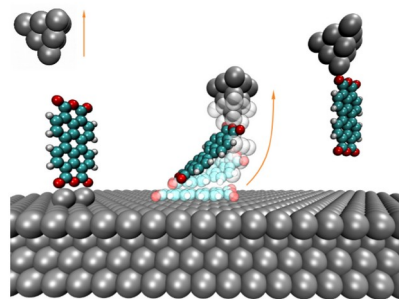


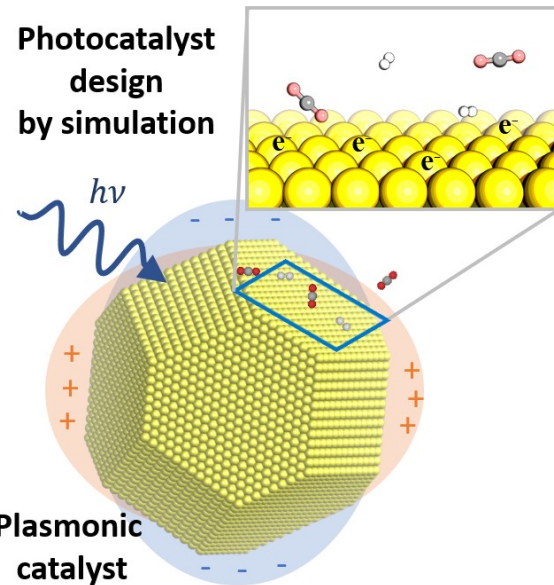


Computational Surface Science

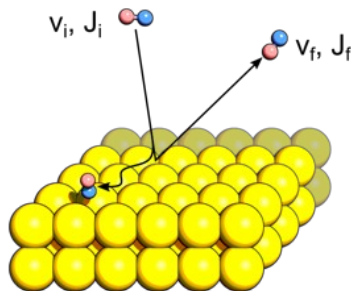
Metastable nanostructures



Photocatalyst design by simulation



Surface Chemistry



Our approach



Electronic Structure Theory

Machine Learning
of Electronic
Structure

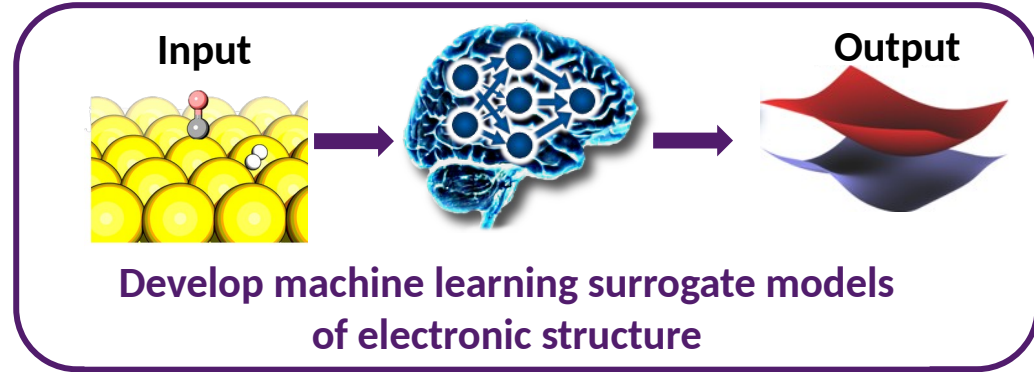
Coupling of Light
and Electrons



Nonadiabatic
Molecular
Dynamics

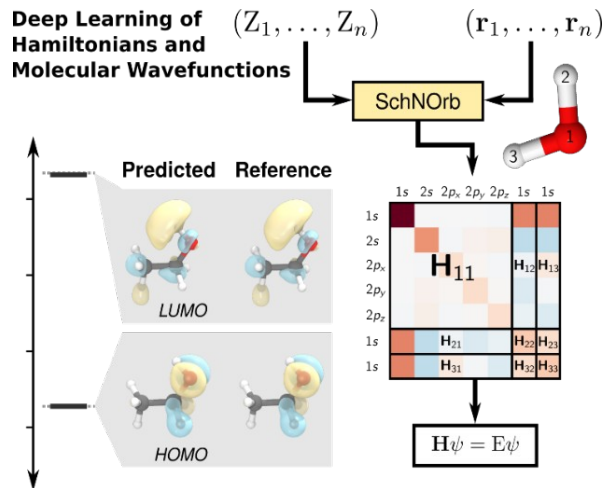


NQCDynamics.jl

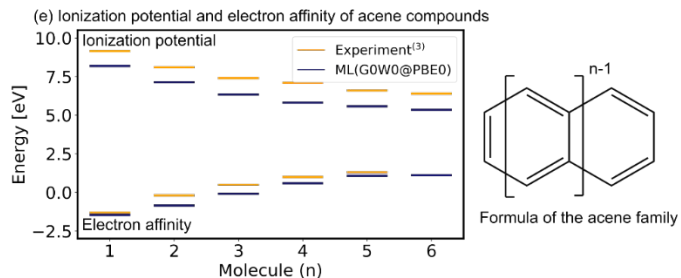


Machine Learning of Electronic Structure

Deep learning of molecular spectroscopy

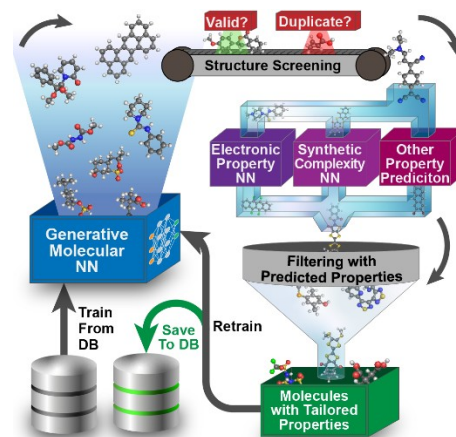


Nature Commun 10, 5024 (2019)



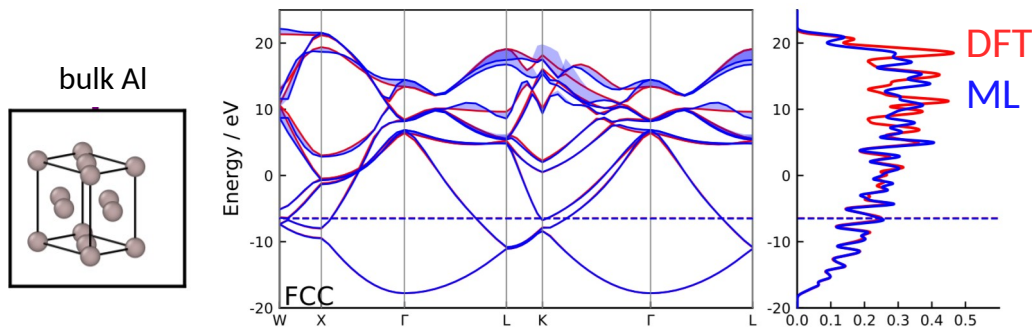
Chem. Sci. 12, 10755-10764 (2021)

Inverse property-driven design of electronic properties



Nature Comp. Sci. (2023)

Machine learning of electronic Hamiltonians



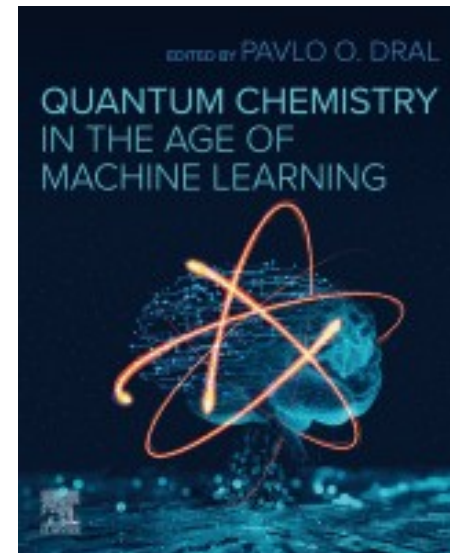
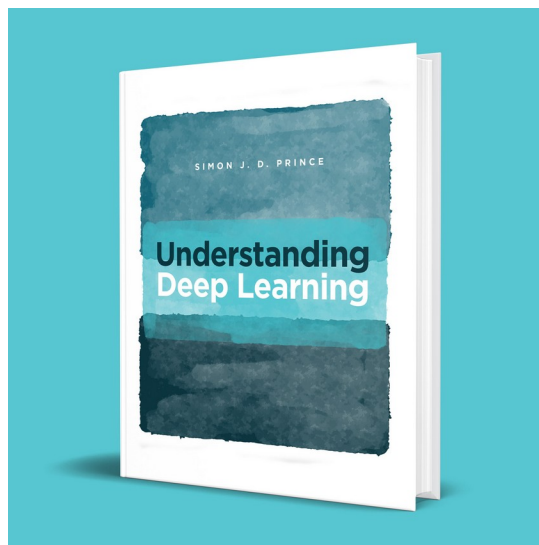
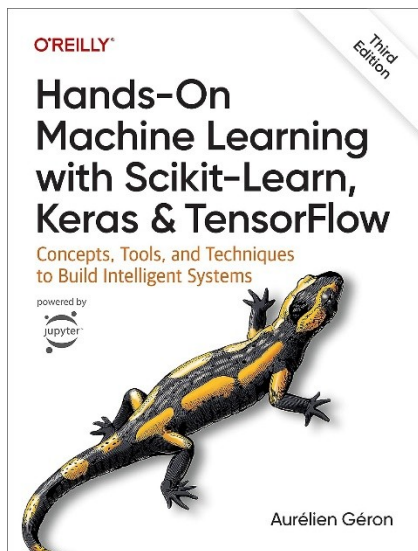
npj Computational Materials (2022)

Goal of this lecture

After this lecture (and accompanying workshop), you should ideally

- Know the **role of machine learning** in the computational physical sciences
- Understand the **basic terminology** of machine learning (“Slang busting”)
- Have an **overview of methodologies** and how they connect
- Understand how to approach a **typical machine learning workflow**
- Know how to **prepare and analyse datasets**
- Be able to validate and optimise models with **cross-validation**
- Know basic approaches to **featurisation and representation** in chemistry
- Know how to evaluate and assess **prediction errors and uncertainties**

Resources

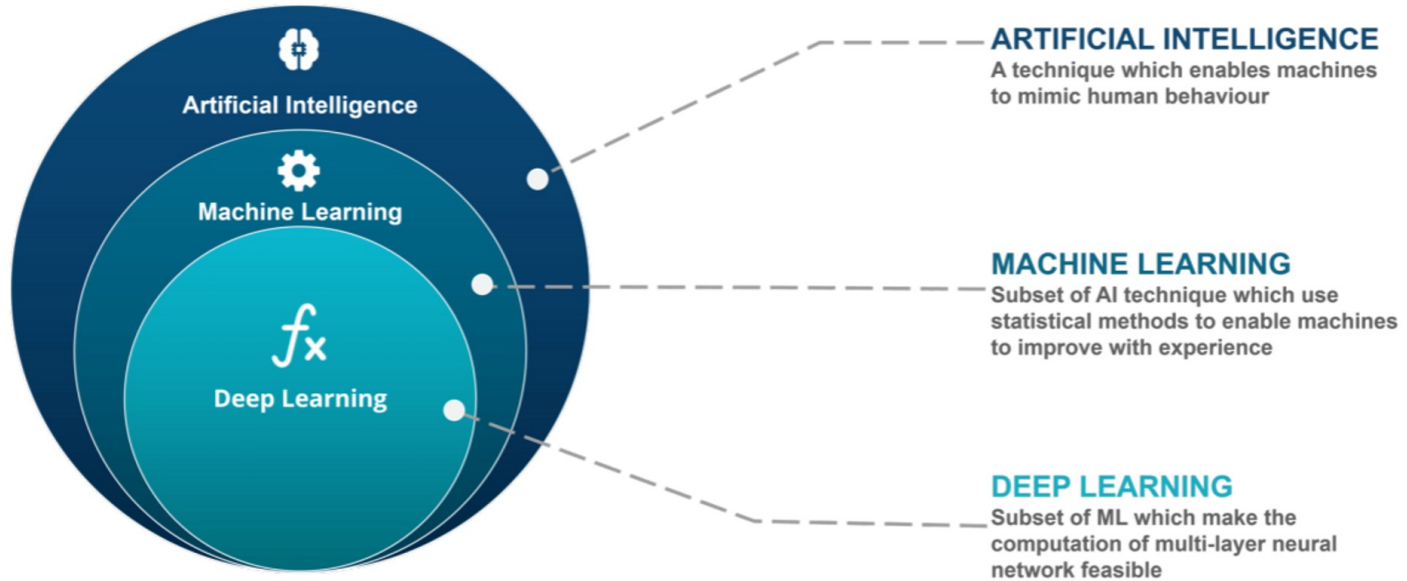


Perspective on integrating machine learning into computational chemistry and materials science

Agenda

- What is ML?
- Basic Definitions
- Data Representations and Features / Descriptors
- Types of ML methods
- Putting it all together
- A research example

What is ML and how can it be useful?

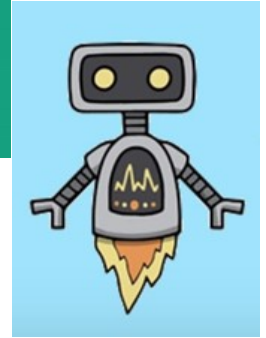


Learning from Data

“Machine Learning is the field of study that gives the computer the Ability to learn without being explicitly programmed.”, Arthur Samuel

ML is everywhere

- Medical outcome analysis
- Robots
- Autonomous driving
- Speech recognition, natural language processing
- Image recognition & creation, deep fakes, ...



