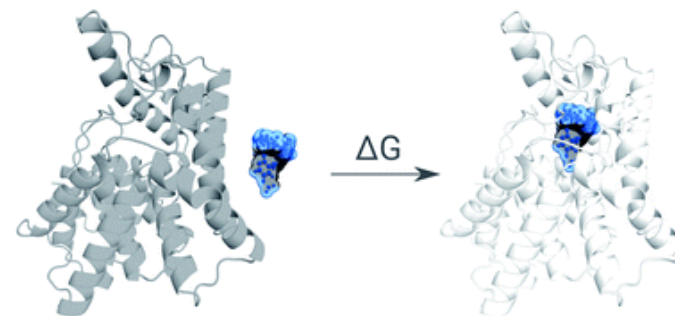
# Machine Learning for Computational Chemistry

## Reaction Mechanisms



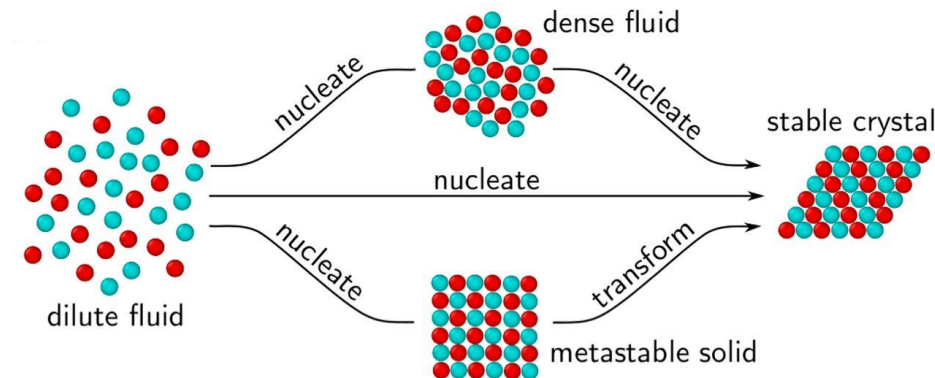
Ali, Y. *Sci. Rep.* **10**, 10995 (2020)

## Thermodynamic Binding Constants

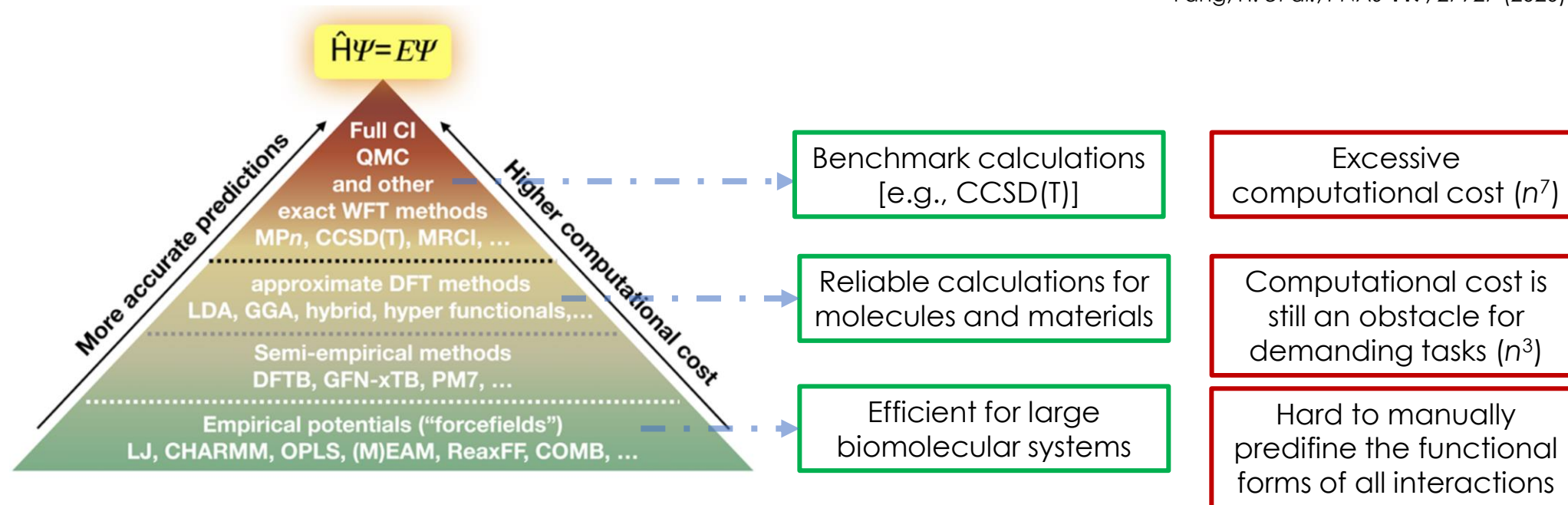


Khalak, Y. *et al.*, *Chem. Sci.* **12**, 13958 (2021)

## Nucleation Events in Phase Transitions

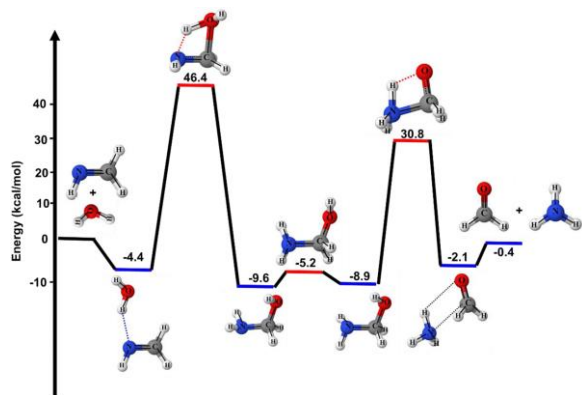


Fang, H. *et al.*, *PNAS* **117**, 27927 (2020)



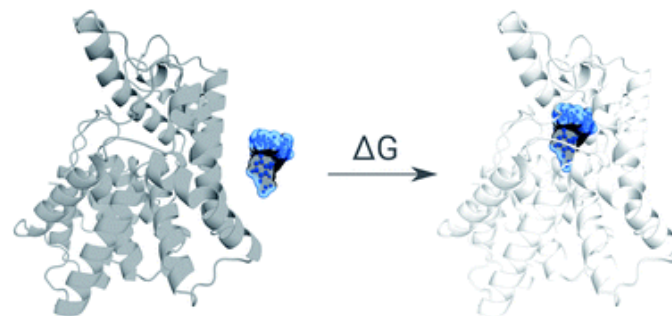
# Machine Learning for Computational Chemistry

## Reaction Mechanisms



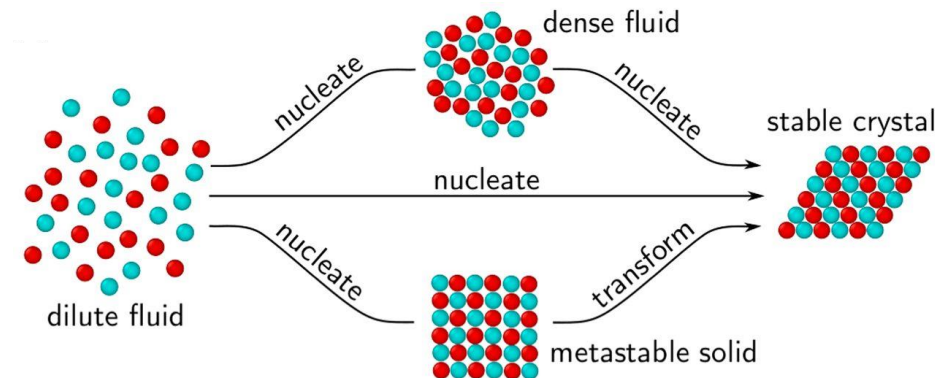
Ali, Y. *Sci. Rep.* **10**, 10995 (2020)

## Thermodynamic Binding Constants

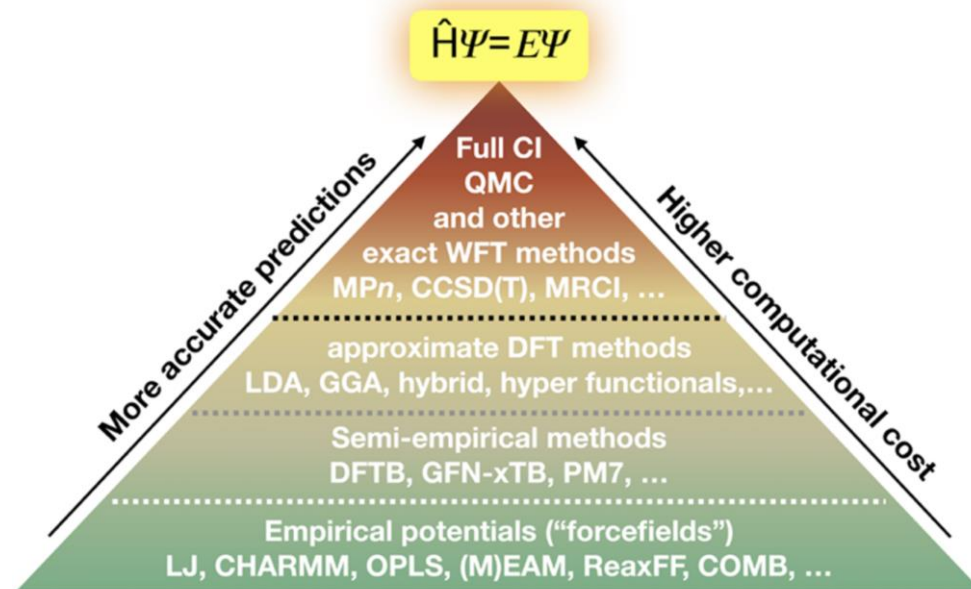


Khalak, Y. *et al.*, *Chem. Sci.* **12**, 13958 (2021)

## Nucleation Events in Phase Transitions



Fang, H. *et al.*, *PNAS* **117**, 27927 (2020)



## Machine Learning

$$\hat{f}: \mathcal{X} \xrightarrow{ML} \mathcal{Y}$$

$\mathcal{X}$  : Molecular configuration

$\mathcal{Y}$  : Property

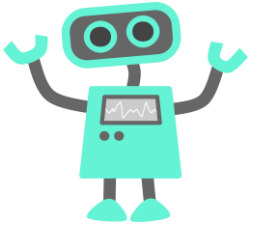
**Ab initio**  
accurate

**Machine**  
Learning

**Force Fields**  
efficient

# Materials Design and Discovery

## High-throughput screening

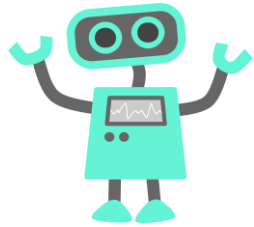


+

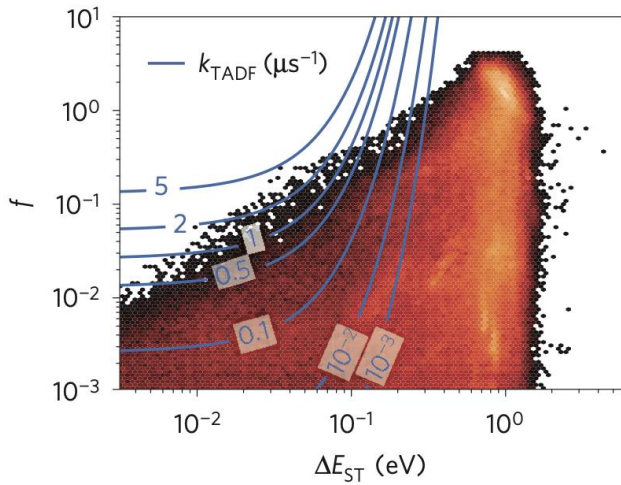
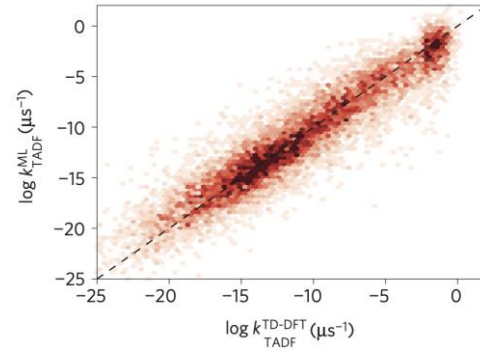