Commonly used in Chemistry and Materials Science research

The ML paradigm and your email spam filter

Classical Computing Paradigm

Determines
category

Predefined set of criteria
that the code checks when
analyzing email content



Mail



Inbox



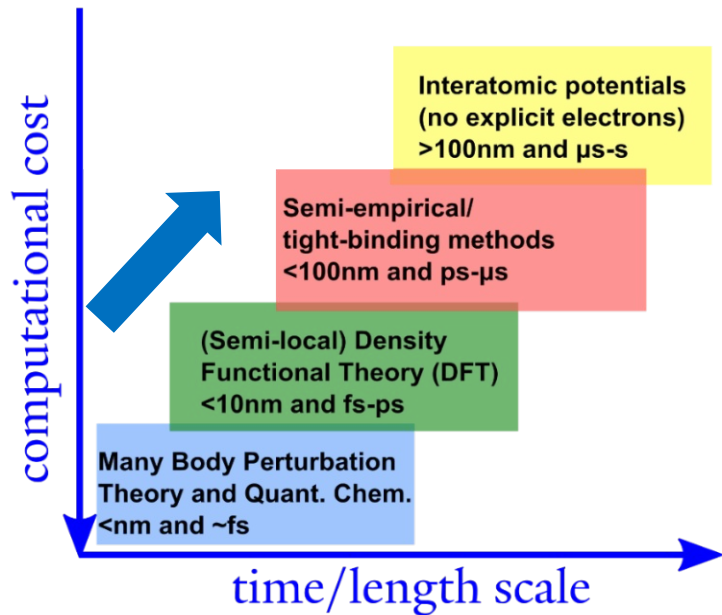
Machine Learning Paradigm

Predicts
category

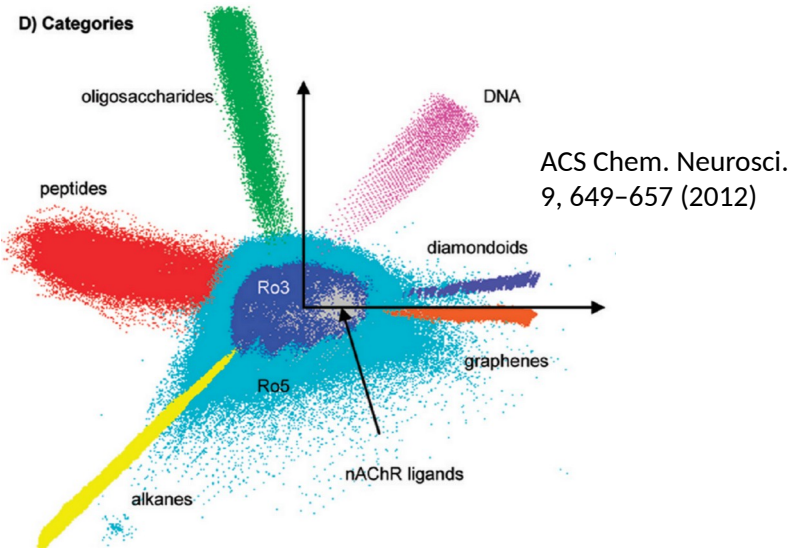
Learns from user/provider
data what features of
emails are most likely
associated with spam

Challenges in Chemistry & Materials Science

Accurate Quantum Chemistry is slow
Fast Quantum Chemistry is inaccurate



Chemical compound space is bigger than the known universe



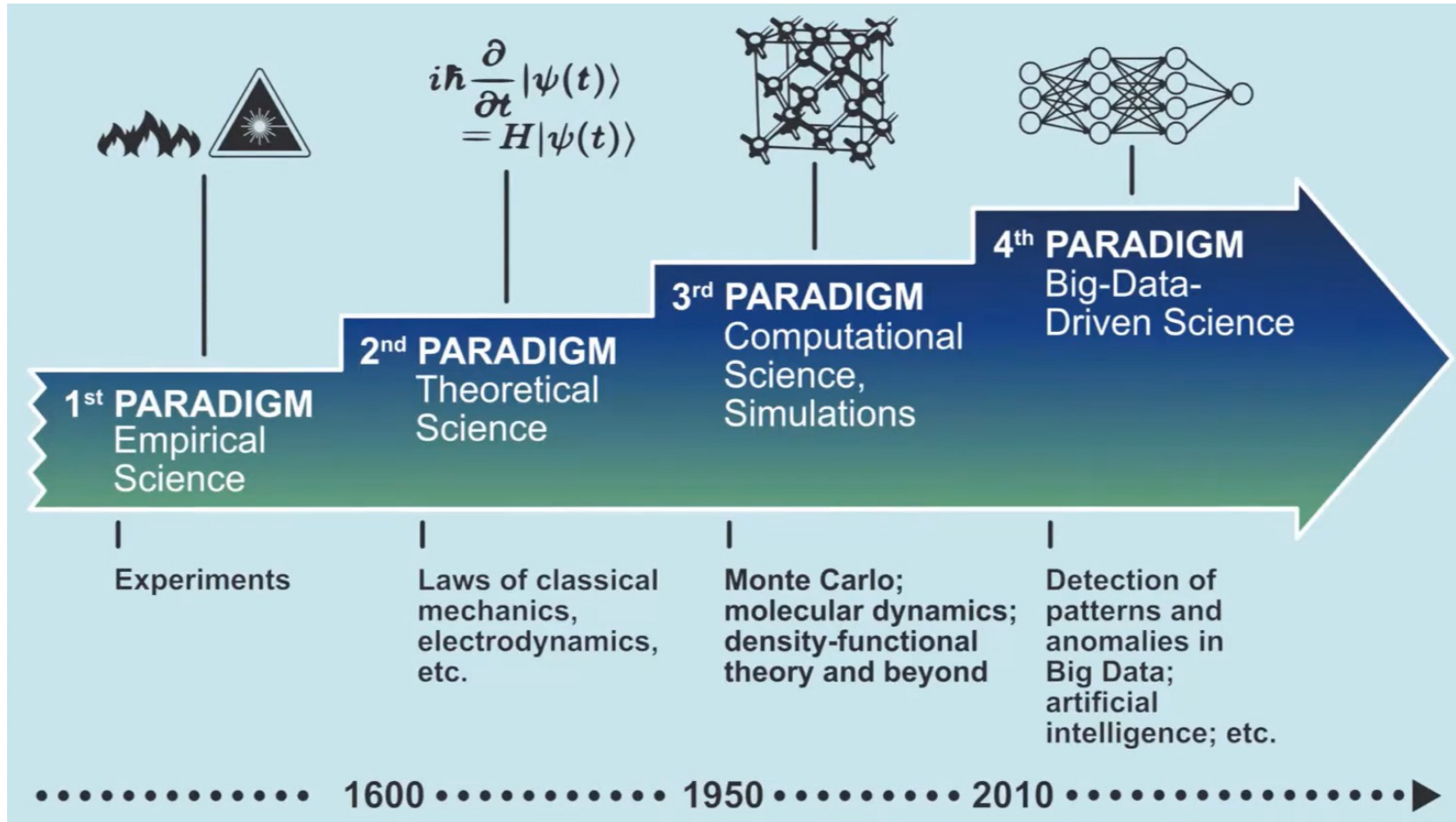
10^{180} possible molecules 10^{80} atoms in universe
 10^8 known molecules

Accelerate property prediction

Machine Learning
can help!

Accelerate materials design

The 4th paradigm of science



Basic Definitions

ML Definitions

ML is concerned with algorithms that improve with increasing amount of available data under some performance measure.

Find a predictive function that connects input space \mathcal{X} to a target space \mathcal{Y}

$$f: \mathcal{X} \rightarrow \mathcal{Y}$$

ML focuses on **universal approximators**, able to represent any function with arbitrary accuracy, when given enough training data and parameters.

The functional relationship to be found is specified by choosing a suitable **loss function** $\ell(f(x), y)$.

If the loss function requires knowledge of targets (labels) $y \in \mathcal{Y}$, we speak of **supervised learning**.

How it works

1. Provide a training dataset, T ← **Training Dataset: data that model sees to learn**
2. The machine is **trained** to capture statistical trends in T and suggests a function f
3. Ideally, $f(x_i) = y_i$ for all values
4. Then, given new (unseen) x' values, the model should **predict** $f(x') = y'$
5. We can validate how much the model has learned with a **test set** of data for which the ground truth is known

Test/Validation Dataset

Data that is unseen to assess learning. Generate by splitting off ~20% of data.

We train a model by minimising risk

The optimal model minimizes the expected risk, $R(f)$, defined as the **expectation value** of $\ell(f(\mathcal{X}), \mathcal{Y})$:

$$R(f) = \langle \ell(f(\mathcal{X}), \mathcal{Y}) \rangle = \int \ell(f(\mathcal{X}), \mathcal{Y}) \, d\mathcal{P}(\mathcal{X}, \mathcal{Y})$$

The ML training process is the process of minimizing R_{emp} or the **loss**.

Different types of loss functions

Different ways to measure distances between the ground truth and the model prediction f by a single number.

L1-norm

$$J(f) = \sum_i^n |f_i(x_i) - y_i| \quad \text{Least absolute deviations}$$

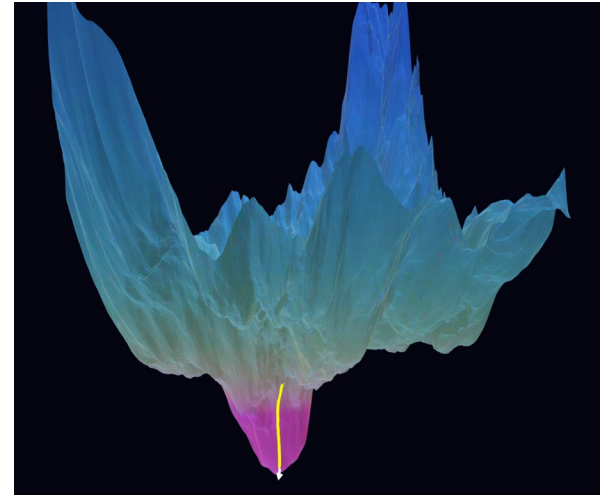
L2-norm

$$J(f) = \sum_i^n (f_i(x_i) - y_i)^2 \quad \text{Least square regression}$$

L1 norm more robust against outliers, but more difficult to optimize

Training and Overfitting

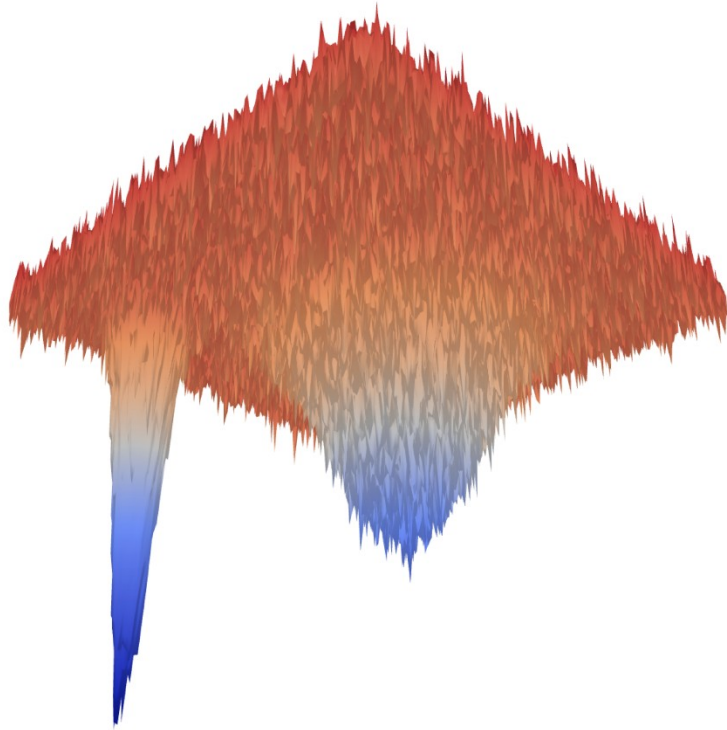
- Usually **no function $f(x)$ perfectly captures the training data**. $f(x_i) = y_i$ does not hold and there is an error.
- Many functions can map $x \rightarrow y$.
- The accuracy of a function is determined by the loss function $\ell(f(x) = y)$.
- **Training is the process of minimising the loss** (e.g. through stochastic gradient descent).
- Some functions are complex, some are simple.
- Overly complex functions lead to unstable predictions → **overfitting**.
- We introduce a **regularisation** during training to punish complex functions and avoid overfitting.



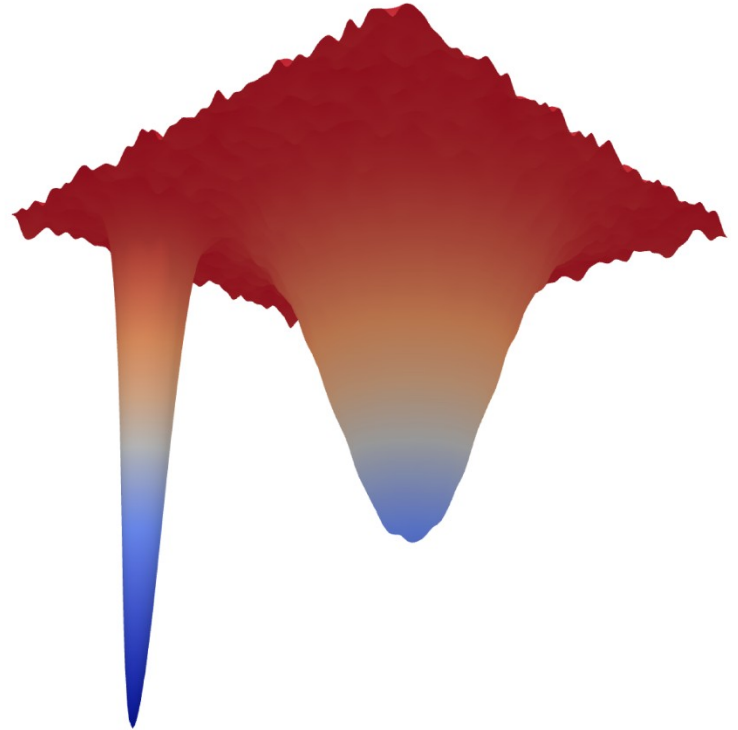
Loss landscape

<https://losslandscape.com/explorer>

Without regularisation

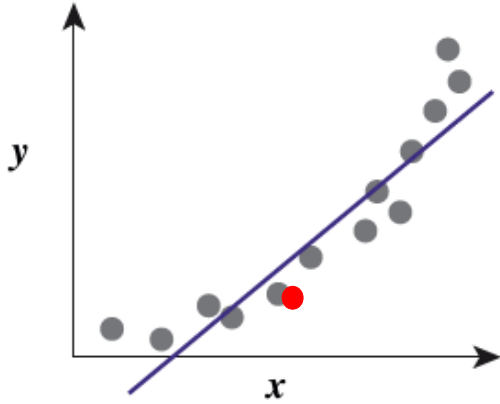


With regularisation



Overfitting yields an increased error on unseen data by approximating a simple functional relationship with an overly complex function on the training set

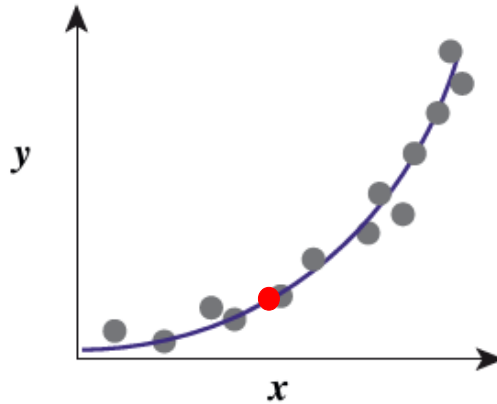
underfitting



training error: **high**
validation error: **high**

Model is not sufficiently complex and flexible to capture data

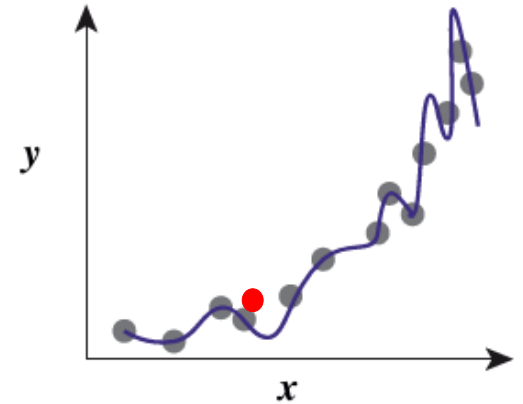
appropriate fitting



training error: **low**
validation error: **low**

Model is suitably complex to capture trends in data

overfitting



training error: **low**
validation error: **high**

Model is overly complex and will likely not **generalize to other samples**

How do we quantify the performance of a model?