# Materials Design and Discovery

## High-throughput screening



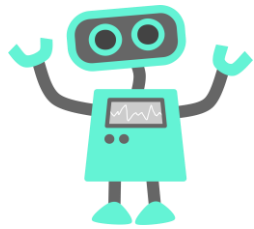
+



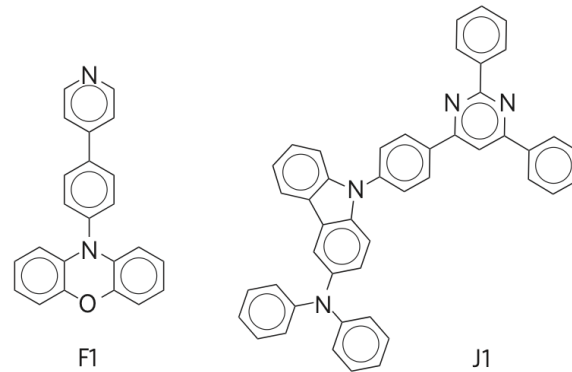
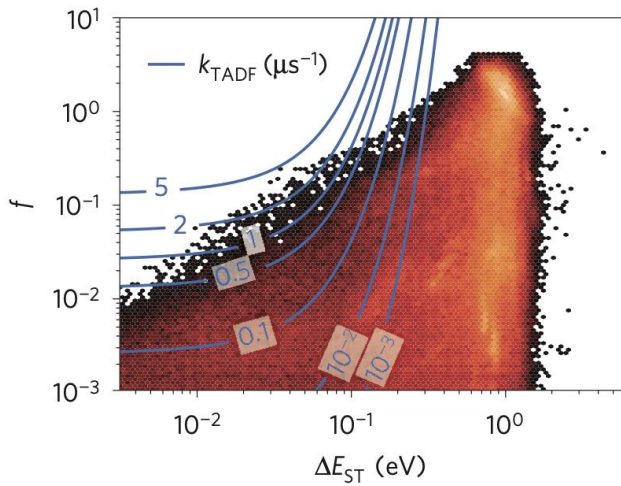
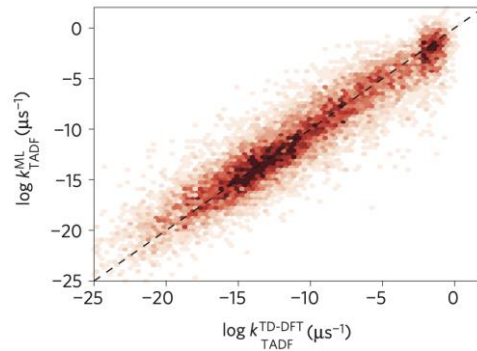
Gómez-Bombarelli, R., et al. *Nature Mater.* **15**, 1120–1127 (2016)

# Materials Design and Discovery

## High-throughput screening



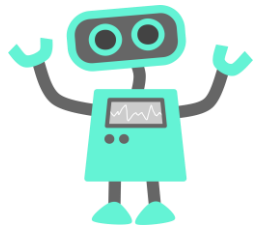
+



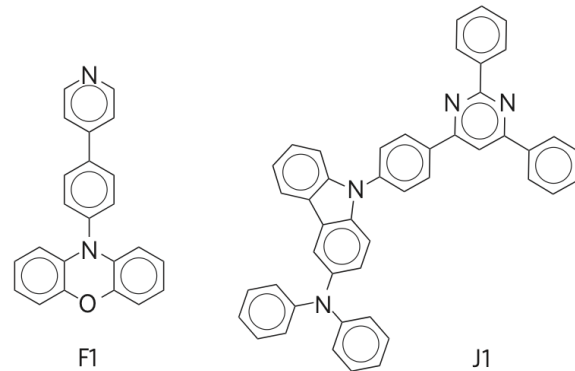
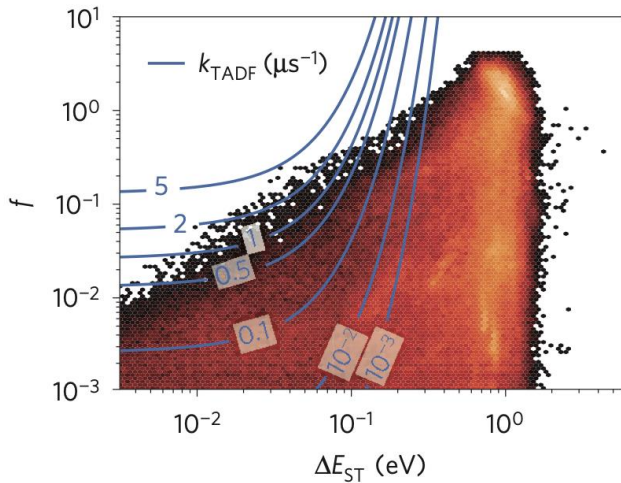
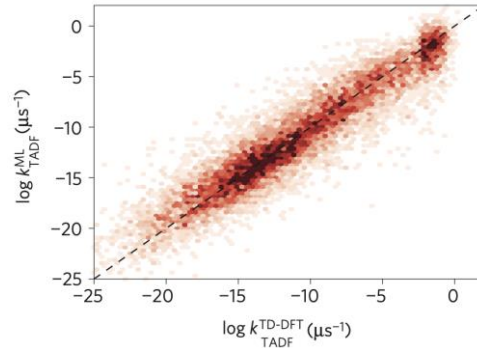
Gómez-Bombarelli, R., et al. *Nature Mater.* **15**, 1120–1127 (2016)

# Materials Design and Discovery

## High-throughput screening



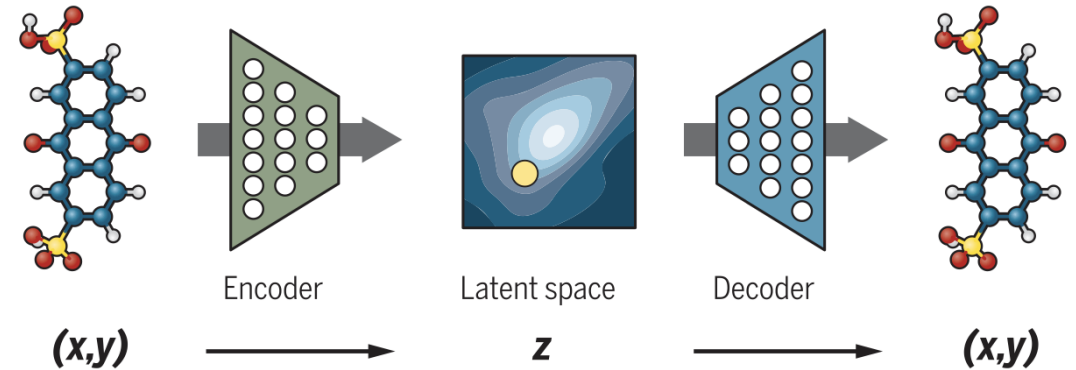
+



Gómez-Bombarelli, R., et al. *Nature Mater.* **15**, 1120–1127 (2016)

## Generative Models

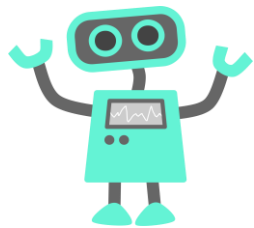
### VAE: Variational autoencoders



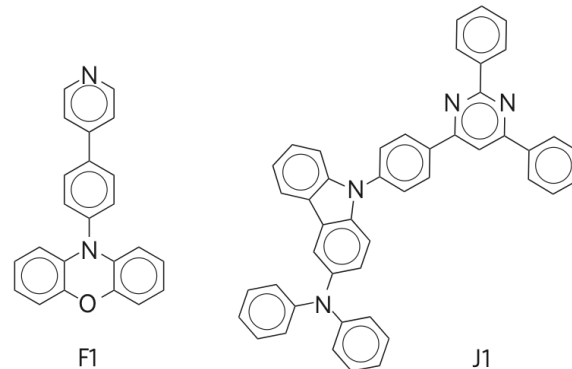
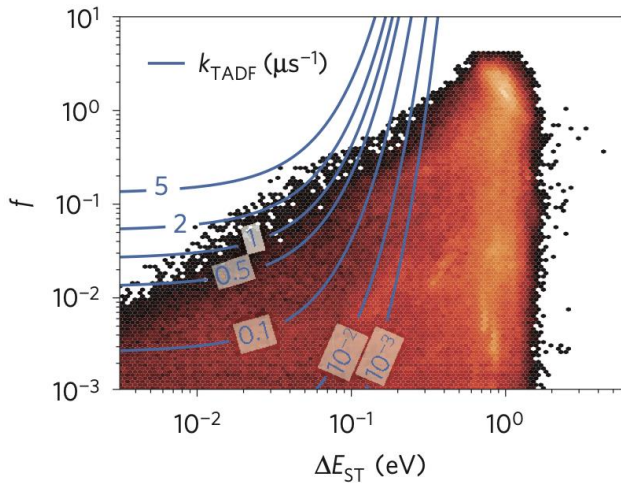
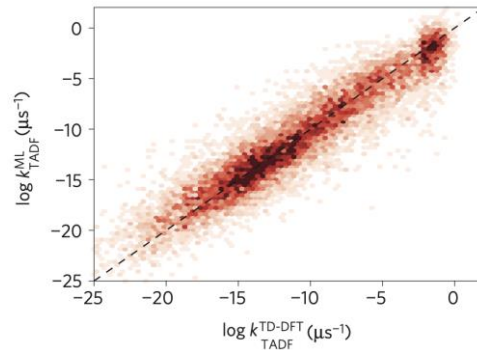
Sanchez-Lengeling, B. & Aspuru-Guzik, A. *Science* **361**, 360–365 (2018).

# Materials Design and Discovery

## High-throughput screening



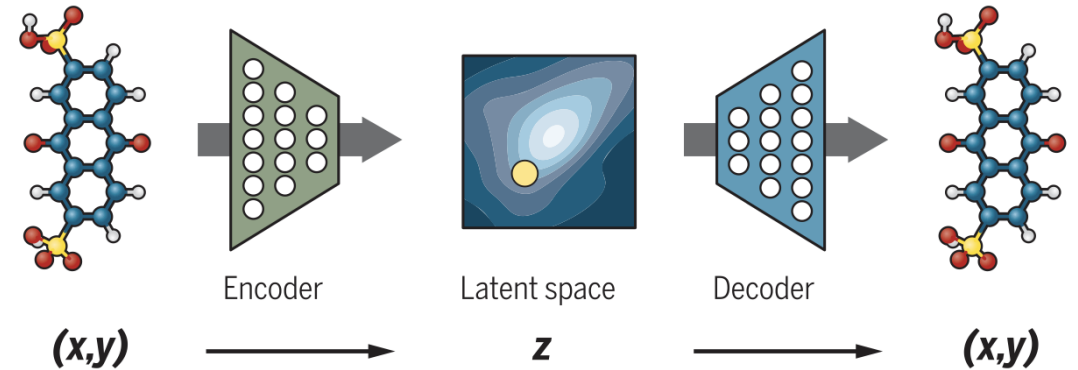
+



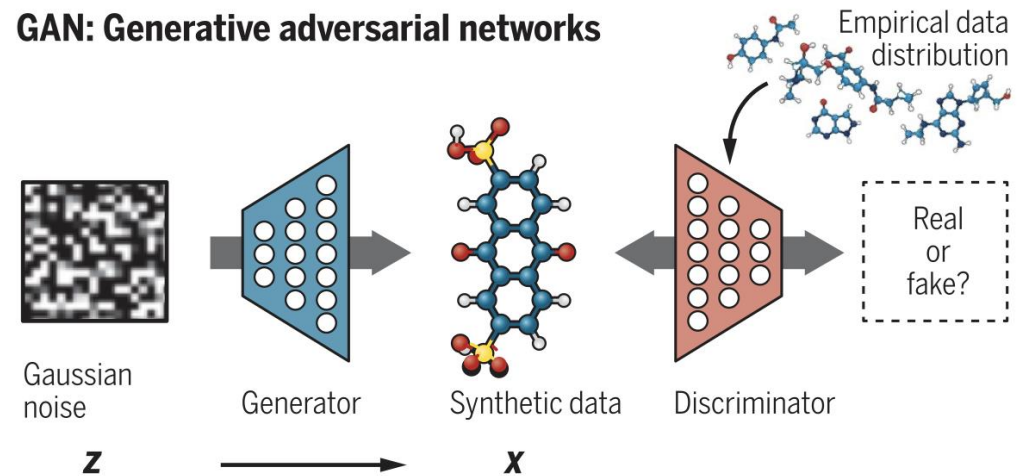
Gómez-Bombarelli, R., et al. *Nature Mater.* **15**, 1120–1127 (2016)

## Generative Models

### VAE: Variational autoencoders



### GAN: Generative adversarial networks



Sanchez-Lengeling, B. & Aspuru-Guzik, A. *Science* **361**, 360–365 (2018).

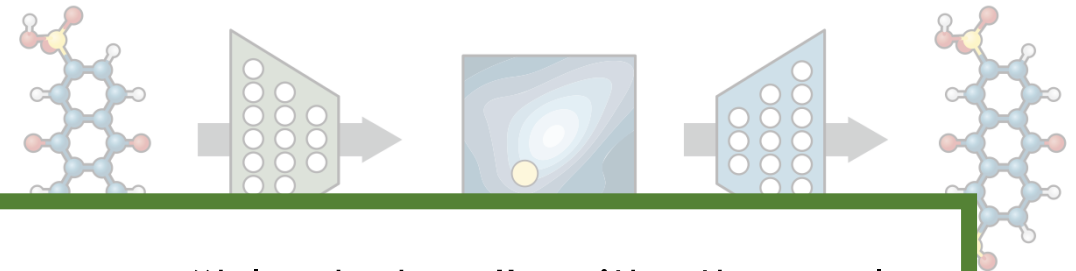
# Materials Design and Discovery

## High-throughput screening

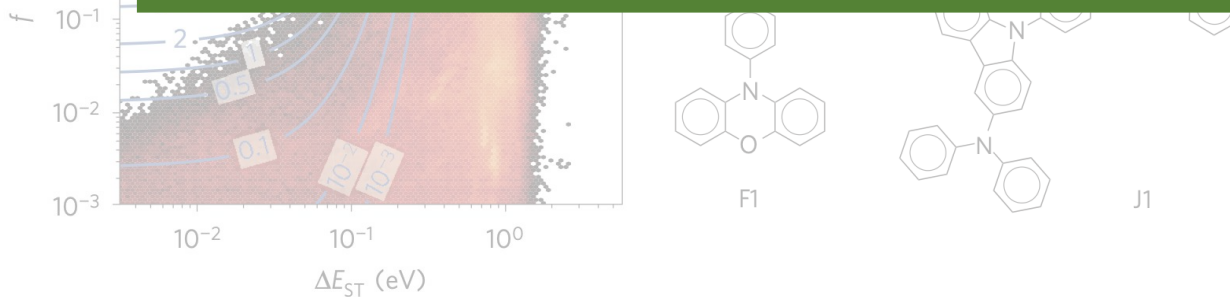


## Generative Models

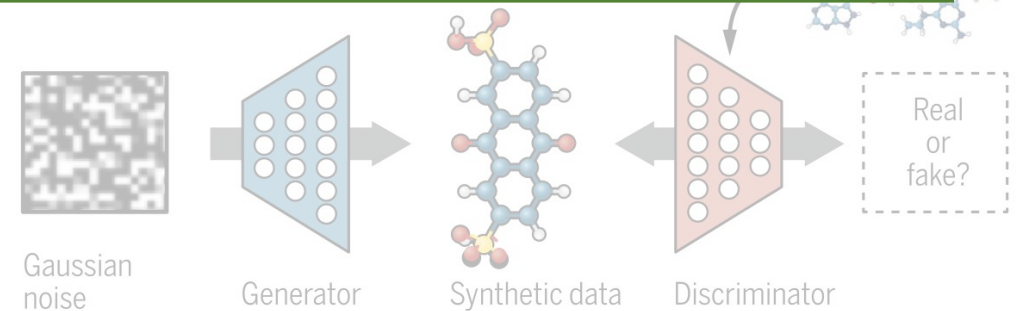
### VAE: Variational autoencoders



These approaches are very often used as a “black box” with the sole objective of obtaining a desired designed target with scarce or no special attention on why a given material is found to be better than others.



Gómez-Bombarelli, R., et al. *Nature Mater.* **15**, 1120–1127 (2016)



Sanchez-Lengeling, B. & Aspuru-Guzik, A. *Science* **361**, 360–365 (2018).

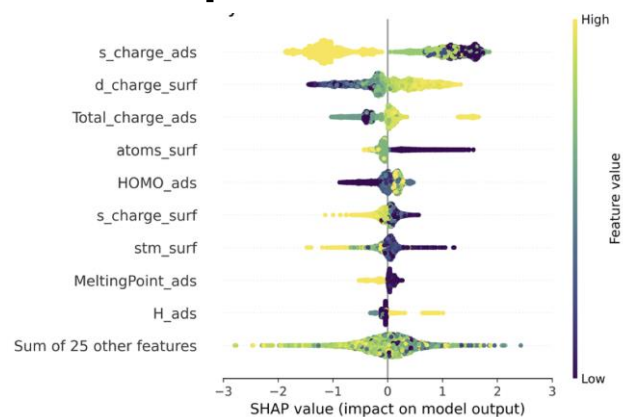
# Explainability

Assessment of prediction patterns to understand why a model makes specific decisions

# Explainability

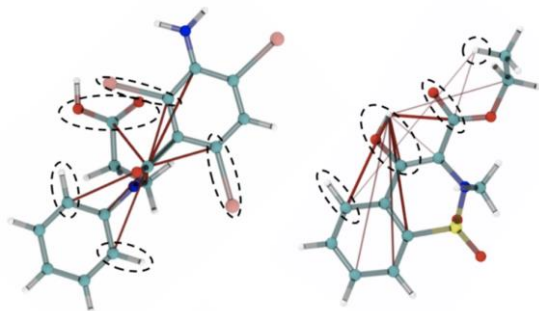
Assessment of prediction patterns to understand why a model makes specific decisions

## Feature importance attribution



Usuga A. F., et al. *J. Mater. Chem. A* **12**, 2708 (2024).

## Attention mechanisms

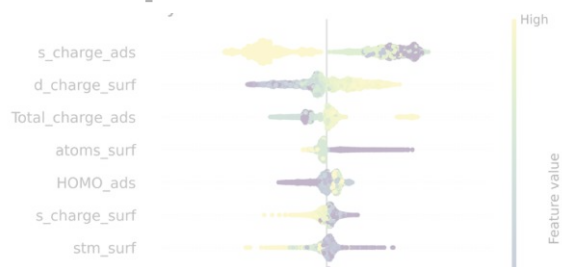


Cremer, J., et al. *Chem. Res. Toxicol.* **36**, 1561-1573 (2023).

# Explainability

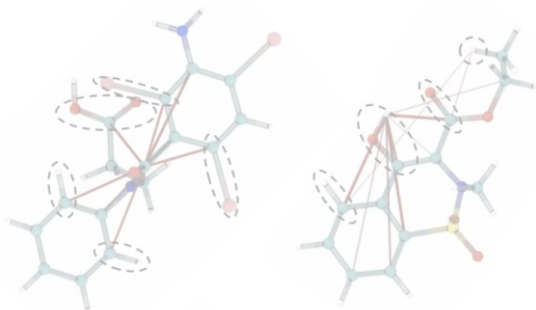
Assessment of prediction patterns to understand why a model makes specific decisions

## Feature importance attribution



These methods are not actionable (i.e., they do not tell how a given input can be changed in order to modify the output).

## Attention mechanisms

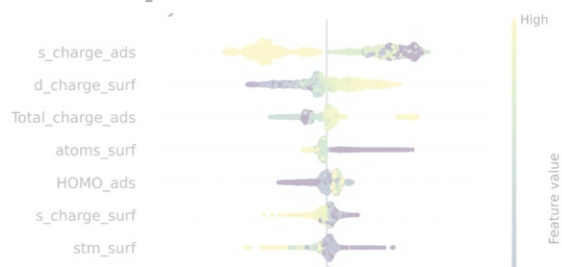


Cremer, J., et al. *Chem. Res. Toxicol.* 36, 1561-1573 (2023).

# Explainability

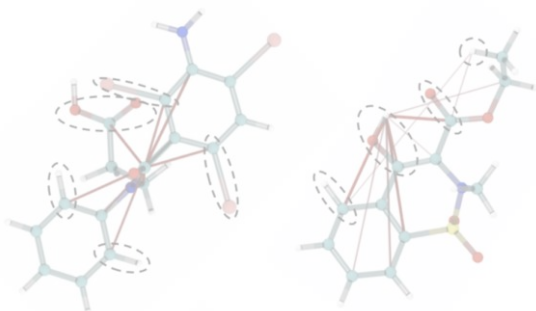
Assessment of prediction patterns to understand why a model makes specific decisions

## Feature importance attribution



These methods are not actionable (i.e., they do not tell how a given input can be changed in order to modify the output).

## Attention mechanisms



Cremer, J., et al. Chem. Res. Toxicol. 36, 1561-1573 (2023).

## Counterfactual explanations

Provide insights of model operation by determining examples or cases that explain the difference between a desired outcome and actual outcome

### Classification

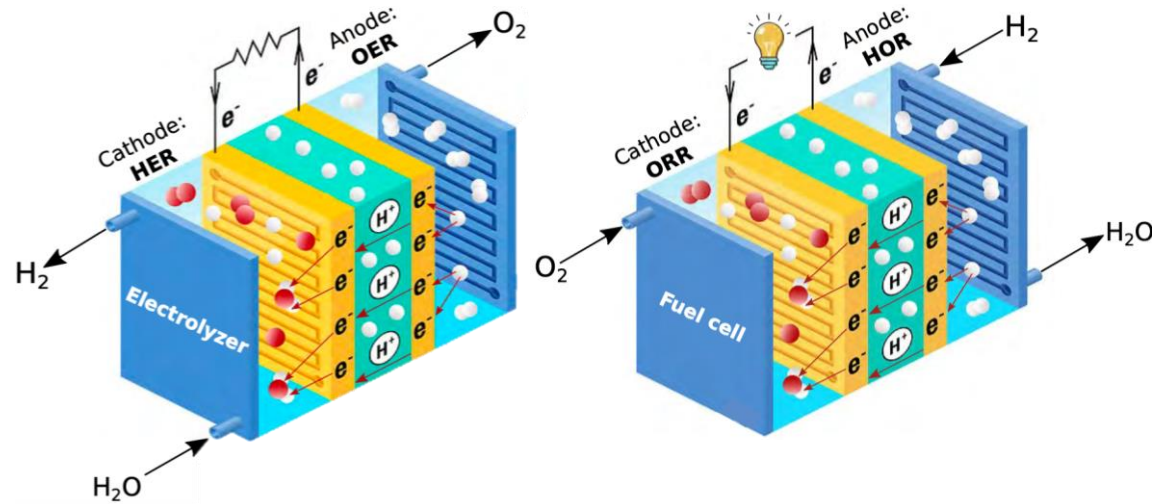
minimize  $d(\mathbf{x}, \mathbf{x}')$   
such that  $\hat{f}(\mathbf{x}) \neq \hat{f}(\mathbf{x}')$

### Regression

minimize  $d(\mathbf{x}, \mathbf{x}')$   
such that  $|\hat{f}(\mathbf{x}) - \hat{f}(\mathbf{x}')| \geq \Delta$

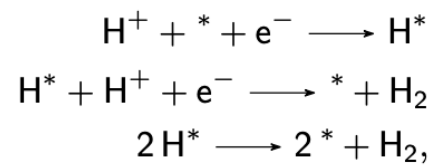
# Catalysts for the HER and ORR

## Hydrogen production and energy generation

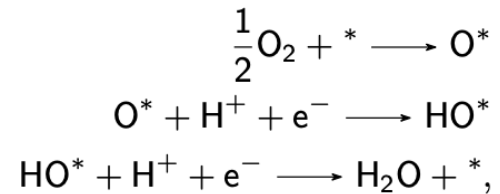


Martínez-Alonso, C., *High-throughput computational strategies to discover new catalysts for the hydrogen economy including elastic strain engineering* (2024).

### HER

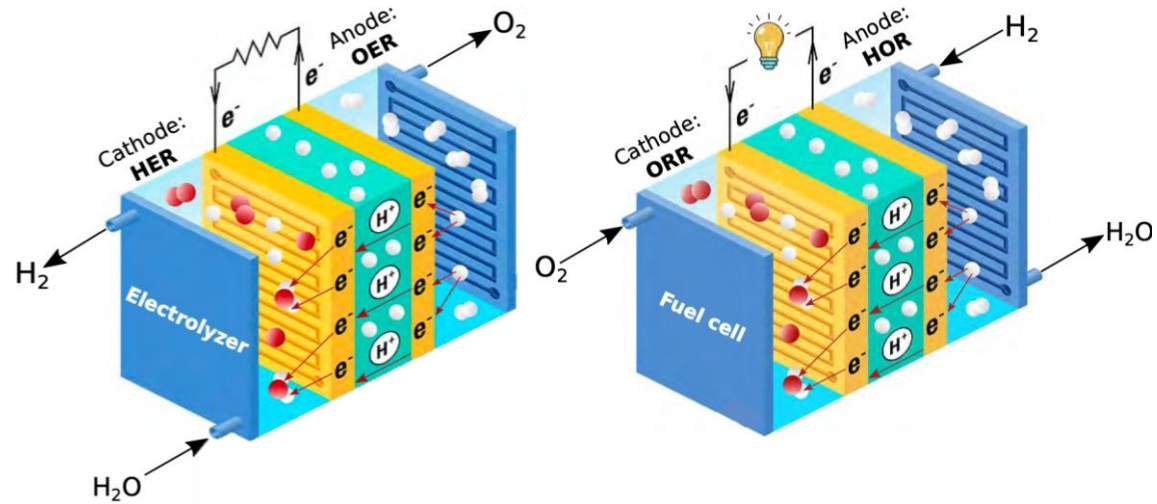


### ORR



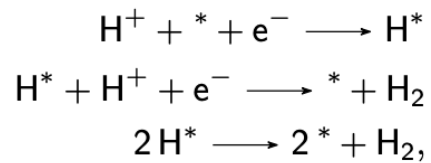
# Catalysts for the HER and ORR

## Hydrogen production and energy generation

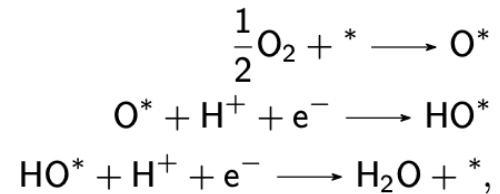


Martínez-Alonso, C., *High-throughput computational strategies to discover new catalysts for the hydrogen economy including elastic strain engineering* (2024).

### HER



### ORR

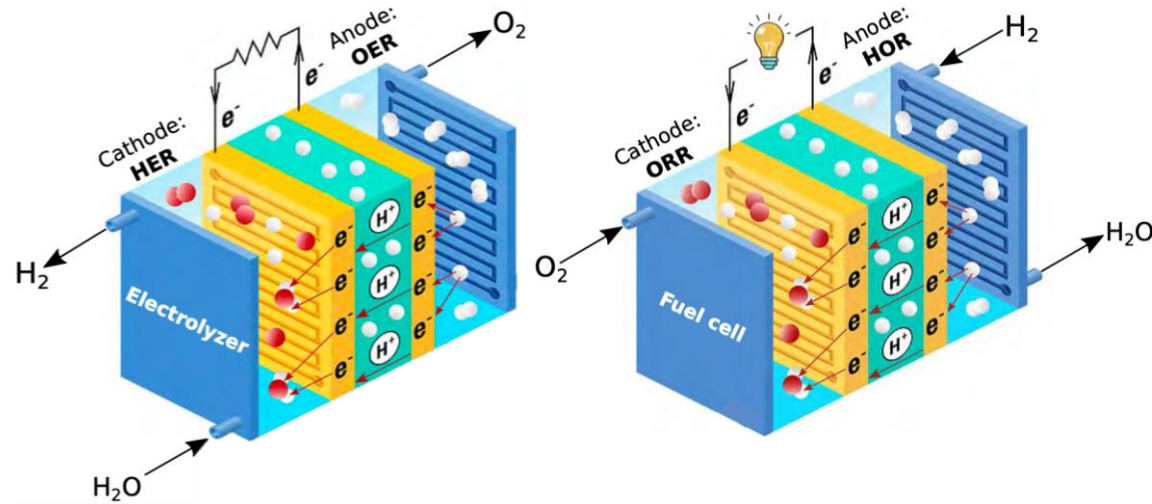


$E_{ads}$



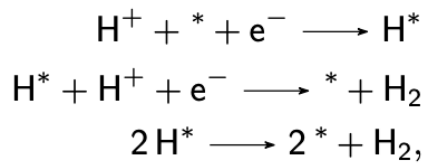
# Catalysts for the HER and ORR

## Hydrogen production and energy generation

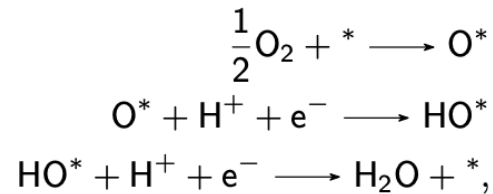


Martínez-Alonso, C., *High-throughput computational strategies to discover new catalysts for the hydrogen economy including elastic strain engineering* (2024).