Residual sum of squares

Coefficient of determination

$$R^2 = 1 - \frac{\sum_i (y_i - f_i)^2}{\sum_i (y_i - u)^2}$$

Total sum of squares (distance from mean u)

Mean Absolute Error

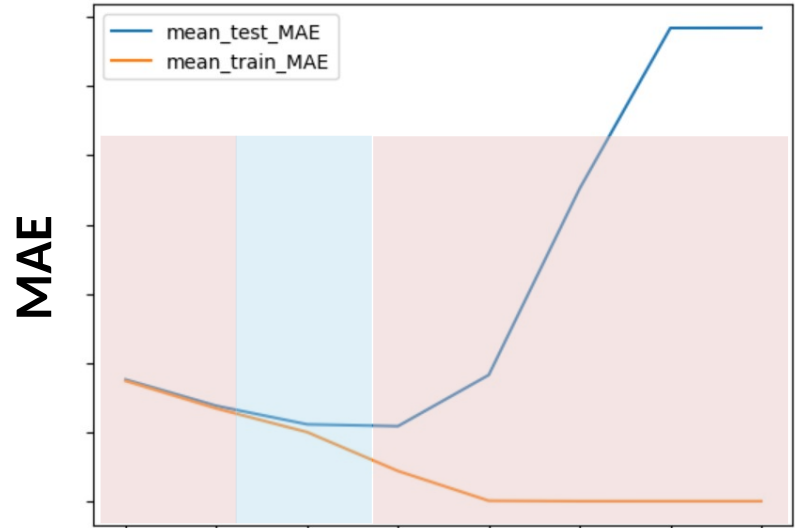
$$\text{MAE} = \frac{\sum_i |y_i - f_i|}{n}$$

Root Mean Squared Error

$$\text{RMSE} = \sqrt{\frac{\sum_i (y_i - f_i)^2}{n}}$$

y_i ... ground truth label for data point i

f_i ... model prediction for data point i



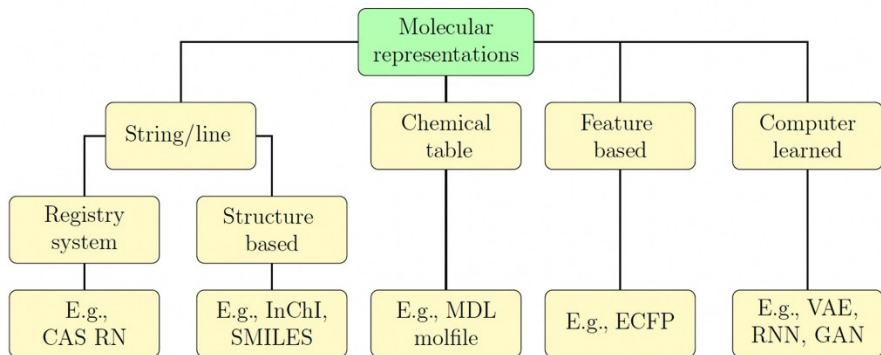
“flexibility of the model”
(e.g. number of parameters)

Try to achieve lowest possible test set error

Data Representation and Features

What is the input x in $y = f(x)$?

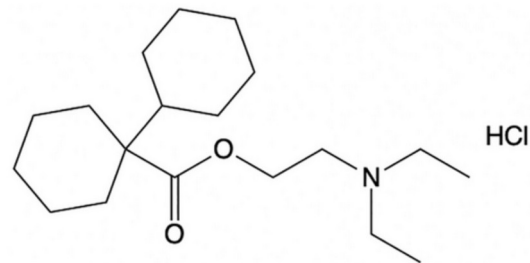
Examples of Molecular Data Representations



- Registry systems

Chemical Abstract Services Registry Number
PubChem CID
ChemSpider
ChEMBL

- Lewis structures (2D graph)



- String representations

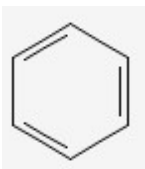
Generic names ⁵	Dicycloverine HCl, benacol, bentyl, dibent, Dyspas, and so on
Mol. formula	C ₁₉ H ₃₆ ClNO ₂
IUPAC name	2-(Diethylamino)ethyl 1-cyclohexylcyclohexane-1-carboxylate hydrochloride
CAS RN	67 – 92 – 5
Canonical SMILES	CCN(CC)CCOC(=O)C1(CCCCC1)C2CCCC2.Cl
InChI	InChI = 1S/C19H35NO2.ClH/c1-3-20(4-2)15-16-22-18(21)19 (13-9-6-10-14-19)17-11-7-5-8-12-17;/h17H,3-16H2,1-2H3;1H
	InChIKey:GUBNMFJOJDCEL-UHFFFAOYSA-N
WLN ⁶	L6TJA-AL6TJAVO2N2&2&GH

SMILES

Simplified Molecular Input Line Entry System
String of ASCII characters that defines a molecule

- Atoms are represented with their chemical symbol
- Single bonds are implicit or (-), double bonds (=), triple bonds (#)
- Rings represented via a number after the initial atom and closing atom (e..g C1CCNCC1)
- Branching: represented with parentheses around the branch
- Aromaticity: aromatically bonded atoms in lower case or alternating -C=C-

C1ccccc1,
C1=C-C=C-C=C1,
C1=CC=CC=C1

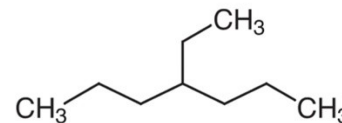


- Hard to represent unsaturated bonds, radicals, or unusual valences/bonding
- No info on molecular conformation
- Not well suited for generative models

DeepSMILES

SELFIES (self-referencing embedded strings)

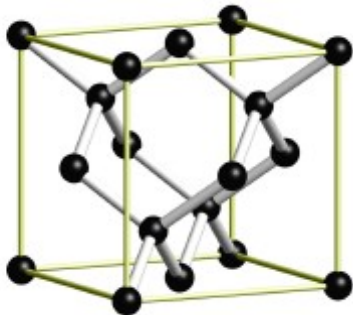
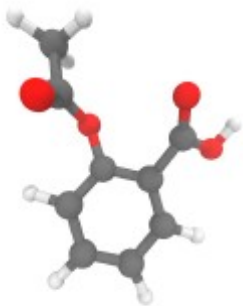
4-ethylheptane:
CCCC(CC)CCC



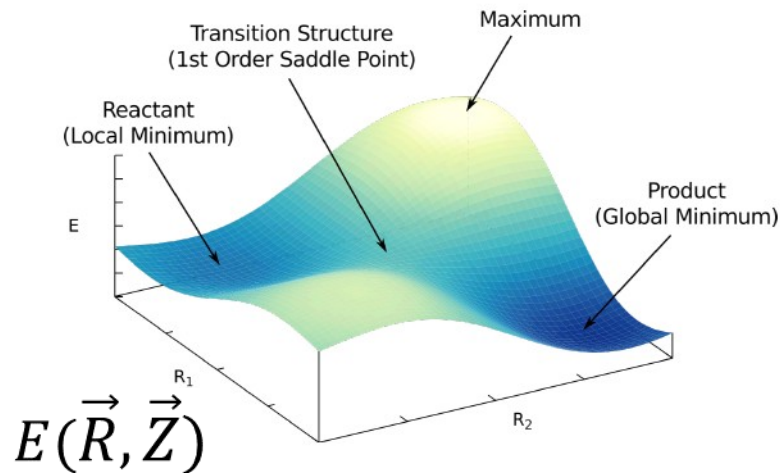
Data Representation / Featurization for Molecules & Materials

$$\hat{H}\Psi = E\Psi$$

Molecules & Materials



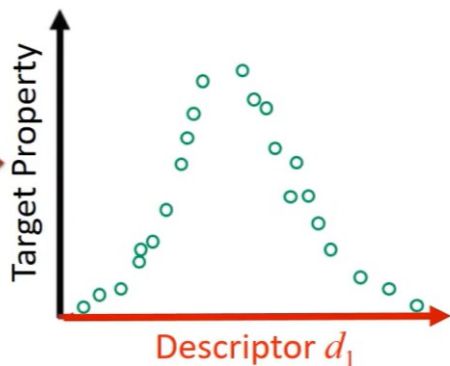
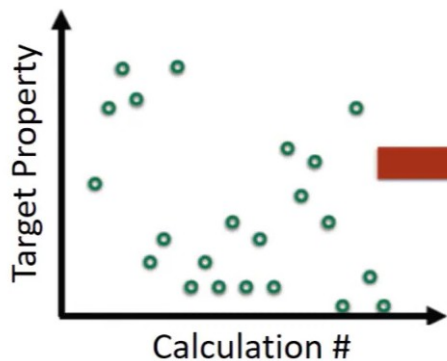
$$S = \{\vec{R}, \vec{Z}\}$$



Energy landscape encodes:

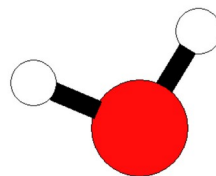
- Rotational/Translational invariance
- Permutation invariance
- Symmetry/Local Features

XYZ atomic positions give a terrible representation



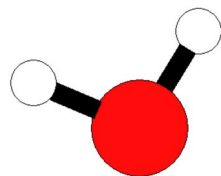
Bad representation for target

Good representation



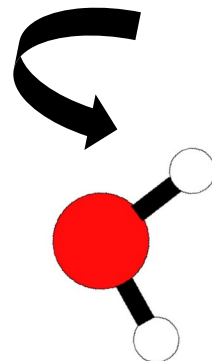
```
[[ 0.      0.      0.119262]
 [ 0.      0.763239 -0.477047]
 [ 0.     -0.763239 -0.477047]]
```

```
[[5.      7.      9.119262]
 [5.      7.763239 8.522953]
 [5.      6.236761 8.522953]]
```

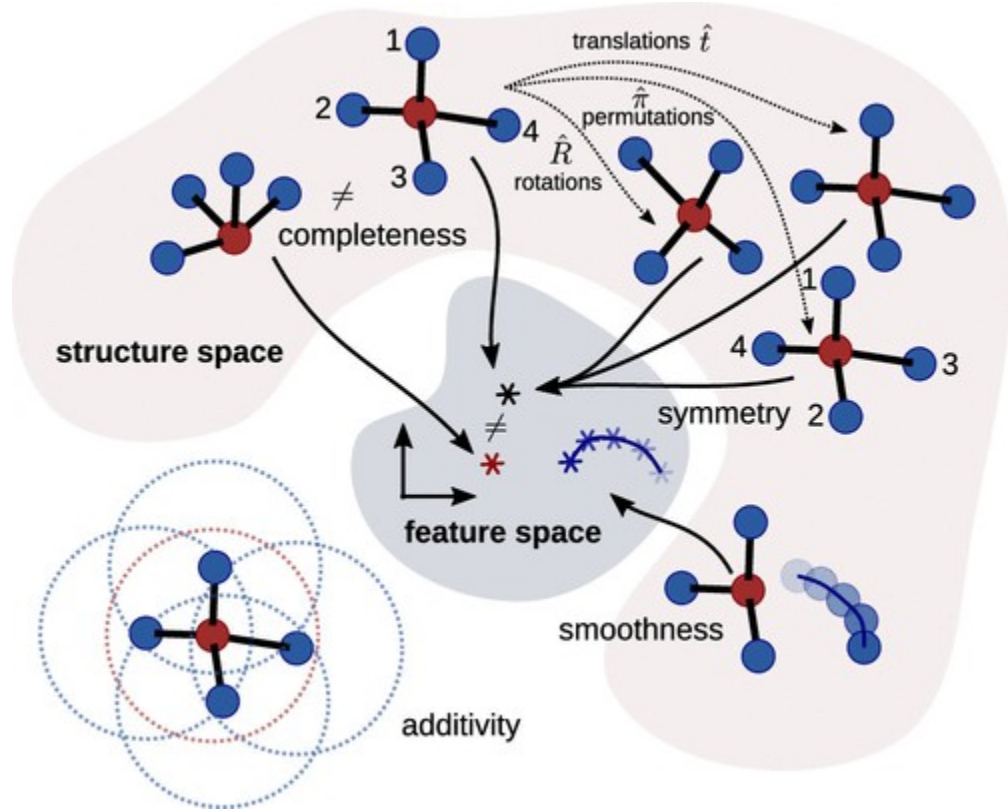


XYZ coordinates do not capture the basic structural symmetry properties of molecules and materials

```
[[ -7.      5.      9.119262]
 [ -7.763239 5.      8.522953]
 [ -6.236761 5.      8.522953]]
```



Requirements on molecular representations (a.k.a. descriptors)



Symmetry Invariance:

Mapping onto the same point in feature space

Completeness & Uniqueness:

Being able to distinguish all symmetry-inequivalent arrangements by mapping onto different points in feature space

Smoothness:

Smooth changes in atomic positions lead to smooth changes in feature space

Nomenclature

- ▶ **Structure:** Atomic positions, chemical species, unit cell
- ▶ **Descriptor:** A specific method for transforming the structure of different molecules/materials into a constant vector with correct symmetry properties
- ▶ **Feature vector:** A vector of numbers produced by converting the structure according to a certain descriptor. The feature vector is the structural representation in “feature space”