### HER



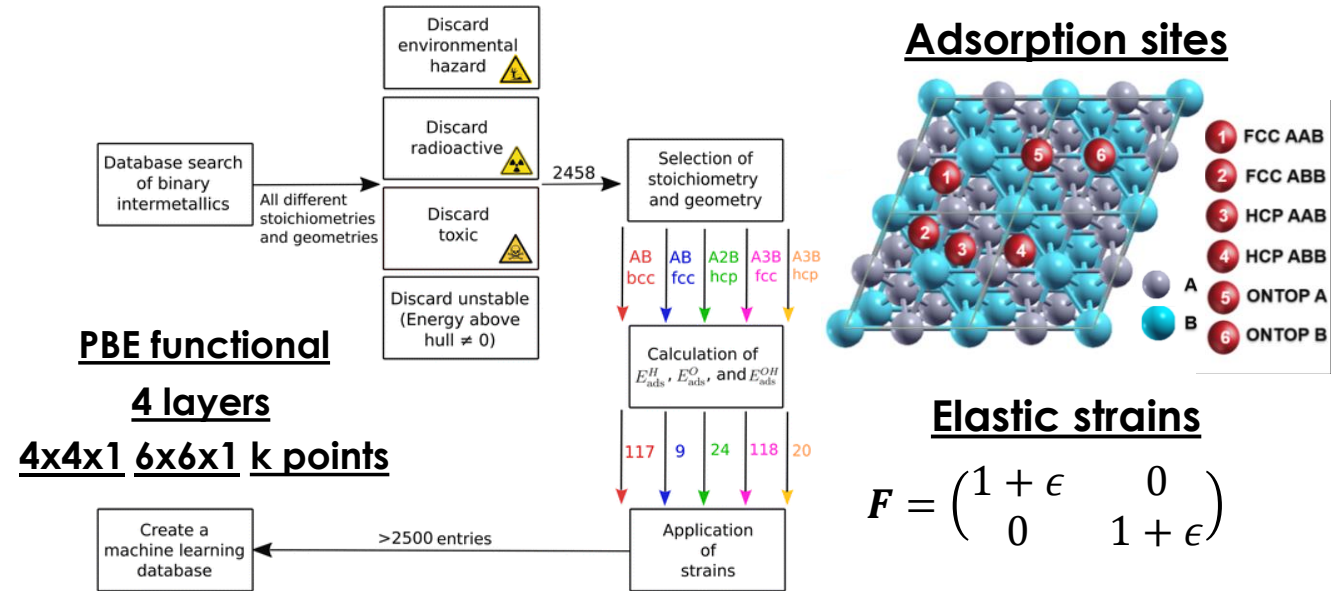
### ORR



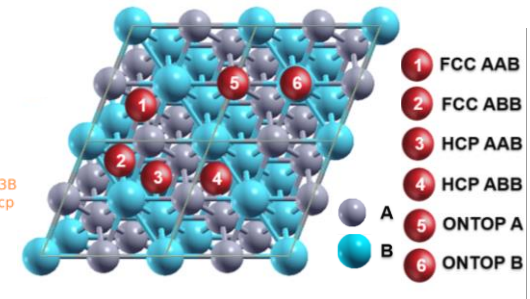
$E_{ads}$



## Dataset of $E_{ads}$



### Adsorption sites



### Elastic strains

$$F = \begin{pmatrix} 1 + \epsilon & 0 \\ 0 & 1 + \epsilon \end{pmatrix}$$

### Geometric

Unit cell volume  
WAR  
GCN

### Electronic

WEN  
WIE  
 $S_A$   
 $S_B$   
 $\Psi$

### Strain

Biaxial strain

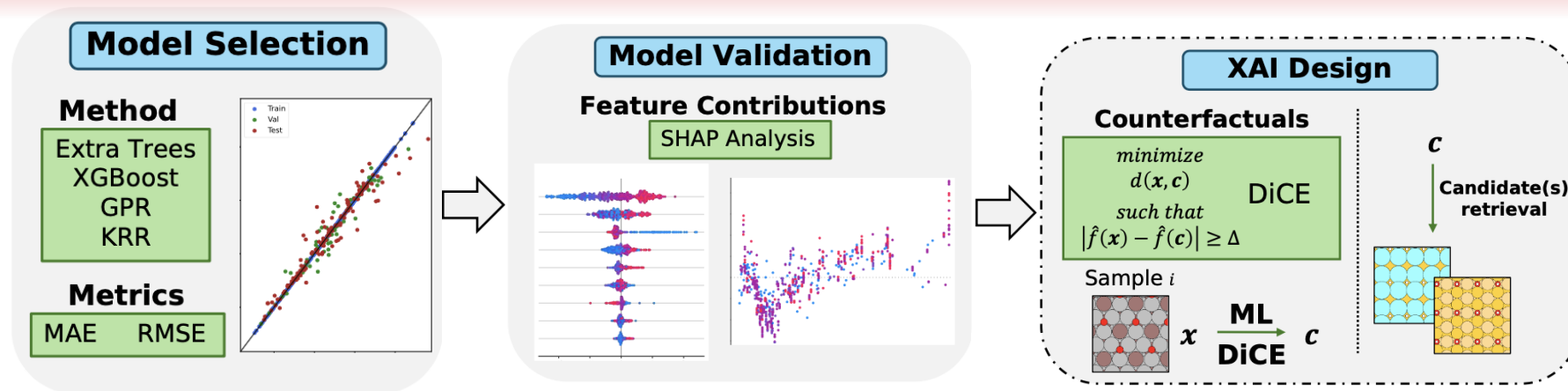
$$GCN = \frac{\sum_{i=1}^N CN_i}{CN_{max}}$$

*Angew. Chem.* **53**, 8316–8319 (2014)

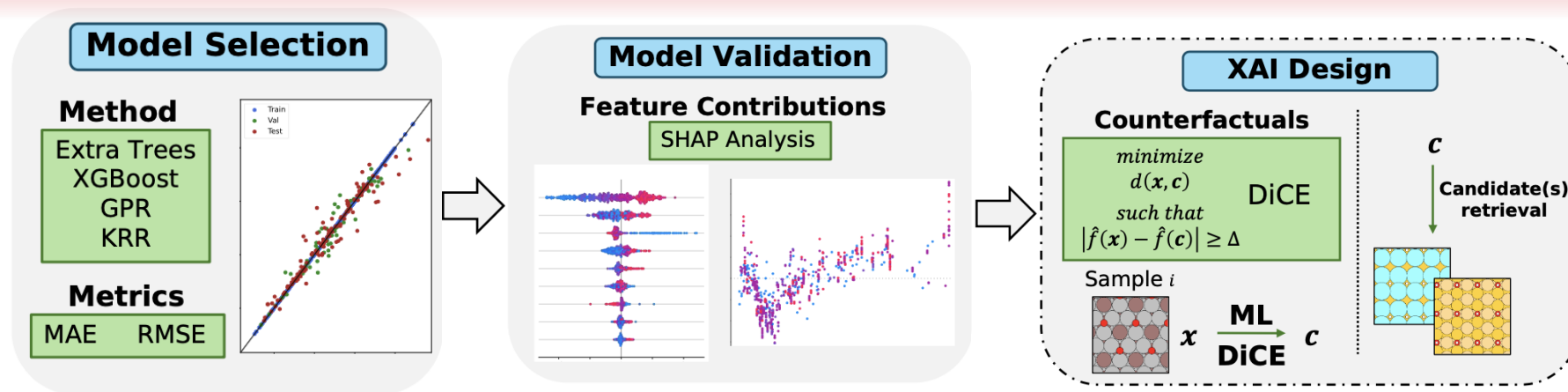
$$\Psi = \frac{\left(\prod_{i=1}^N S_i\right)^{\frac{2}{N}}}{\left(\prod_{i=1}^N EN_i\right)^{\frac{1}{N}}}$$

*Nat. Commun.* **11**, 1196 (2020)

# Explainable Artificial Intelligence (XAI) Strategy



# Explainable Artificial Intelligence (XAI) Strategy

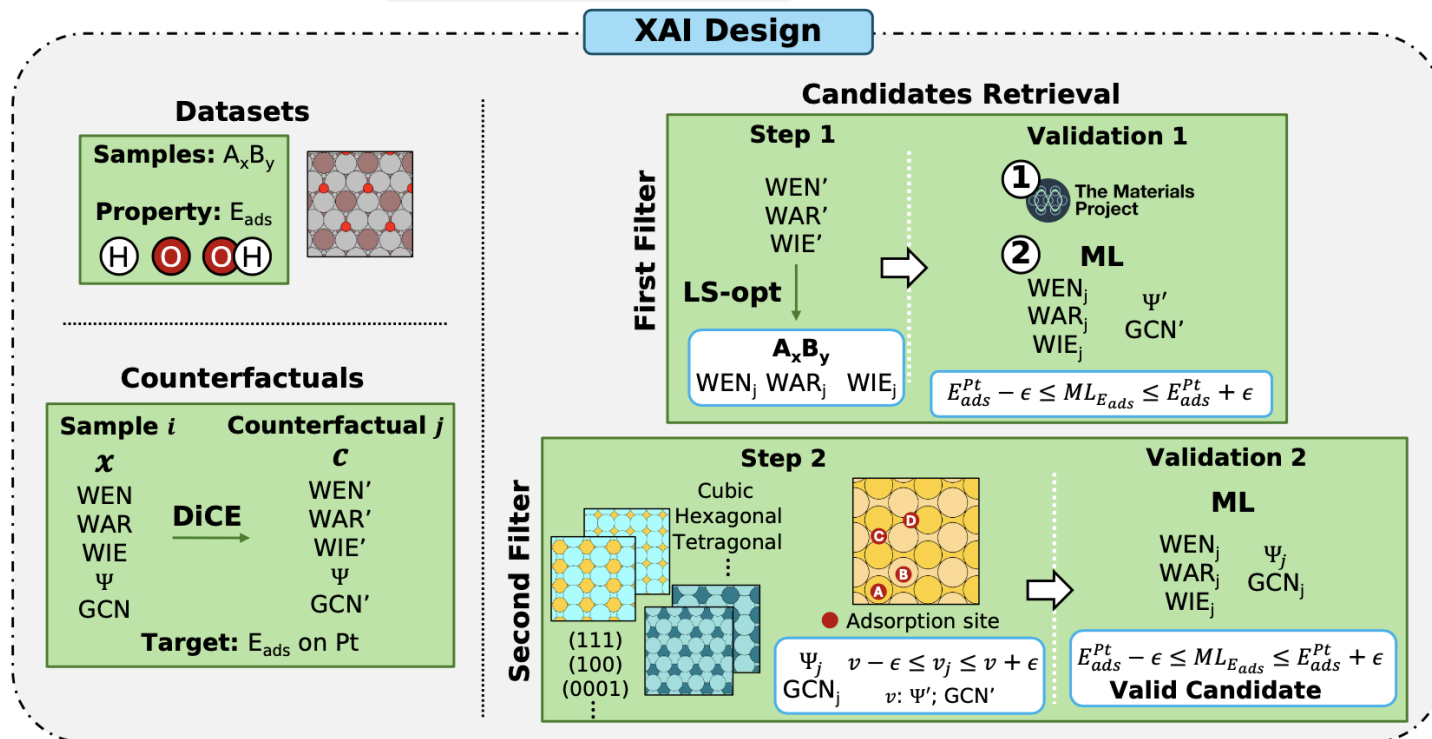
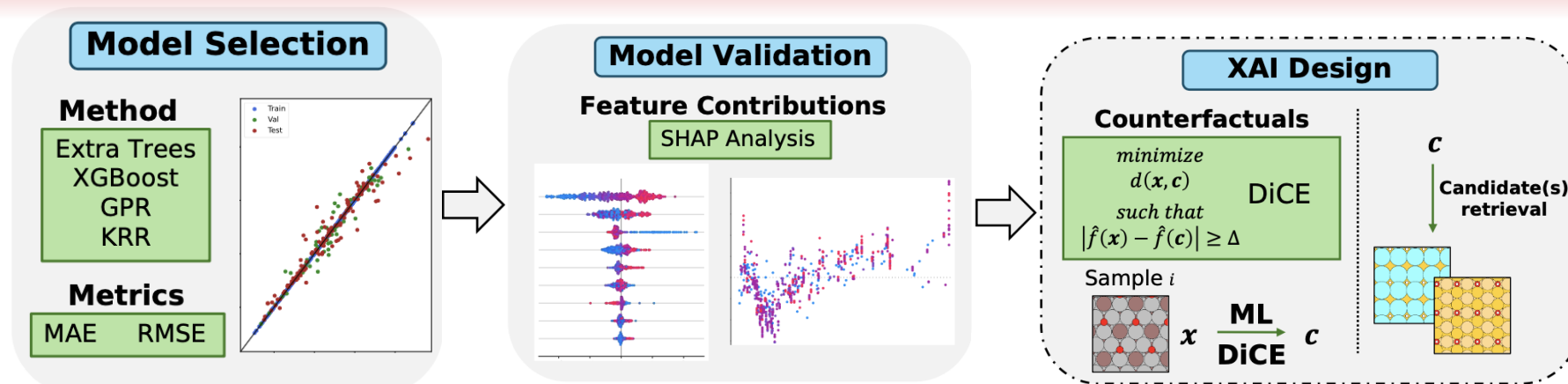