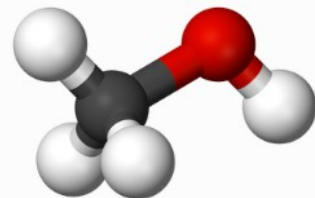
Global descriptors

Atom-wise descriptors

Coulomb Matrix

$$M_{ij}^{\text{Coulomb}} = \begin{cases} 0.5Z_i^{2.4} & \text{for } i = j \\ \frac{Z_i Z_j}{R_{ij}} & \text{for } i \neq j \end{cases}$$

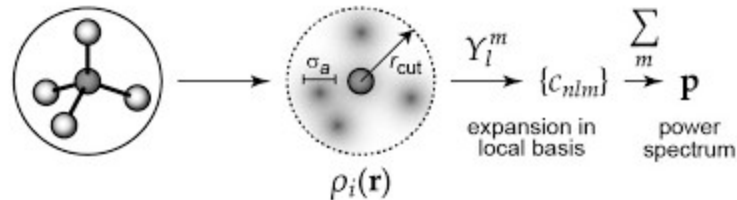
- A **global** representation
- Nuclear charges and atomic distances
- Molecules of different sizes require padding with zeros
- Vanilla version doesn't work for periodic systems
- But adaptations for periodic systems exist



36.9	33.7	5.5	3.1	5.5	5.5
33.7	73.5	4.0	8.2	3.8	3.8
5.5	4.0	0.5	0.35	0.56	0.56
3.1	8.2	0.35	0.5	0.43	0.43
5.5	3.8	0.56	0.43	0.5	0.56
5.5	3.8	0.56	0.43	0.56	0.5

Smooth Overlap of Atomic Positions (SOAP)

- Smooth Overlap of Atomic Positions (SOAP) encodes regions of atomic geometries (“atomic environments”)
- One set of features for each atom in the system with a cutoff
- Satisfies all relevant symmetries and feature requirements
- Perfect for interatomic potentials, but can be slow to calculate



Density smearing with Gaussians

$$\rho(\mathbf{r}) = \sum_i \exp\left(-\alpha\|\mathbf{r} - \mathbf{r}_i\|_2^2\right) f_{\text{cut}}(|\mathbf{r}_i|).$$

Representation via radial funcs and spherical harmonics

$$\rho(\mathbf{r}) = \sum_{nlm} c_{nlm} R_n(r) Y_{lm}(\hat{\mathbf{r}}).$$

Compute the power spectrum to ensure rotation invariance

$$p_{nn'l} = \sum_m c_{nlm}^* c_{n'lm}$$

Distance metric for each component of SOAP vector

$$k(\rho, \rho') = \sum_{nn'l} p_{nn'l} p'_{nn'l}$$

Atomic Cluster Expansion (ACE)

Atomic energy contributions

Designed for interatomic potential construction

$$E = \sum_i \varepsilon_i$$

Many-body cluster expansion

$$\varepsilon_i = V^{(0)}(Z_i) + \sum_{j_1} V^{(1)}(\mathbf{x}_{ij_1}) + \sum_{j_1 < j_2} V^{(2)}(\mathbf{x}_{ij_1}, \mathbf{x}_{ij_2}) + \dots + \sum_{j_1 < \dots < j_{\tilde{v}}} V^{(\tilde{v})}(\mathbf{x}_{ij_1}, \dots, \mathbf{x}_{ij_{\tilde{v}}})$$

Atomic basis: radial functions and spherical harmonics

$$\phi_{znlm}(\mathbf{r}_{ij}, Z_i, Z_j) = R_{nl}(r_{ij}, Z_i, Z_j) Y_l^m(\hat{\mathbf{r}}_{ij}) \delta_{zZ_j}$$

$$A_{znlm}^i = \sum_{j \in \mathcal{N}(i)} \phi_{znml}(\mathbf{r}_{ij}, Z_j, Z_i)$$

where $\mathcal{N}(i)$ denotes the set of indices of all atoms within the cutoff radius from atom i .

Expand each body order component in a basis

$$V^{(1)}(\mathbf{x}_{ij_1}) = \sum_{k_1} c_{k_1}^{(Z_i)} \phi_{k_1}(\mathbf{x}_{ij_1})$$

$$V^{(2)}(\mathbf{x}_{ij_1}, \mathbf{x}_{ij_2}) = \sum_{k_1, k_2} c_{k_1 k_2}^{(Z_i)} \phi_{k_1}(\mathbf{x}_{ij_1}) \phi_{k_2}(\mathbf{x}_{ij_2})$$

\vdots \vdots

$$V^{(\tilde{v})}(\mathbf{x}_{ij_1}, \dots, \mathbf{x}_{ij_{\tilde{v}}}) = \sum_{k_1, \dots, k_{\tilde{v}}} c_{k_1 \dots k_{\tilde{v}}}^{(Z_i)} \phi_{k_1}(\mathbf{x}_{ij_1}) \dots \phi_{k_{\tilde{v}}}(\mathbf{x}_{ij_{\tilde{v}}})$$

$$E = \text{[Diagram: sum of atomic basis functions]} + \dots + \text{[Diagram: sum of two-body basis functions]} + \dots + \text{[Diagram: sum of many-body basis functions]} + \dots$$

Body/correlation order products of basis functions

Product basis: for lexicographically ordered tuples $(z, \mathbf{n}, \mathbf{l}, \mathbf{m}) = (z_t, n_t, l_t, m_t)_{t=1}^v$ we define

$$A_{znlm}^i = \prod_{t=1}^v A_{z_t n_t l_t m_t}^i \quad (\text{A2})$$

$$E = \text{[Diagram: sum of many-body basis functions]} + \dots$$

Atomic Cluster Expansion (ACE)

Making basis functions rotationally invariant

$$B_{in}^{(1)} = A_{in00},$$

$$B_{in_1n_2l}^{(2)} = \sum_{m=-l}^l (-1)^m A_{in_1lm} A_{in_2l-m},$$

$$B_{in_1n_2n_3}^{(3)} = \sum_{l_1l_2l_3} \sum_{m_1=-l_1}^{l_1} \sum_{m_2=-l_2}^{l_2} \sum_{m_3=-l_3}^{l_3} \begin{pmatrix} l_1 & l_2 & l_3 \\ m_1 & m_2 & m_3 \end{pmatrix} \\ \times A_{in_1l_1m_1} A_{in_2l_2m_2} A_{in_3l_3m_3},$$

Atomic energy contributions become a series expansion in correlation order products of basis functions

$$\varepsilon_i = \sum_n c_n^{(1)} B_{in}^{(1)} + \sum_{n_1n_2l} c_{n_1n_2l}^{(2)} B_{in_1n_2l}^{(2)}$$

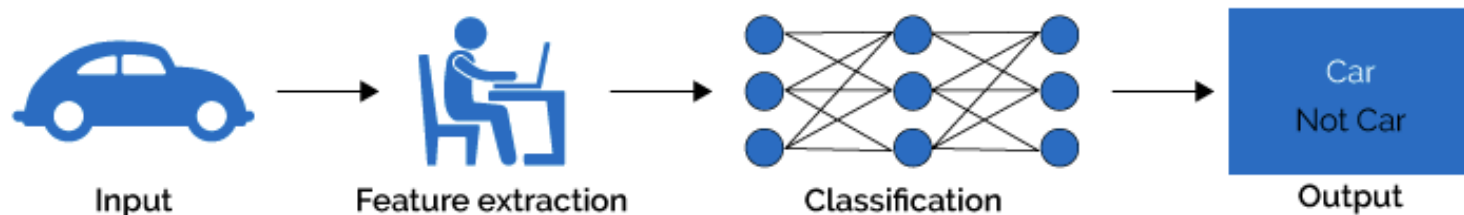
$$+ \sum_{n_1n_2n_3} \sum_{l_1l_2l_3} c_{n_1n_2n_3}^{(3)} B_{in_1n_2n_3}^{(3)}$$

$$\varepsilon_i = \sum_{Knl} c_{nl}^{(K)} B_{inl}^{(K)}$$

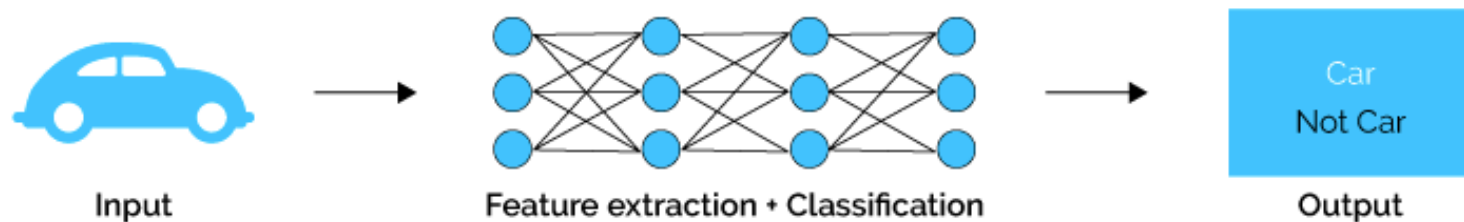
ACE is a generalisation of many possible atom-centred descriptors
(SOAP, ACSFs (Behler), SOAPS, Steinhardt parameters)

Side note on features

Machine Learning

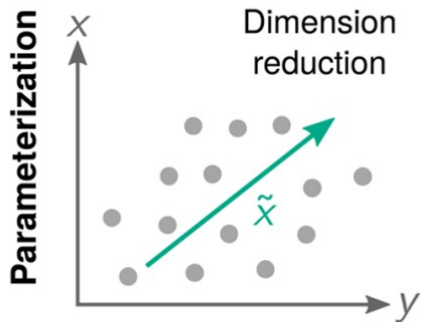


Deep Learning

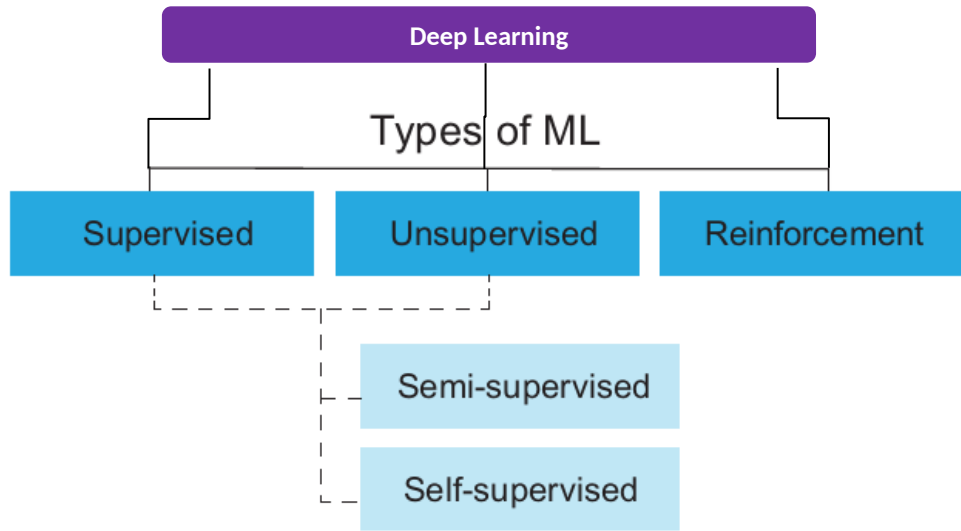
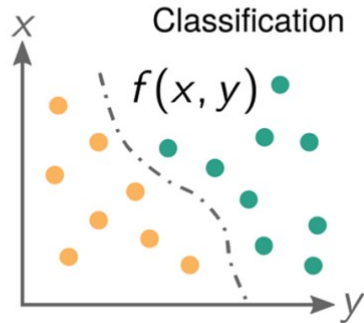
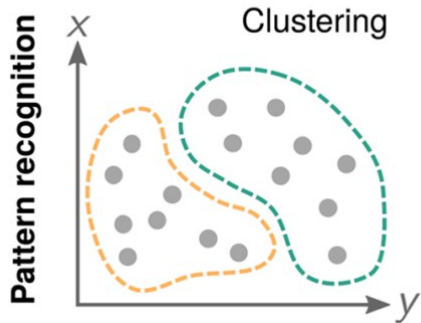
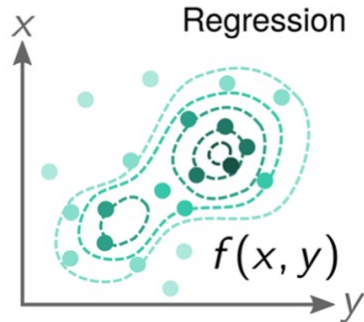


Types of ML methods

Unsupervised ML (unlabeled data)



Supervised ML (labeled data)



Exploratory analysis:
(leading to hypothesis)

- Clustering
- Dimensionality reduction

Confirmatory analysis:
(models that test ideas)

- Regression
- Classification

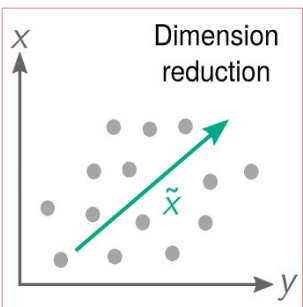
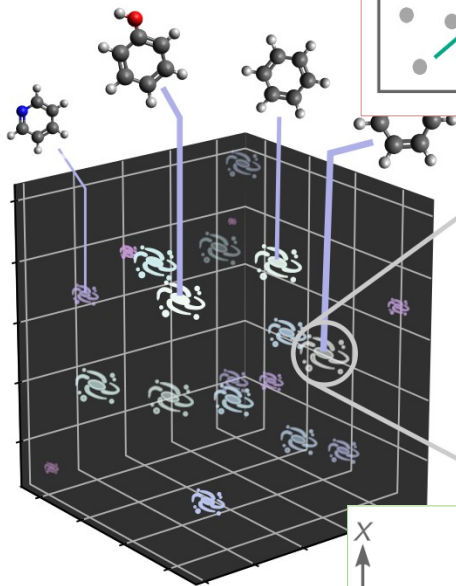
Deep learning approaches
exist for all types of ML

Machine learning methods in chemistry

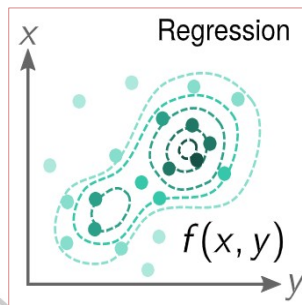
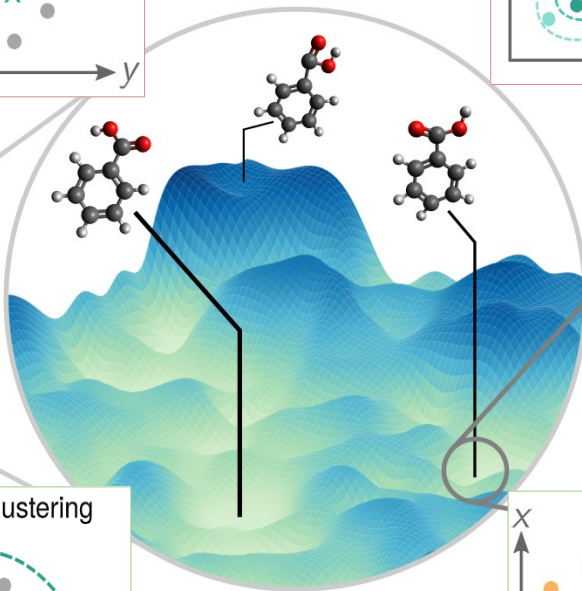
Parametrization

Composition

Chemical space

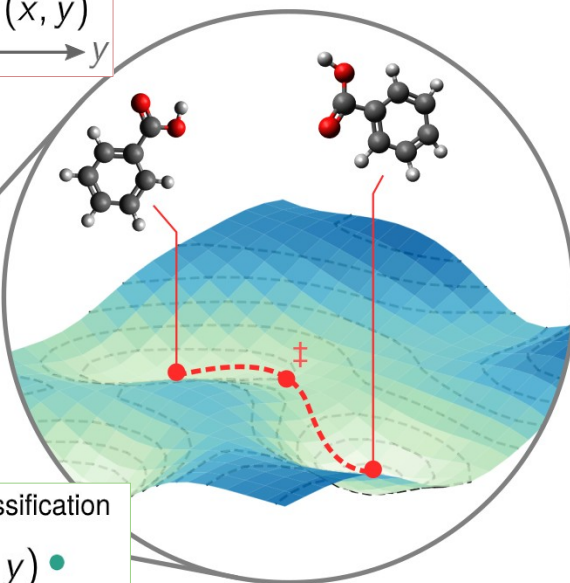


Global exploration

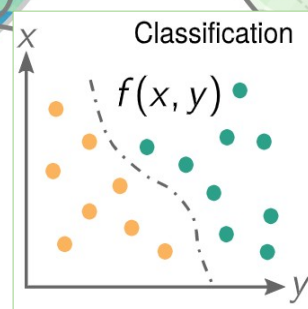
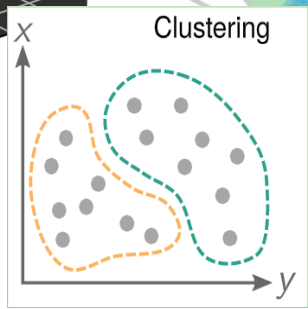


Configuration

Local exploration



Pattern recognition



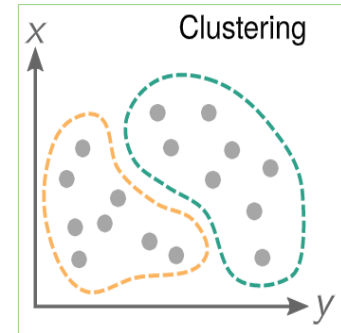
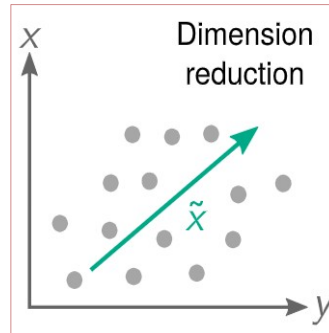
J. Chem. Phys. 154,
230903 (2021)

Unsupervised Machine Learning

Constructing a model from input data without corresponding output labels

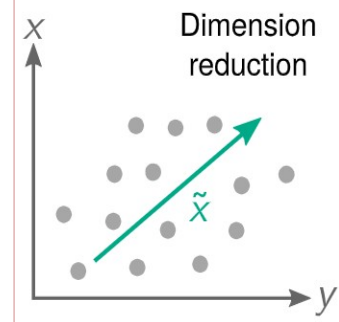
Goal: describe or understand the structure of the data.

- Dimensionality Reduction
- Clustering
- Outlier detection
- Generative Machine Learning



Unsupervised ML: Dimensionality Reduction

Dimensionality Reduction is crucial for identifying the smallest number of features that contain as much information as possible



Principal Component Analysis (PCA)

Start with input vectors \mathbf{x} and build covariance matrix $\mathbf{Q} = \mathbf{x}^T \mathbf{x}$

Covariance matrix tells us how similar each of the inputs x_i are (all pairwise dot products of x_i)

Diagonalize $\mathbf{Q} = \mathbf{W} \mathbf{\Lambda} \mathbf{W}^T$

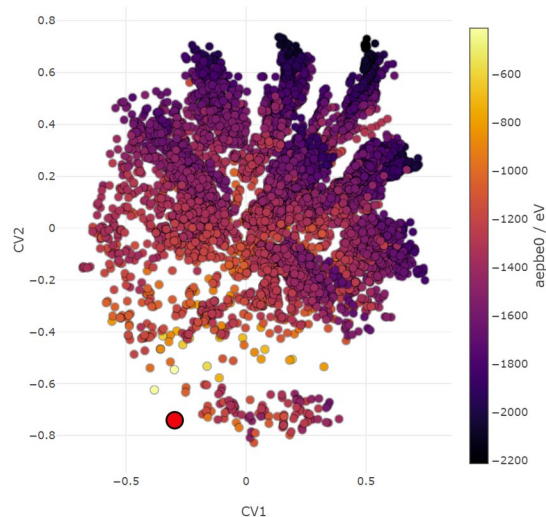
Principal Component eigenvalues $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$

Biggest eigenvalues are responsible for most of the data variance

Principal Component eigenvectors \mathbf{W} tell us how inputs x_i mix

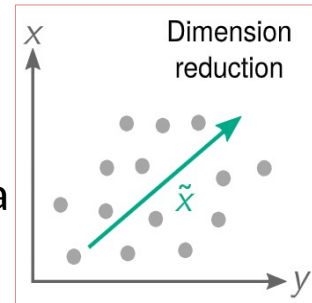
Select a subset of L principal components and transform into lower dimensional space

$$\tilde{\mathbf{X}}_L = \mathbf{X} \mathbf{W}_L$$



Unsupervised ML: Dimensionality Reduction

Dimensionality Reduction can help to recognize and visualize patterns in data

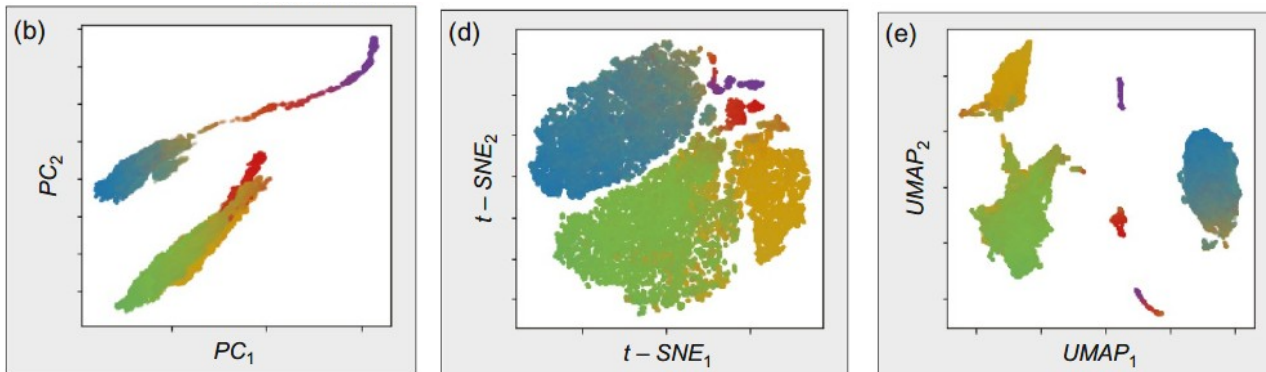


T-SNE: t-distributed stochastic neighbor embedding

t-SNE focuses on preserving the pairwise similarities between data points in a lower-dimensional space

UMAP: Uniform Manifold Approximation and Projection

Similar to t-SNE but uses tricks of topological data analyses to reduce comp. overhead



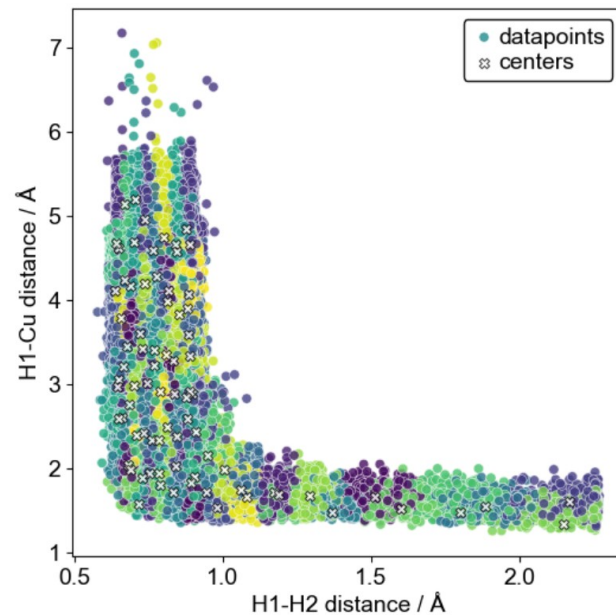
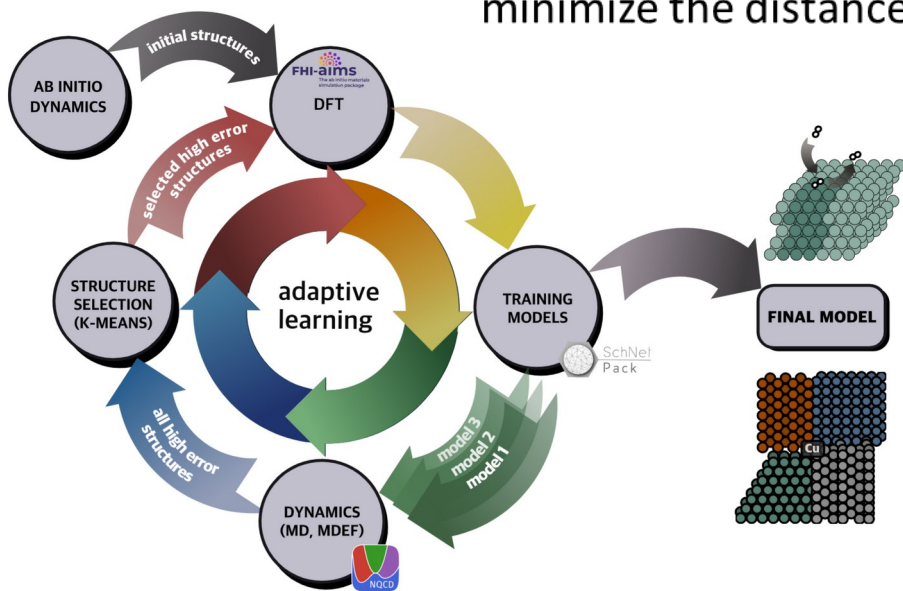
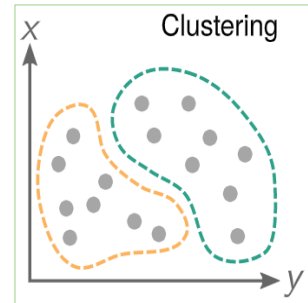
t-SNE preserves small pairwise similarities whereas, PCA maintains large pairwise distances to maximize variance.

Unsupervised ML: Clustering

Example: **K-Means** algorithm

$$\sum_{i=0}^n \min_{\mu_j \in C} (||x_i - \mu_j||^2)$$

- Clusters n data points by separating them into k clusters of equal variance
- Requires number of clusters as input
- The algorithm chooses cluster centroids that minimize the distance to each point in the cluster

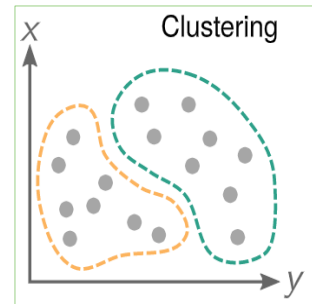


Unsupervised ML: Clustering

Partition-based clustering

vs.

Density-based clustering



Find clusters of equal size

Find variable size clusters around regions of high density

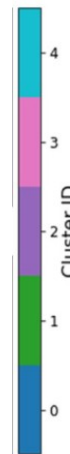
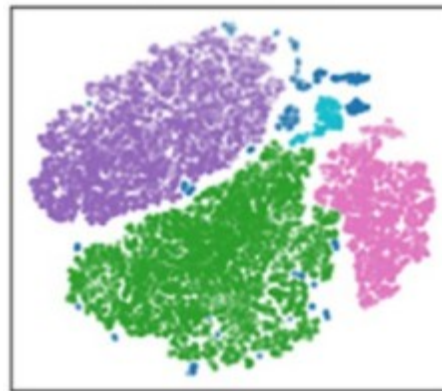
e.g. **K-Means**

e.g. **DBSCAN, HDBSCAN**

K-Means

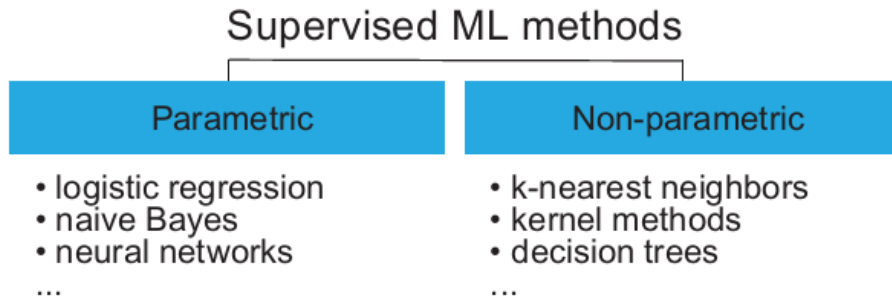


DBSCAN



Supervised Machine Learning

Create model: $y = f(x)$



Discrete output space (classification)
Continuous output space (regression)

Parametric model:

Number of model parameters are independent of number of training datapoints.
Model has a fixed size

Non-parametric model:

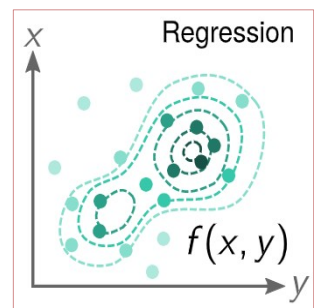
Number of model parameters depend on the number of training data points.