## DiCE Diverse Counterfactual Explanations

$$C(\mathbf{x}) = \arg \min_{\mathbf{c}_1, \dots, \mathbf{c}_k} \frac{1}{k} \sum_{i=1}^k \text{yloss}(f(\mathbf{c}_i), y) + \frac{\lambda_1}{k} \sum_{i=1}^k \text{dist}(\mathbf{c}_i, \mathbf{x}) - \lambda_2 \text{dpp-diversity}(\mathbf{c}_1, \dots, \mathbf{c}_k)$$

Mothilal, R. K., et al. in *Proceedings of the 2020 conference on fairness, accountability, and transparency*, 607–617 (2020)

# Explainable Artificial Intelligence (XAI) Strategy



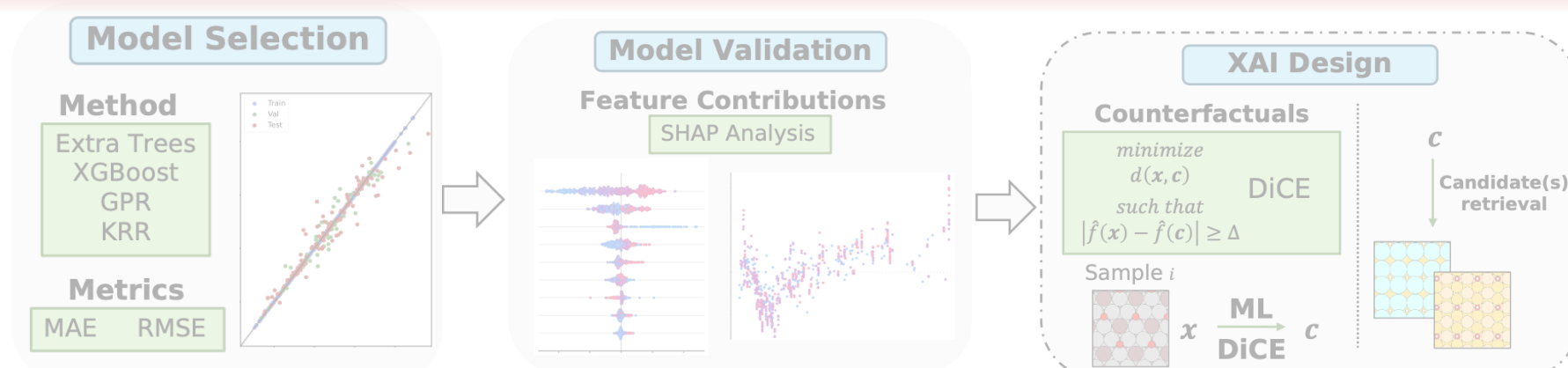
## DiCE

### Diverse Counterfactual Explanations

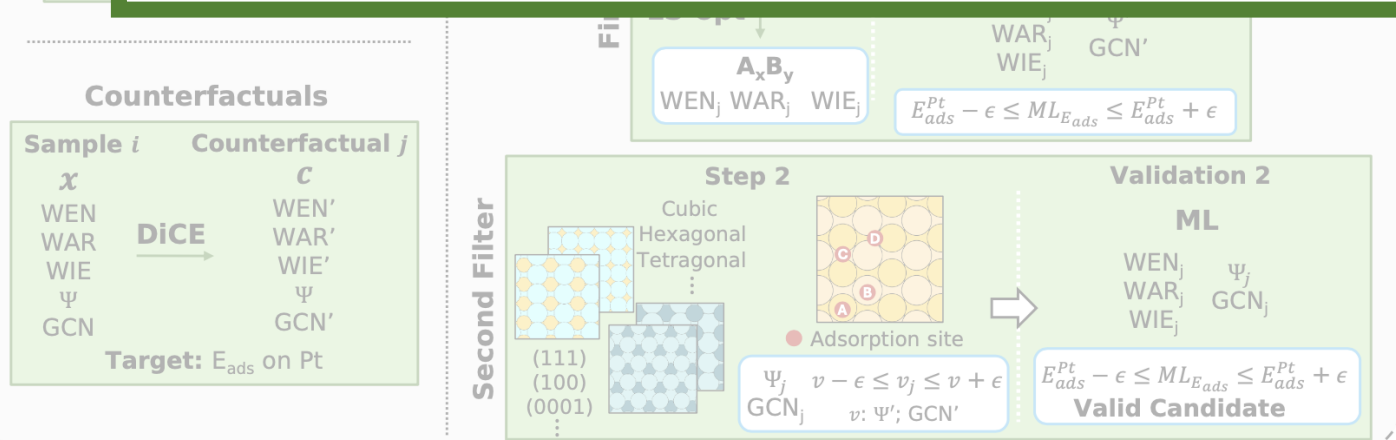
$$C(\mathbf{x}) = \arg \min_{\mathbf{c}_1, \dots, \mathbf{c}_k} \frac{1}{k} \sum_{i=1}^k y_{loss}(f(\mathbf{c}_i), y) + \frac{\lambda_1}{k} \sum_{i=1}^k dist(\mathbf{c}_i, \mathbf{x}) - \lambda_2 dpp\_diversity(\mathbf{c}_1, \dots, \mathbf{c}_k)$$

Mothilal, R. K., et al. in Proceedings of the 2020 conference on fairness, accountability, and transparency, 607–617 (2020)

# Explainable Artificial Intelligence (XAI) Strategy



Explainability is ensured by construction since the discovered materials can be linked directly to the original sample of the dataset from which the counterfactual was generated.



$$C(x) = \arg \min_{c_1, \dots, c_k} \frac{1}{k} \sum_{i=1}^k y_{loss}(f(c_i), y) + \frac{\lambda_1}{k} \sum_{i=1}^k dist(c_i, x) - \lambda_2 dpp\_diversity(c_1, \dots, c_k)$$

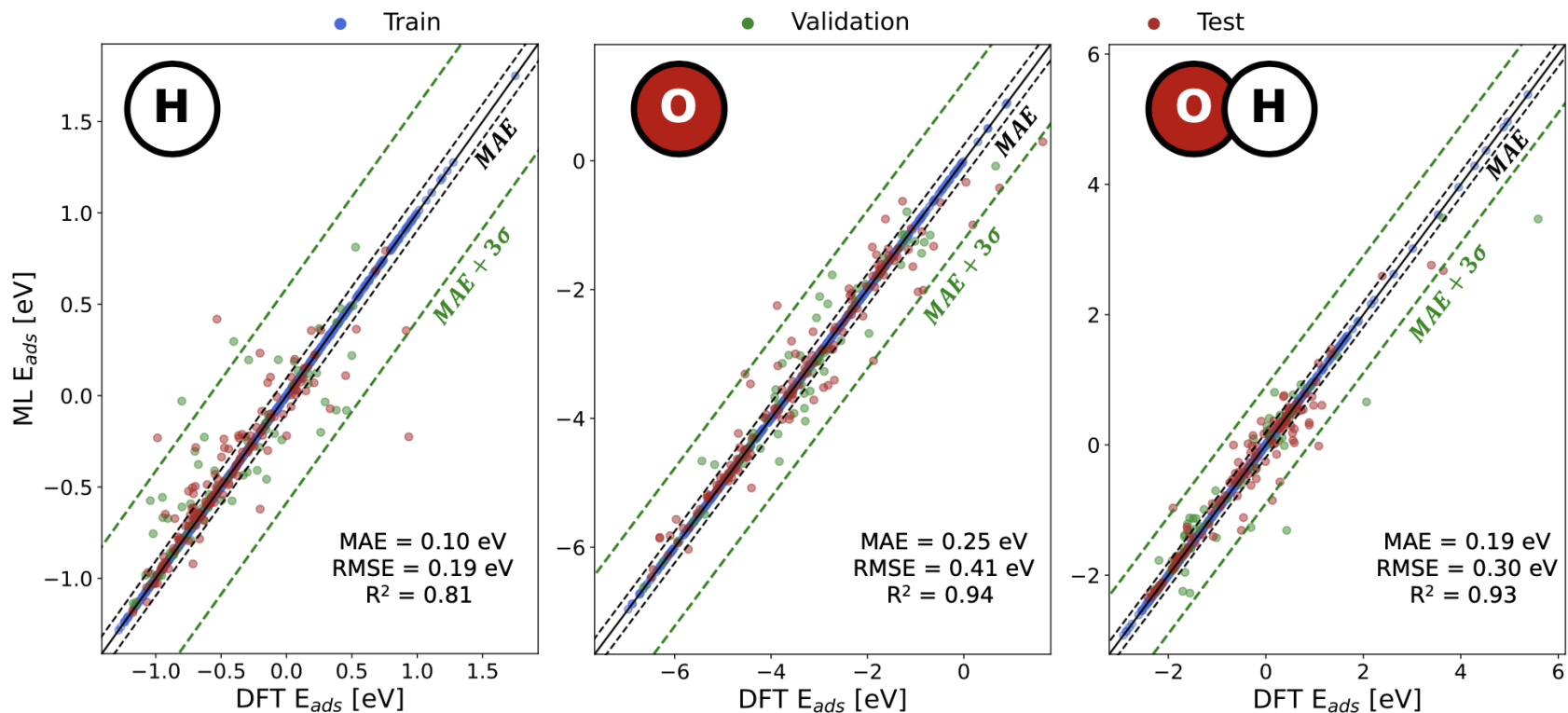
Mothilal, R. K., et al. in Proceedings of the 2020 conference on fairness, accountability, and transparency, 607–617 (2020)

# Model Selection

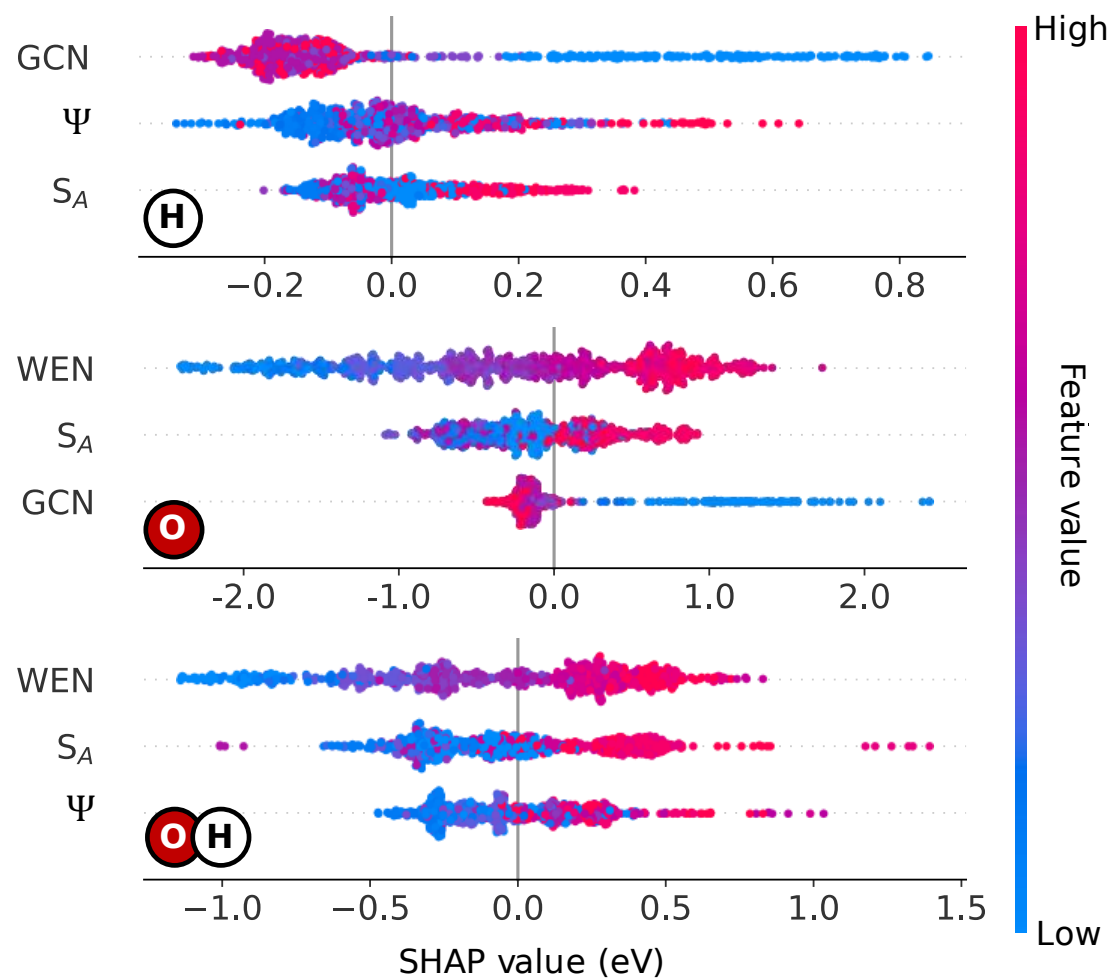
## Training settings

- 15% samples for testing
- 85% samples for training/validation
- 10-fold “pseudo-random” cross-validation