Models we will discuss:

- Multivariate Linear Regression
- Kernel Ridge Regression
- Decision Trees

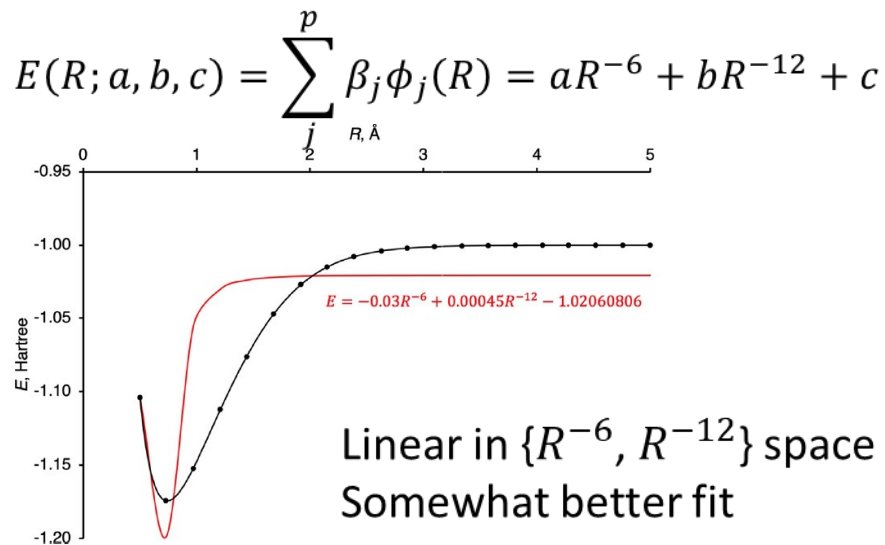
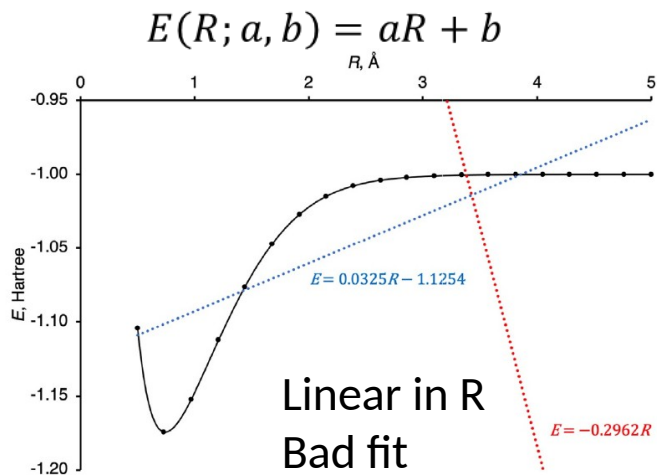
Multivariate Linear Regression

- Linear fit in high dimensional space
- Find a set of regression coefficients β
- Important that input data x_i is represented in a way that shows close to linear relationship with y_i
- Many ways to fit a large set of parameters
- Not a universal estimator so NOT ML

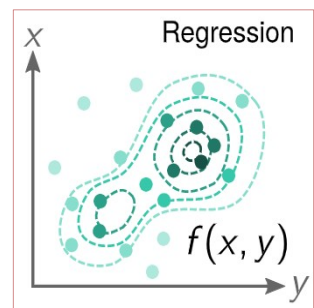


$$f(\mathbf{x}; \beta) = \sum_{j=1}^p \beta_j x_j = \mathbf{x}^T \beta$$

Example: Fit H_2 dissociation curve, $E(R)$



Kernel Ridge Regression



Step 1:
Start from linear regression

$$f(\mathbf{x}; \beta) = \sum_{j=1}^p \beta_j x_j = \mathbf{x}^T \beta$$

Step 2:
Expand coeffs. in high-dimensional space
spanned by training data

$$\beta_j = \sum_{i=1}^{N_{tr}} \alpha_i x_{ij} \longrightarrow f(\mathbf{x}') = \sum_{i=1}^{N_{tr}} \alpha_i \langle \mathbf{x}_i, \mathbf{x}' \rangle$$

Step 3:
Do the same in space of nonlinear basis functions $\{\phi\}$

$$f(\phi(\mathbf{x}')) = \sum_{i=1}^{N_{tr}} \alpha_i k(\mathbf{x}_i, \mathbf{x}')$$

Dot product measures similarity
between data point x_i
and point of prediction x'

This is called the “**kernel trick**”
Many different kernels/basis functions:
Gaussian, Laplacian, polynomial, ...

If $k(\mathbf{x}_i, \mathbf{x}') = \langle \mathbf{x}_i, \mathbf{x}' \rangle$
then we are back to linear regression but in
training data space that can be expanded
-> Universal approximator

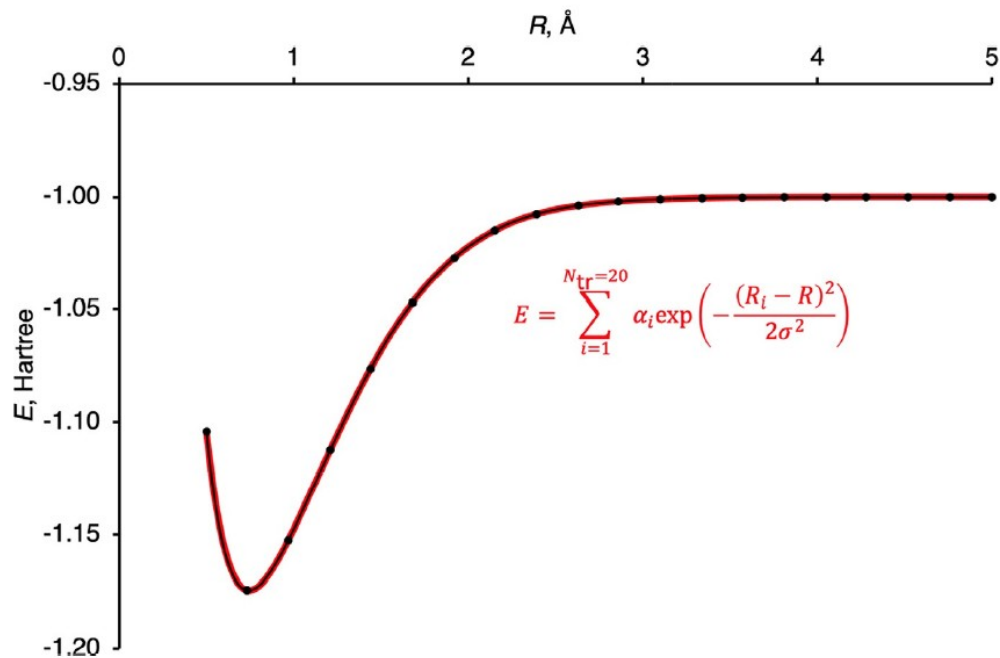
Example: Fit H_2 dissociation curve, $E(R)$ with Kernel Ridge Regression (KRR)

Gaussian kernel

$$k(\mathbf{x}_i, \mathbf{x}') = \exp\left(-\frac{1}{2\sigma^2} \sum_j^p (x_{ij} - x'_j)^2\right)$$

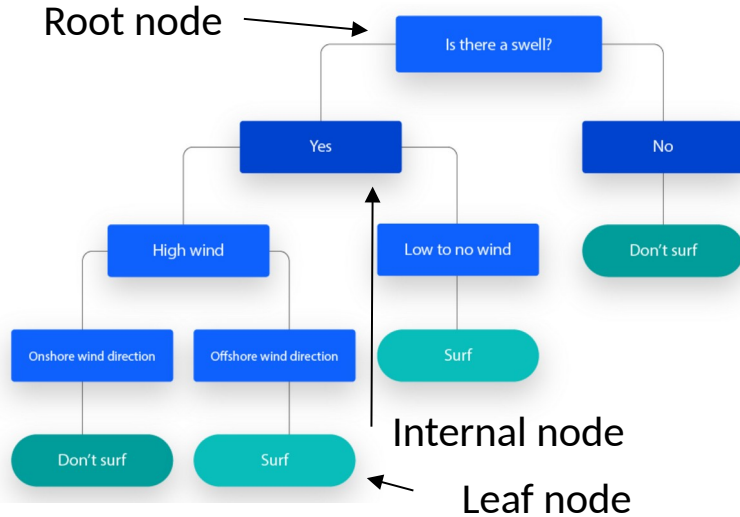
“Gaussian KRR puts Gaussian basis function with width σ on each data point and multiplies it with coefficient α_i .”

We fit coefficients α_i
(least squares fit)



This is now a universal ML method. There is a risk of overfitting so we require regularization during fitting (fitting along the “Ridge” of solutions where coefficients remain as small as possible)

Decision Trees



Internal nodes represent attribute tests,

Branches represent attribute values

Leaf nodes represent final decisions or predictions

Start with dataset of (\mathbf{X}, Y) where Y is the label (surf, don't surf) and \mathbf{X} is a set of attributes

swell (Y/N),
wind (numerical)
Wind direction (onshore/offshore)

What are trees

- ✓ Non-parametric supervised learning methods for classification and regression
- ✓ Simple to understand and to visualize
- ✓ Can handle numerical and categorical data
- Predictions are non-smooth
- Large trees can be unstable -> overfitting
- Tackle via ensembles of trees -> Random Forests

Decision Trees - Classification

How are trees built?

- Finding splits (splits represent decision nodes or attribute tests)
- Identify attribute tests that maximise the “Information Gain” or minimise the “Gini Index”

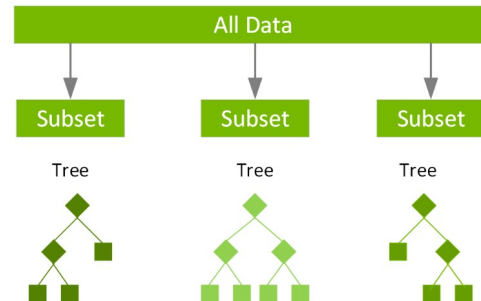
Gini Index: measures how often a randomly chosen data point would be incorrectly identified by a certain attribute test.

Attributes with lower Gini indices are preferred as they better split the data

$$\text{Gini} = 1 - \sum_{i=1}^n p_i^2 \quad p_i^2 \text{ probability of a certain outcome } i$$

Ensembles of trees provide improved generalizability and robustness

- Gradient-boosted trees (e.g. XGBoost)
- Random forests



Reinforcement Learning

- **Agent** performs certain **actions** in an **environment** at each time step in a sequential decision-making framework
 - Actions change the **state** and can provide positive or negative feedback -> goal is to learn a **policy** that provides maximum **reward**
- ▶ Uses rewards instead of labels to learn



- Temporal Credit Assignment Problem: associating a reward with an action
- Finding Trade-off between Exploitation vs. Exploration

Examples

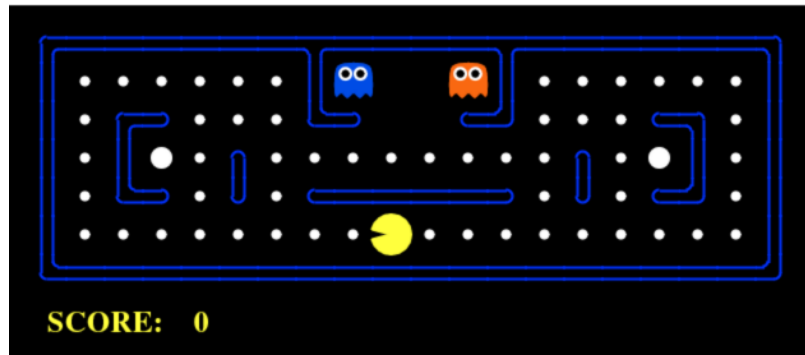
- ▶ Video games
- ▶ Training robots
- ▶ AlphaGo

Markov process

Sequence of states $s_1, s_2, s_3, \dots, s_t, \dots$

Transition probability $t \rightarrow t + 1$

$$Pr(s_{t+1}|s_t)$$



Markov reward process

Sequence of rewards $r_1, r_2, r_3, \dots, r_t, \dots$ discount factor $\gamma^k \in (0,1]$

$$\text{return } G_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

Agent performs actions a_t

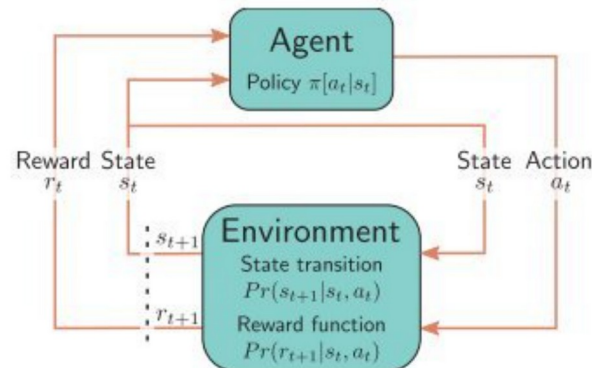
Markov decision process

Sequence of states and actions $s_1, a_1, s_2, a_2, s_3, a_3, \dots, s_t, a_t, \dots$

Actions affect transition probability $Pr(s_{t+1}|s_t, a_t)$

Actions affect reward probability $Pr(r_{t+1}|s_t, a_t)$

Policy: $\pi[a|s]$ determines the action, stochastic or deterministic, stationary or time-dependent



Assign value functions to states, and actions to determine optimal reward G_t

State value function

$$v[s_t|\pi]$$

Expected return for being in state t

Action value function

$$q[s_t, a_t|\pi]$$

Expected return for executing action in state t

If we know the optimal action values, we can derive the optimal policy $\pi[a_t|s_t]$

$$v[s_t] = \sum_{a_t} \pi[a_t|s_t] q[s_t, a_t]$$



Bellman equations to
define optimal policy

$$q[s_t, a_t] = r[s_t, a_t] + \gamma \cdot \sum_{s_{t+1}} Pr(s_{t+1}|s_t, a_t) v[s_{t+1}]$$

Types of RL

- Tabular reinforcement learning (methods that don't rely on function approximation/ML)
- Fitted Q-learning (action value function replaced by machine learning model)
- Policy gradient methods (directly learn a stochastic policy $\pi[a_t|s_t]$)
- and many more...