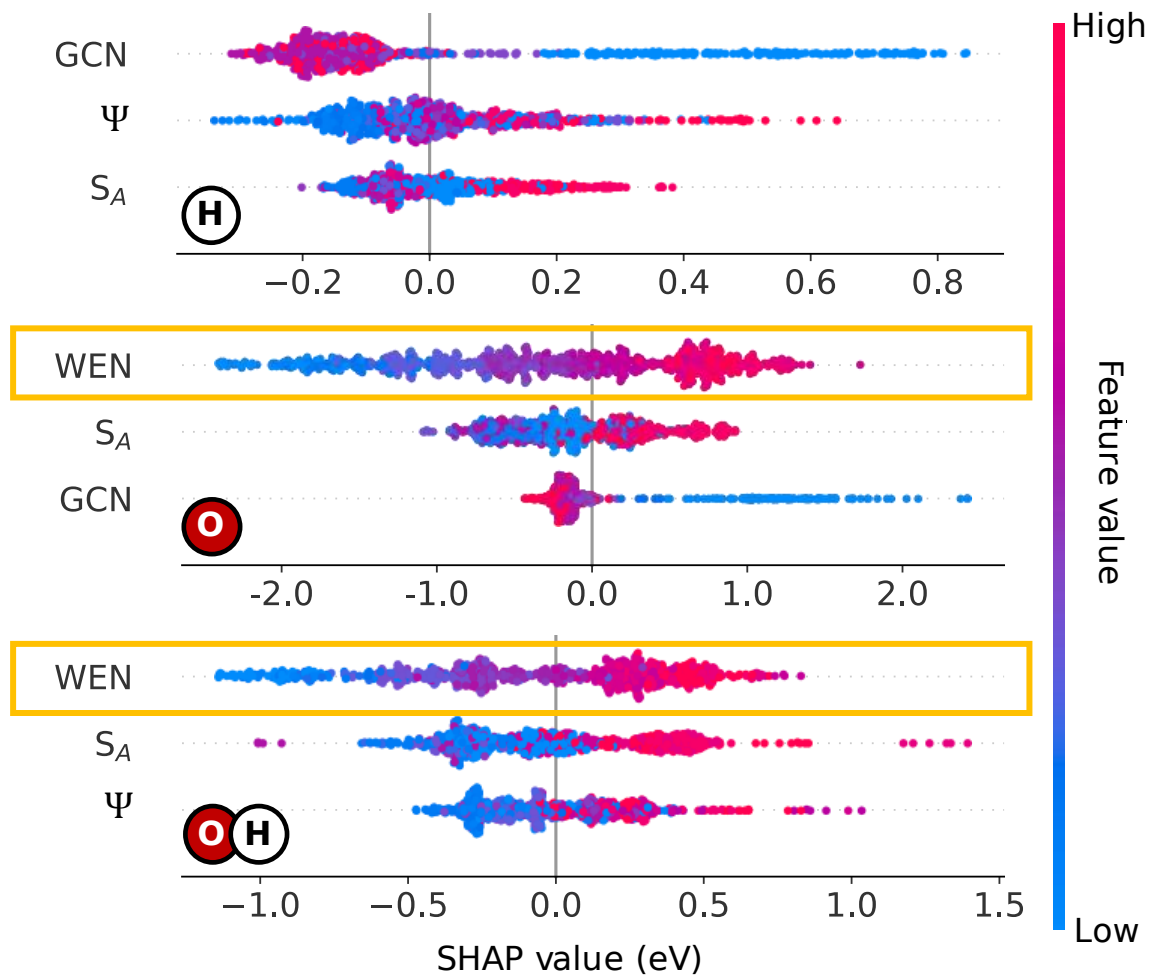
Adsorbate	Metric	ET	XGBoost	KRR	GPR
H	MAE	<b>0.10</b>	0.11	0.13	0.15
	RMSE	<b>0.21</b>	<b>0.21</b>	0.24	0.25
O	MAE	<b>0.25</b>	0.26	0.34	0.34
	RMSE	<b>0.43</b>	0.47	0.54	0.55
OH	MAE	<b>0.20</b>	0.28	0.27	0.41
	RMSE	<b>0.33</b>	0.45	0.42	0.60



# Model Validation

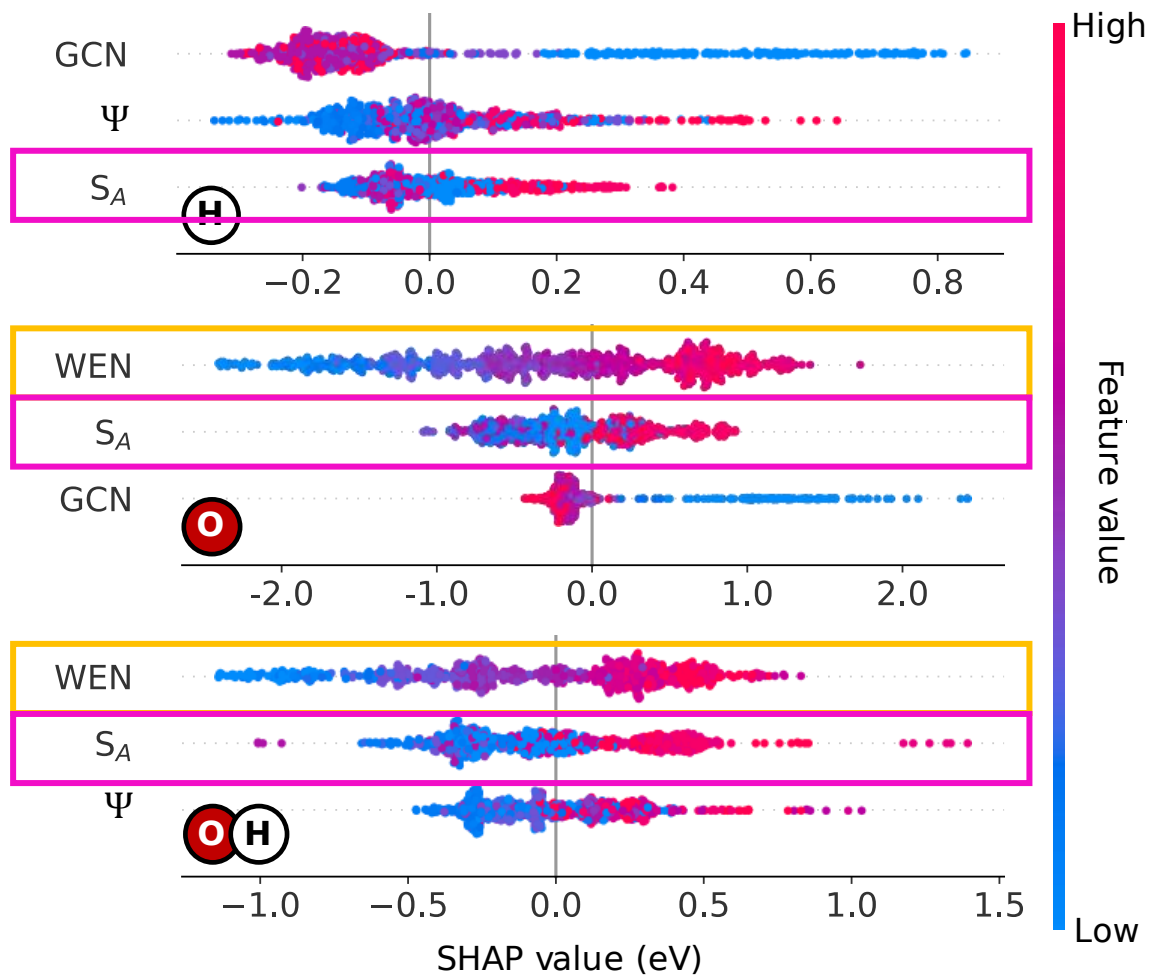


# Model Validation



O is a very electronegative atom, it will tend to bind stronger to surfaces with low electronegativity.

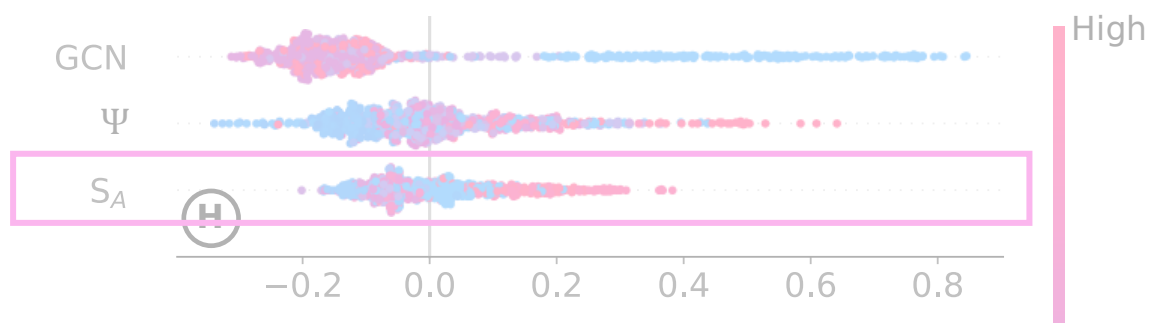
# Model Validation



O is a very electronegative atom, it will tend to bind stronger to surfaces with low electronegativity.

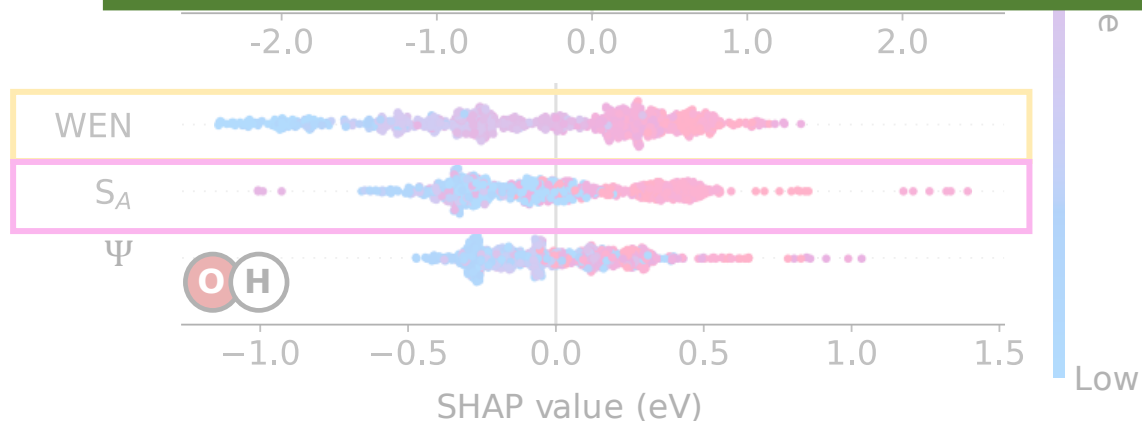
A large number of valence electrons downshifts the d-band center. The more negative the d-band centers the less favorable the adsorption.

# Model Validation



O is a very electronegative atom, it will tend to bind stronger to surfaces with low electronegativity.

Our ML models provide accurate and reliable predictions by leveraging rigorous chemical and physical concepts.



A large number of valence electrons downshifts the d-band center. The more negative the d-band centers the less favorable the adsorption.

# Discovered Materials

Target  $E_{ads}$       **(H)** -0.49 eV      **(O)** -1.79 eV      **(OH)** -1.19 eV

Adsorbate	Material	Crystal	Facet	Site	ML $E_{ads}$ (eV)
<b>(H)</b>	CaNi	Cubic	101	bridge-CaNi	-0.46
	YIr	Cubic	101	ontop-Ir	-0.45
	<b>YAu</b>	<b>Cubic</b>	<b>100</b>	<b>hollow</b>	<b>-0.47</b>
	YAu	Cubic	101	hollow	-0.53
<b>(O)</b>	<b>ZnIr</b>	<b>Hexagonal</b>	<b>10-11</b>	<b>hollow</b>	<b>-1.72</b>
	<b>ZnPt</b>	<b>Tetragonal</b>	<b>101</b>	<b>hollow</b>	<b>-1.76</b>
<b>(OH)</b>	<b>CdAu</b>	<b>Cubic</b>	<b>101</b>	<b>longbridge-Au</b>	<b>1.10</b>
	CdAu	Hexagonal	0001	bridge-Au	1.10

# Discovered Materials

Target  $E_{ads}$       **(H)** -0.49 eV      **(O)** -1.79 eV      **(OH)** -1.19 eV

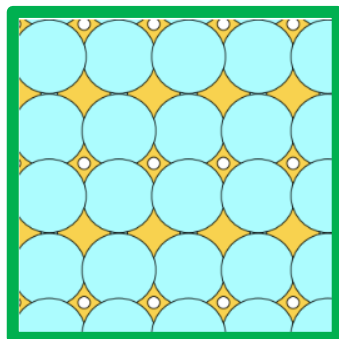
Adsorbate	Material	Crystal	Facet	Site	ML $E_{ads}$ (eV)	DFT $E_{ads}$ (eV)
<b>(H)</b>	CaNi	Cubic	101	bridge-CaNi	-0.46	-
	YIr	Cubic	101	ontop-Ir	-0.45	-0.68
	<b>YAu</b>	<b>Cubic</b>	<b>100</b>	<b>hollow</b>	<b>-0.47</b>	<b>-0.42</b>
	YAu	Cubic	101	hollow	-0.53	-
<b>(O)</b>	<b>ZnIr</b>	<b>Hexagonal</b>	<b>10-11</b>	<b>hollow</b>	<b>-1.72</b>	<b>-1.93</b>
	<b>ZnPt</b>	<b>Tetragonal</b>	<b>101</b>	<b>hollow</b>	<b>-1.76</b>	<b>-1.81</b>
<b>(OH)</b>	<b>CdAu</b>	<b>Cubic</b>	<b>101</b>	<b>longbridge-Au</b>	<b>1.10</b>	<b>1.06</b>
	CdAu	Hexagonal	0001	bridge-Au	1.10	0.89

# Discovered Materials

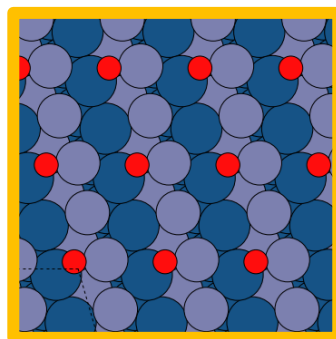
Target  $E_{ads}$        $\textcircled{\text{H}}$  -0.49 eV       $\textcircled{\text{O}}$  -1.79 eV       $\textcircled{\text{O}}\textcircled{\text{H}}$  -1.19 eV

Adsorbate	Material	Crystal	Facet	Site	ML $E_{ads}$ (eV)	DFT $E_{ads}$ (eV)
$\textcircled{\text{H}}$	CaNi	Cubic	101	bridge-CaNi	-0.46	-
	YIr	Cubic	101	ontop-Ir	-0.45	-0.68
	<b>YAu</b>	<b>Cubic</b>	<b>100</b>	<b>hollow</b>	<b>-0.47</b>	<b>-0.42</b>
	YAu	Cubic	101	hollow	-0.53	-
$\textcircled{\text{O}}$	<b>ZnIr</b>	<b>Hexagonal</b>	<b>10-11</b>	<b>hollow</b>	<b>-1.72</b>	<b>-1.93</b>
	<b>ZnPt</b>	<b>Tetragonal</b>	<b>101</b>	<b>hollow</b>	<b>-1.76</b>	<b>-1.81</b>
$\textcircled{\text{O}}\textcircled{\text{H}}$	CdAu	Cubic	101	longbridge-Au	1.10	1.06
	CdAu	Hexagonal	0001	bridge-Au	1.10	0.89

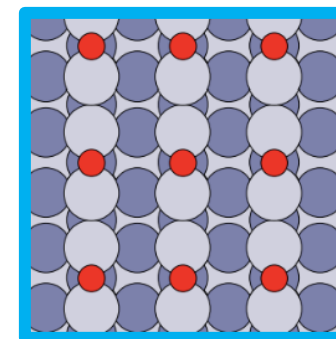
YAu bcc(100)



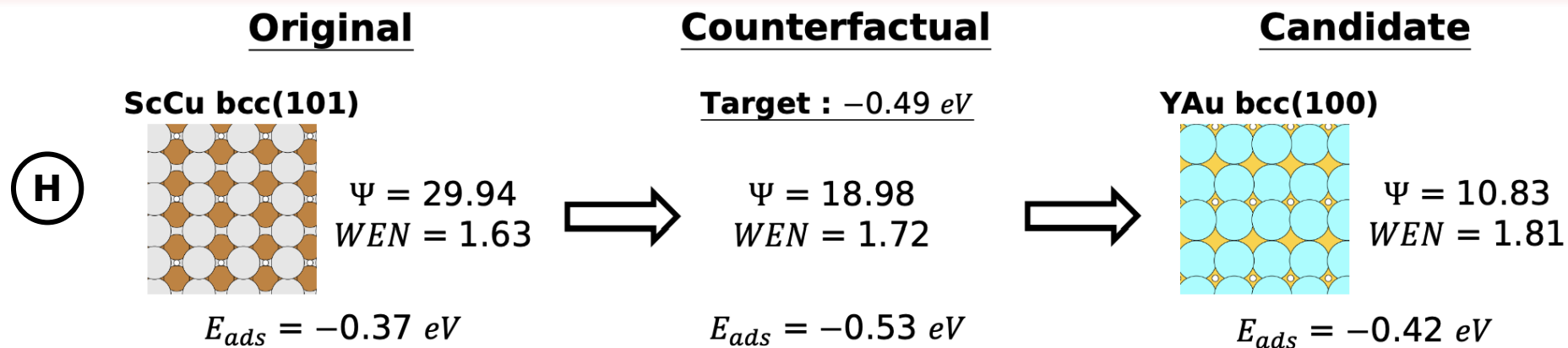
ZnIr hcp(10-11)



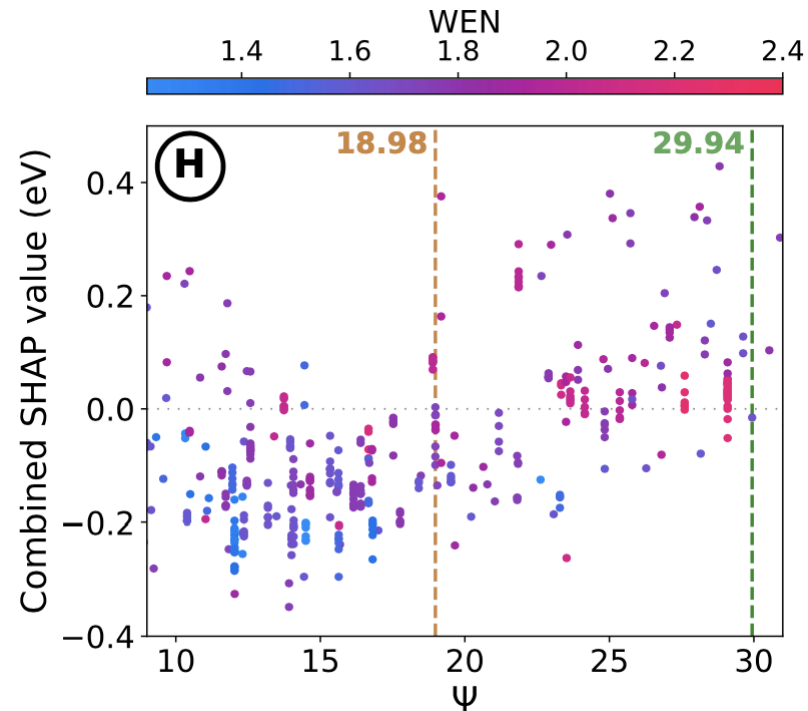
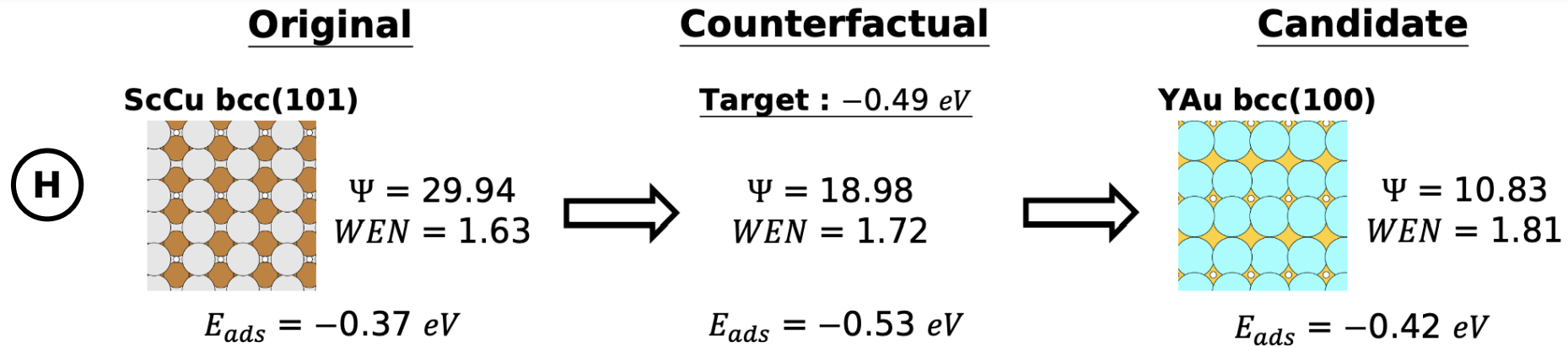
ZnPt bct(101)



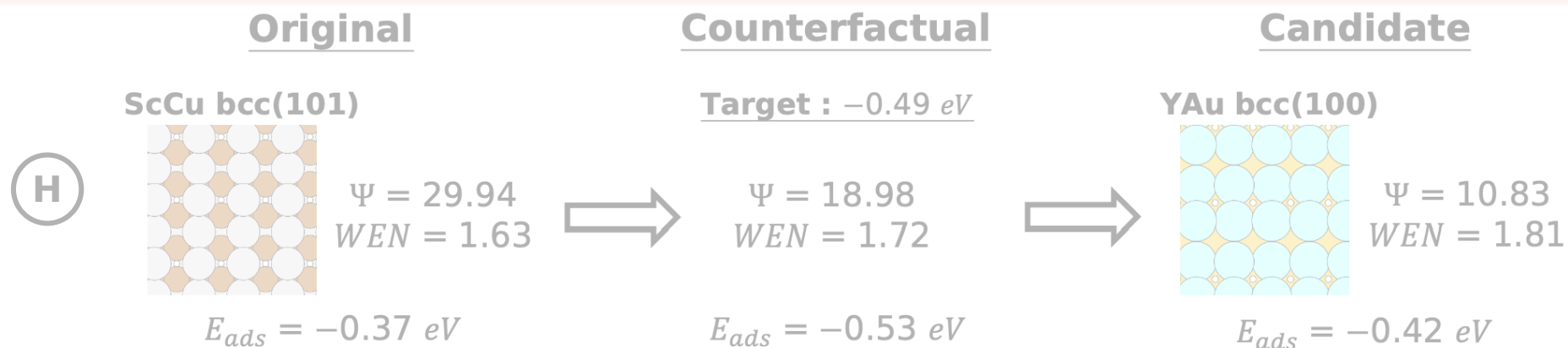
# Retrieving Explanations



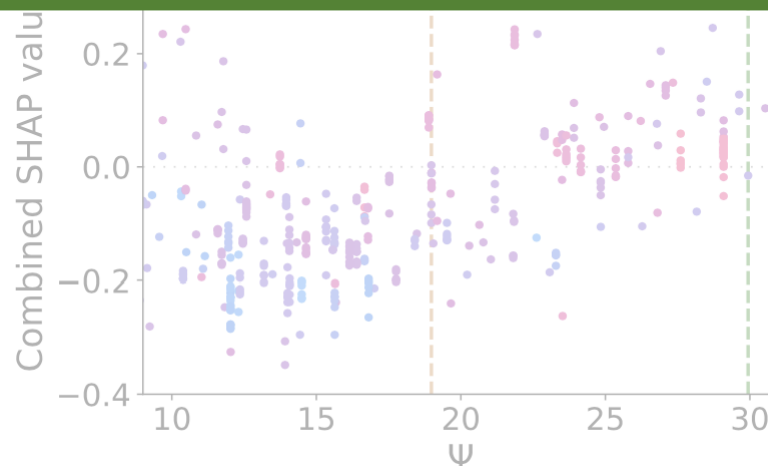
# Retrieving Explanations



# Retrieving Explanations

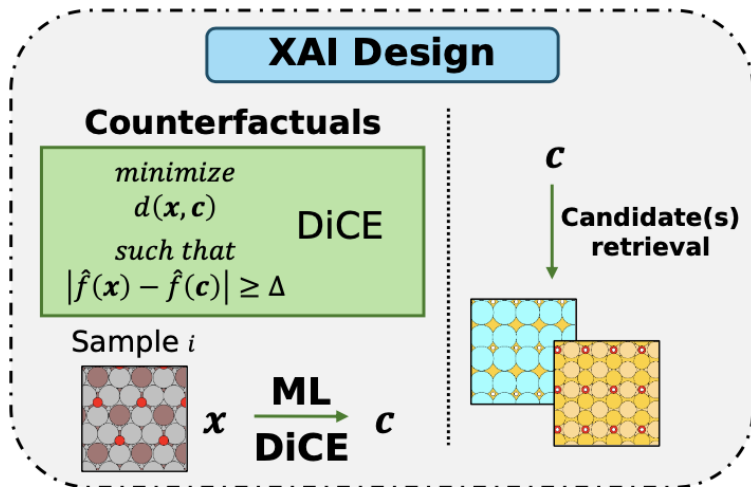


Counterfactuals unveiled subtle relationships between the most relevant features, other, in principle, less important features, and the  $E_{ads}$ .