Assessing a Model

Uncertainty Quantification

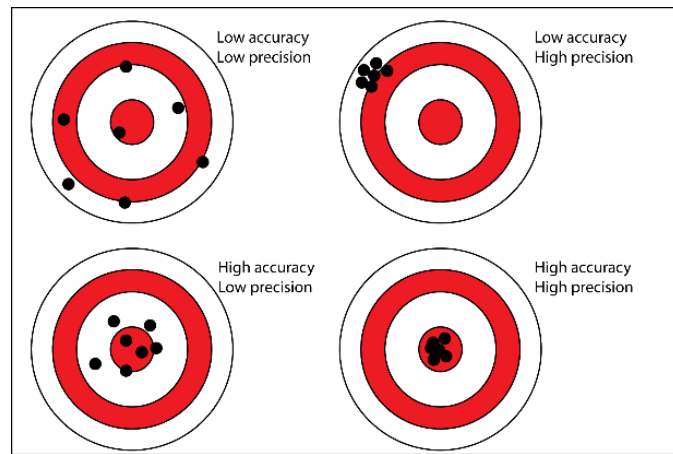
Accuracy of prediction: MAE or RMSE

Precision of prediction: ?????

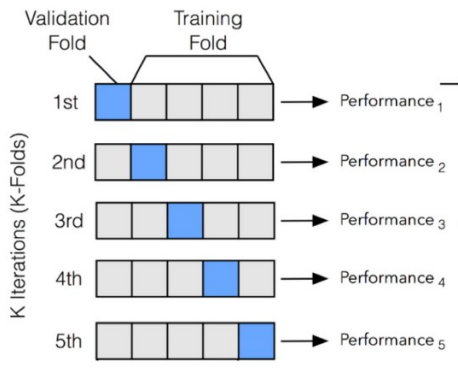
2 Sources of uncertainty in prediction:

- Aleatoric uncertainty (statistical error, noise in data)
- Epistemic uncertainty (intrinsic to model)

How to calculate uncertainty / standard deviation?

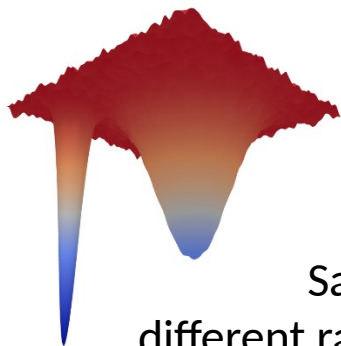


Bootstrapping
(e.g. subsampling)



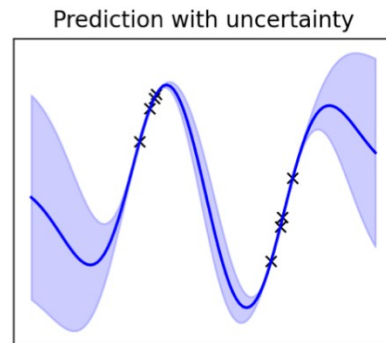
Mean
&
StDev

Ensemble Learning
(stacking, bagging, boosting) (e.g. Gaussian Process Regression)



e.g.
Random Ensembles
Same data,
Same model,
different random seed

Bayesian Uncertainty



Validating and Testing

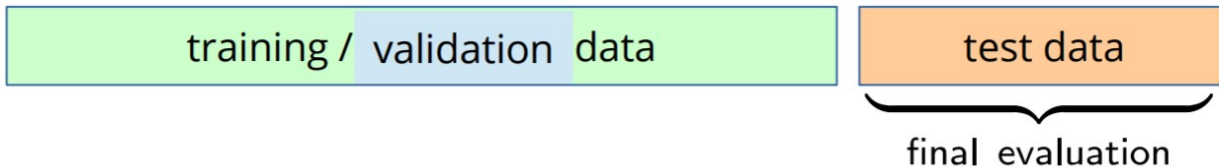
Training is the process where a model analyzes the data to learn underlying patterns and adjust its internal weights to minimize the difference between its predictions and the actual results.

Validation is the process of evaluating a model during development on a subset of data to tune hyperparameters and select the best-performing version of the algorithm and avoid overfitting.

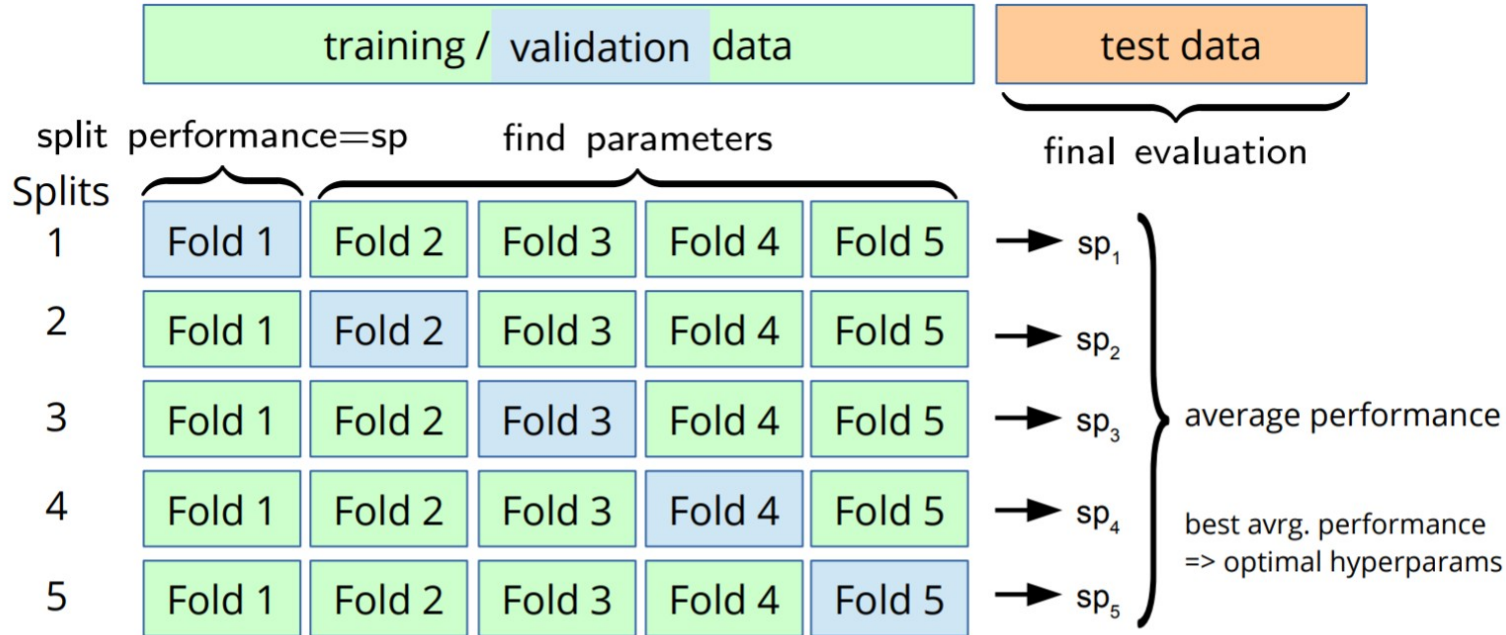
Testing is the final assessment of a completed model using a completely independent dataset to provide an unbiased measurement of how it will perform in the real world.

The Data Split

- **Training Set (60-80%):** The model sees this data during the learning to adjust its parameters.
- **Validation Set (10-20%):** Used to assess how the model handles data it wasn't trained on, allowing tweaking hyperparameters.
- **Test Set (10-20%):** Used to assess the model at the very end.



Validating your model: K-fold Cross Validation



K-fold CV

- used for model validation (calculate accuracy and standard deviation of prediction)
- used for Hyperparameter optimisation (Grid search, Random search)
- Avoids overfitting
- Can help to identify outliers and unbalanced datasets

Putting it all together

Typical workflow in ML project

- Define the task and the objective (informing the loss func. and data gen.)
- Generate and clean the data (e.g. find patterns, find outliers) **Clustering**
- Discover and visualize the data to gain insights (find correlations) **Classification**
- Prepare data for training (e.g. train/test split, scaling) **Dimensionality Reduction**
- Find good data representation: “Featurization”
- Select, train, and evaluate model (e.g. calculate MAE, RMSE) **e.g. Classification or Regression**
- Optimise and fine-tune model with cross-validation
- Assess accuracy and uncertainty of trained model
- Generate more data to improve accuracy/uncertainty (e.g. active learning)
- When ready, deploy model (£££ and/or manuscript) **Clustering**

Research Example: Generative molecular design

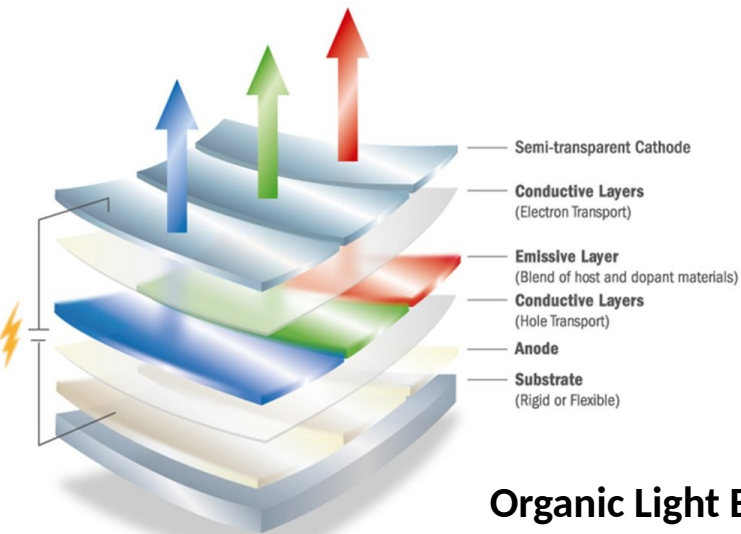
Westermayr, Maurer, Chem. Sci. 12, 10755 (2021)

Westermayr et al., Nature Computational Science 3, 139-148 (2023)

Koczor-Benda et al, arXiv: 2503.14748

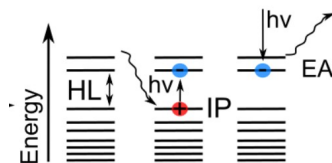
Koczor-Benda et al., arXiv: 2503.21328

Example: Designing molecules with tailormade properties?

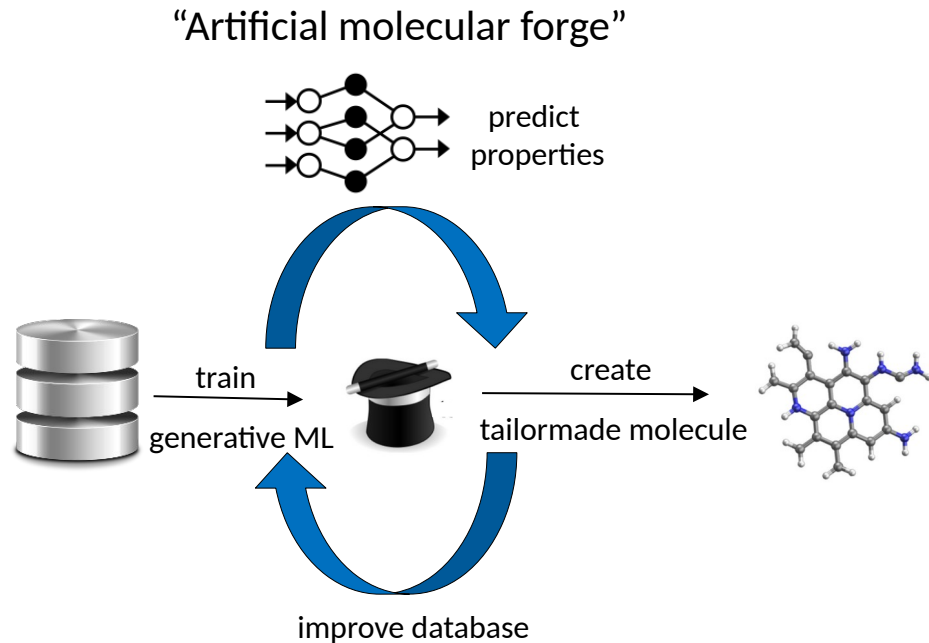


Organic Light Emitting Diodes (OLEDs)

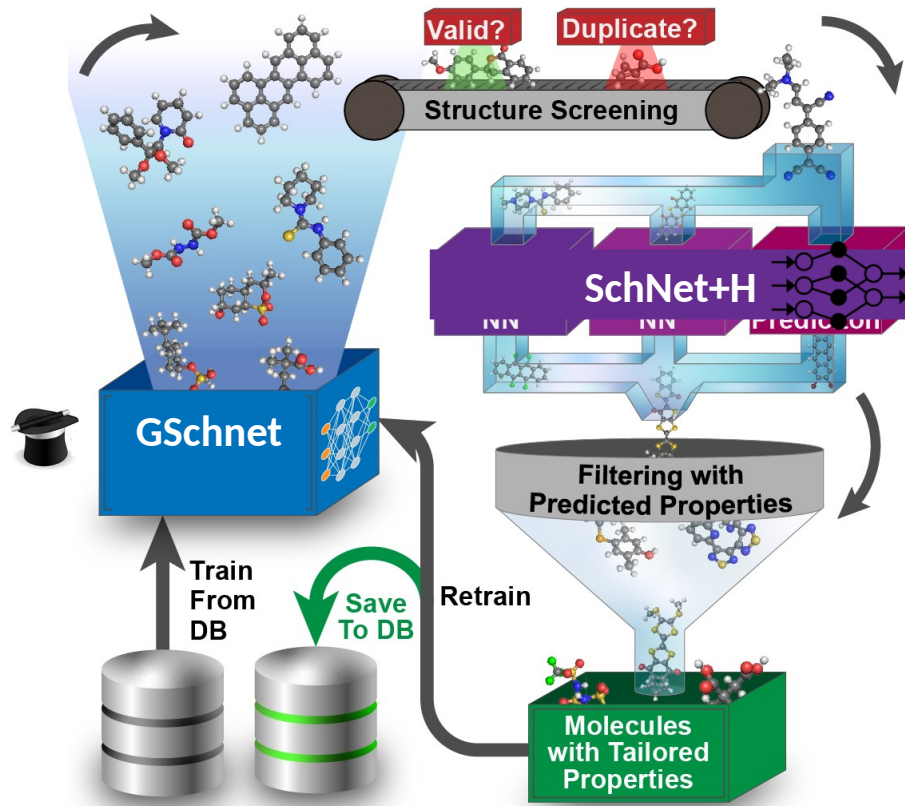
**Application:
Generative molecular design
of organic electronics**