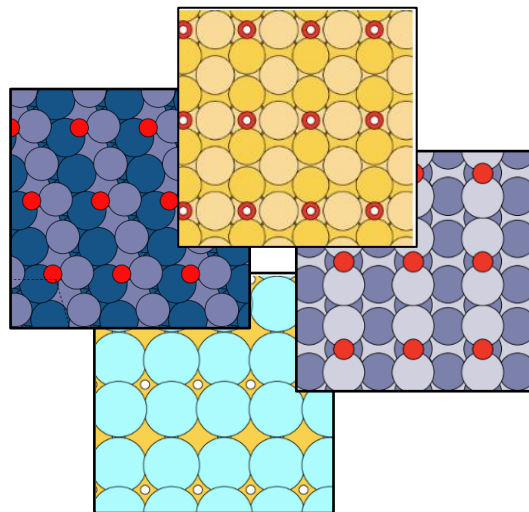
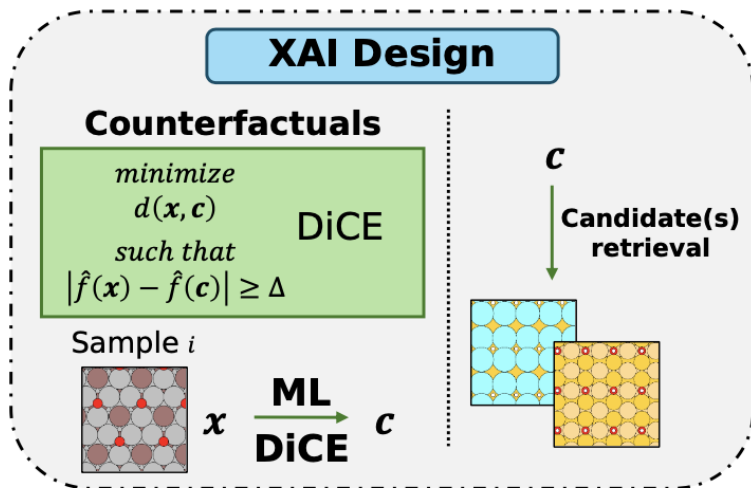
# Conclusions

- A **novel strategy** for the **discovery and design of new materials** based on **counterfactual explanations** was proposed



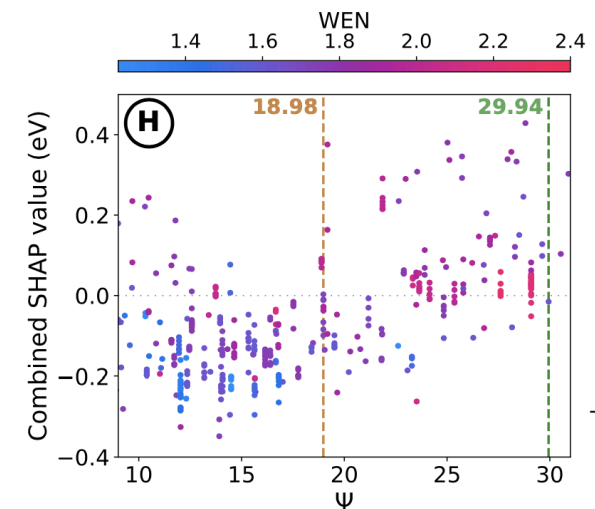
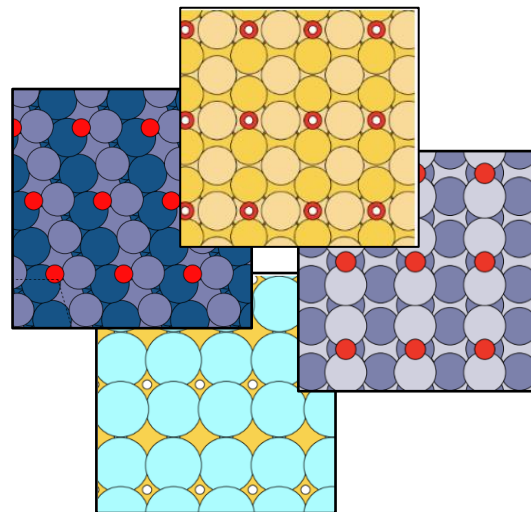
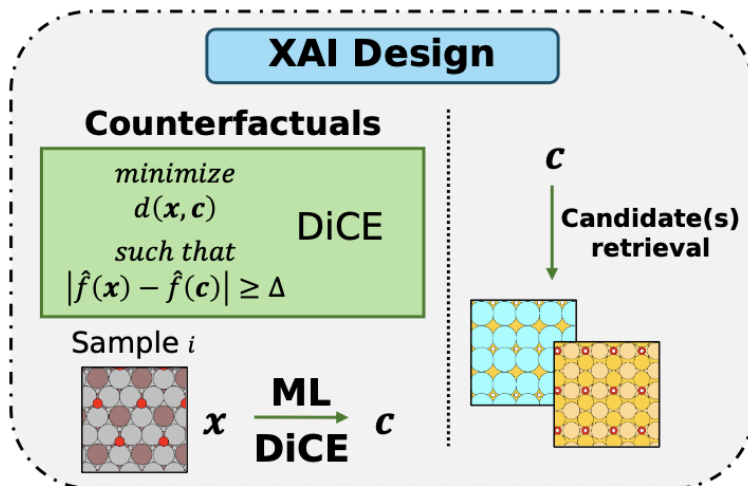
# Conclusions

- A **novel strategy** for the **discovery and design of new materials** based on **counterfactual explanations** was proposed
- **Four** of the **discovered materials** were **confirmed** by reference **DFT calculations, validating** our strategy



# Conclusions

- A **novel strategy** for the **discovery and design of new materials** based on **counterfactual explanations** was proposed
- **Four** of the **discovered materials** were **confirmed** by reference **DFT calculations, validating** our strategy
- By **comparing original samples, counterfactuals, and discovered candidates**, subtle **relationships** were **unveiled** between the most relevant features, other, in principle, less important features, and the  $E_{\text{ads}}$



# Acknowledgements

***XAI will not replace human experts, but human experts who embrace XAI will replace those who do not***

Oviedo, F., et al. *Acc. Mater. Res.* **3**, 597-607 (2022).

# Acknowledgements

***XAI will not replace human experts, but human experts who embrace XAI will replace those who do not***

Oviedo, F., et al. *Acc. Mater. Res.* **3**, 597-607 (2022).

Check the paper here!



Python Routines



# Acknowledgements

***XAI will not replace human experts, but human experts who embrace XAI will replace those who do not***

Oviedo, F., et al. *Acc. Mater. Res.* **3**, 597-607 (2022).



Universidad  
de Navarra

DATAI  
INSTITUTE OF DATA SCIENCE  
AND ARTIFICIAL INTELLIGENCE

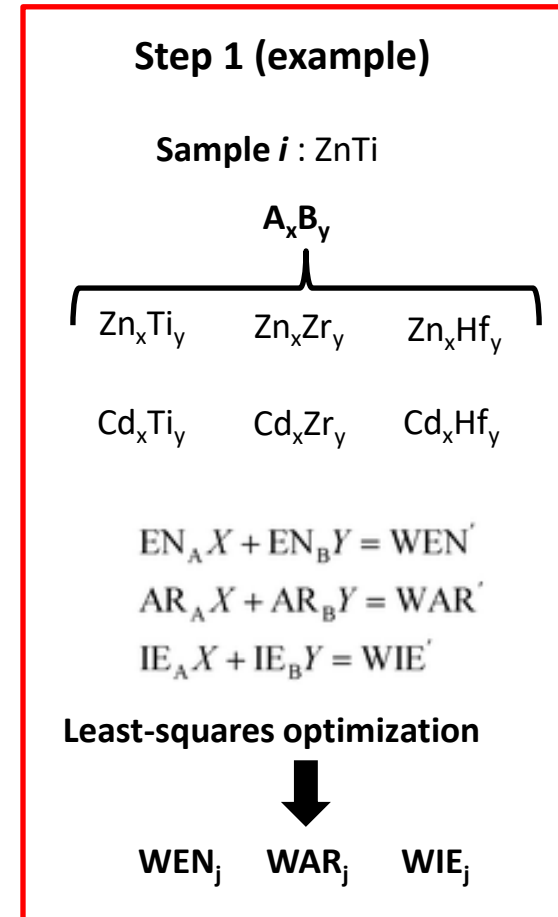
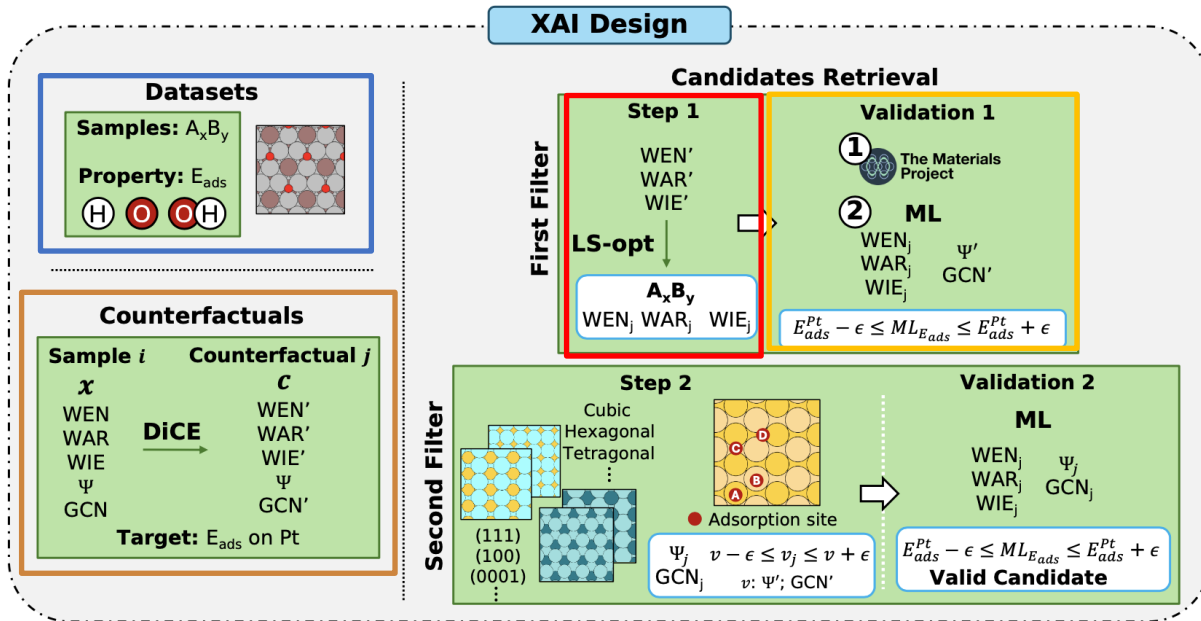
Check the paper here!



Python Routines



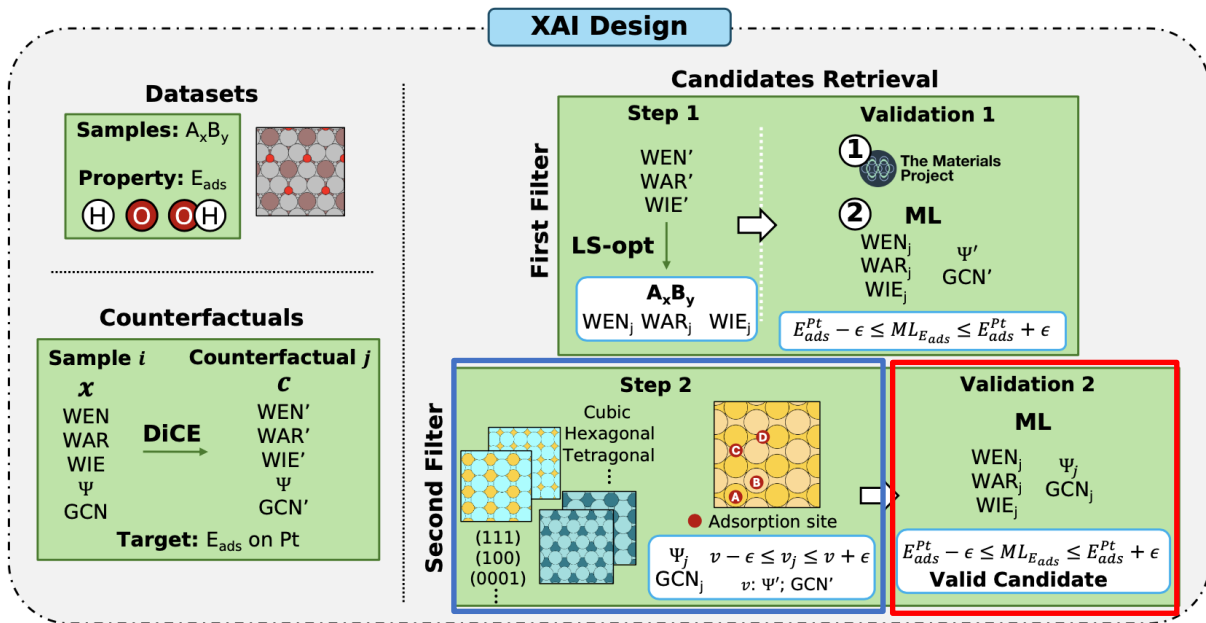




🧠 We start with our databases of  $E_{ads}$  for H, O, and OH

🧠 We generate the counterfactuals with DICE for each sample in the database

- 🧠 Check that the  $A_xB_y$  intermetallics found in Step 1 exist in Materials Project
- 🧠 Use ML model to verify that the  $E_{ads}$  of these intermetallics is still close to the target



- 👤 All possible crystal lattices are considered and all surface slabs for different  $(hk(i)l)$  facets are constructed ( $h, k, l \leq 3$ )
- 👤 All available adsorption sites are determined and values of  $\Psi'$  and GCN' are updated accordingly ( $\Psi_j, GCN_j$ )
- 👤 We only keep those samples for which  $\Psi_j$  and GCN<sub>j</sub> are close to the counterfactual.

👤 Use ML model to verify that the  $E_{ads}$  of these intermetallics is still close to the target

👤 Explainability is ensured by construction since the discovered materials can be linked directly to the original sample of the dataset from which the counterfactual was generated.