IP ... Ionization potential
EA ... Electron affinity
HL ... HOMO-LUMO gap

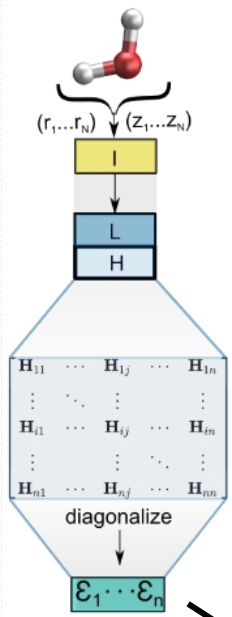


Generative Design of Molecules with Tailored Properties



Electronic structure ML: one model for many molecules

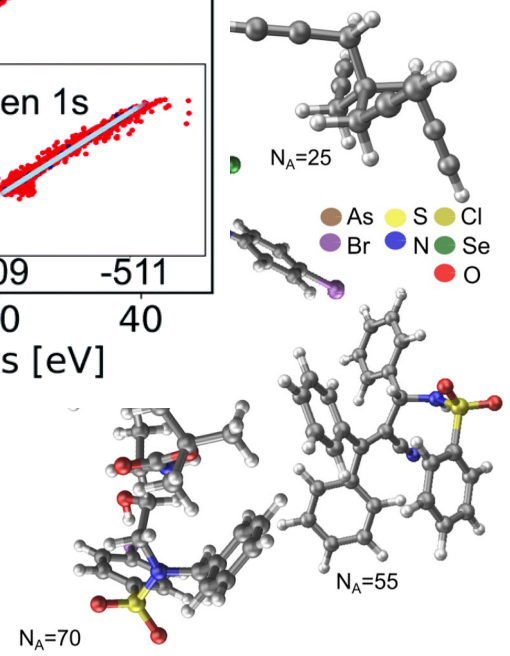
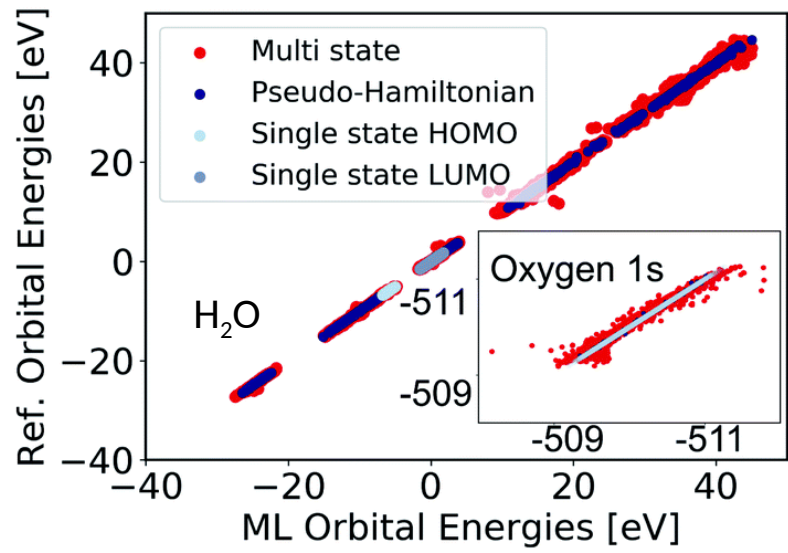
pseudo-Hamiltonian approach: SchNet+H



SchNet + H
Model builds an internal Hamiltonian representation

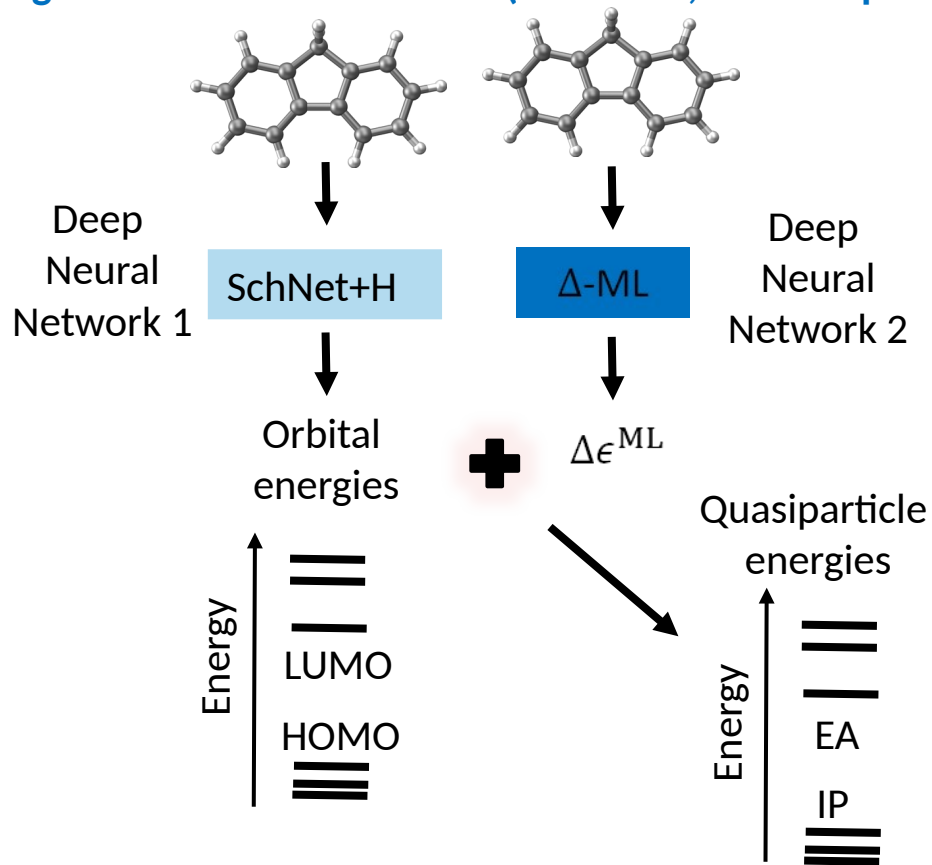
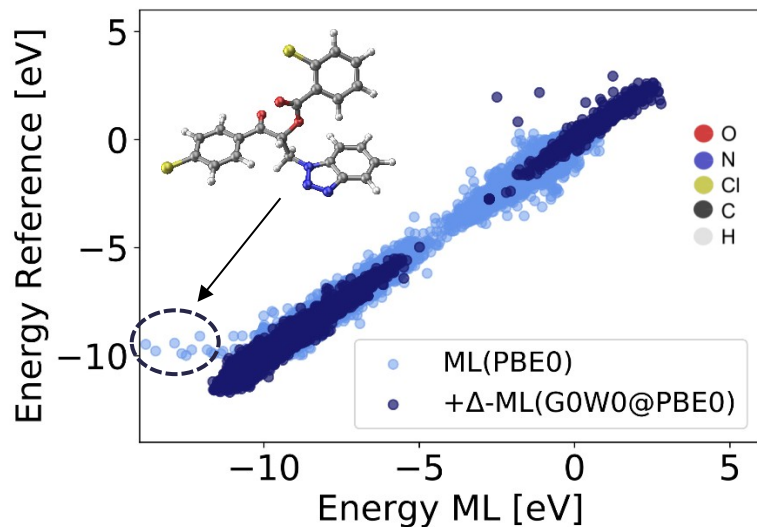
- I Input: Embedding of SchNet
- L Layer(s) to map the input to output(s)
- H Machine learning pseudo-Hamiltonian
- E Output layer: n Eigenvalues ϵ

predicts a vector of N electronic energy levels

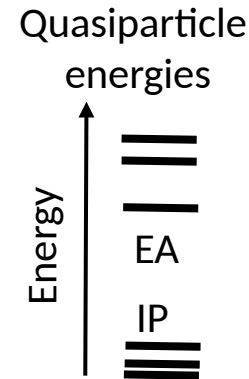
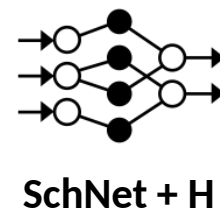
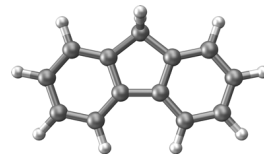
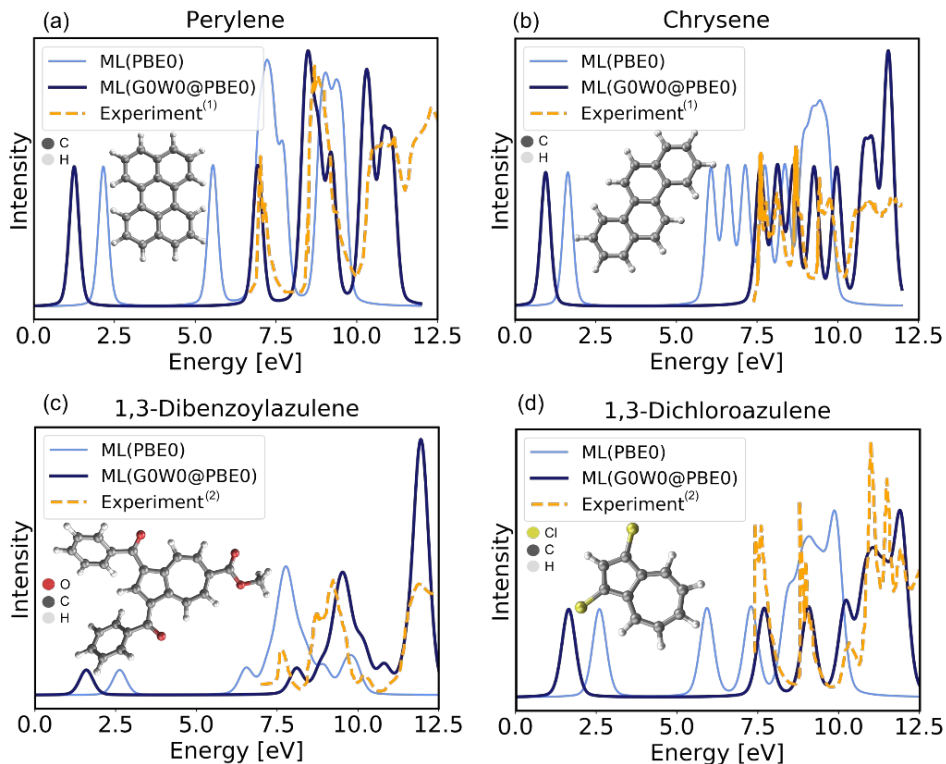


Step 1: SchNet+H predicts DFT orbital energies (train on 50,000 data points)

Step 2: Delta-ML provides accurate ionisation energies and electron affinities (train on 4,000 data points)



ML prediction of molecular photoemission spectra



Potential applications for ML model:

- (Inverse) UV/Vis Photoemission Spectra
- X-ray Photoemission Spectra (XPS)
- **High-Throughput Screening**
- **Inverse Chemical Design**

Exp: (1) Dougherty et al., *J. Electron Spectrosc. Relat. Phenom.* 19, 21-33 (1980)

(2) Deleuze, *J. Chem. Phys.* 116, 7012-7026 (2002)

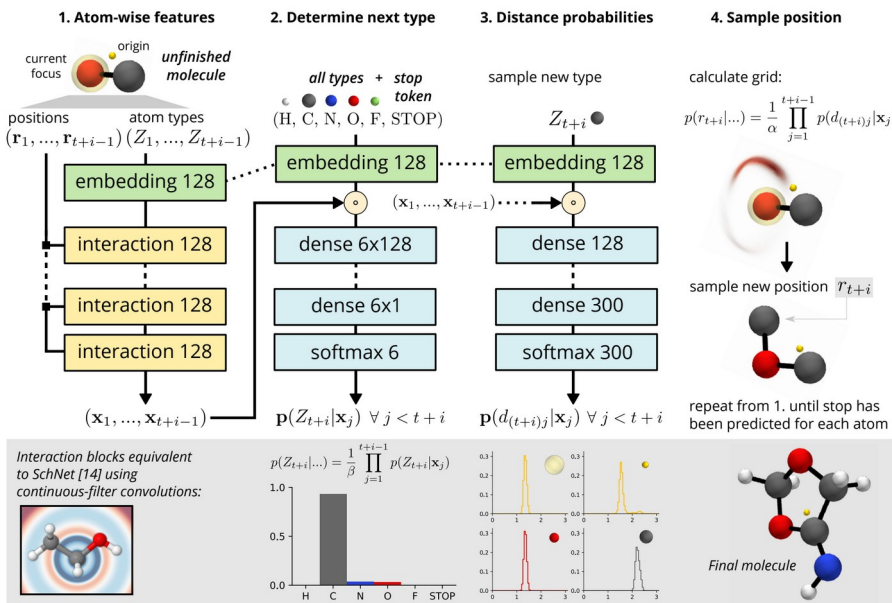
(3) Rangel et al., *Phys. Rev B* 93, 115206 (2016)

Generative deep learning of 3D molecular structures

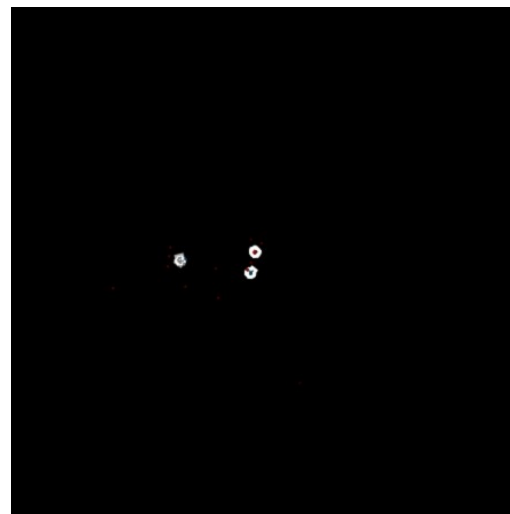
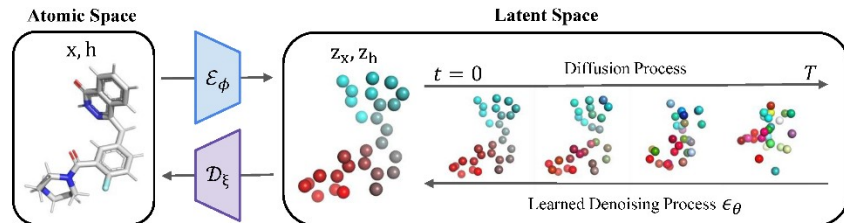


Autoregressive atom-by-atom construction of molecules

GSchnet



Geometric diffusion models



Model trained
by MChem student
Abudalla Al-Fekaiki

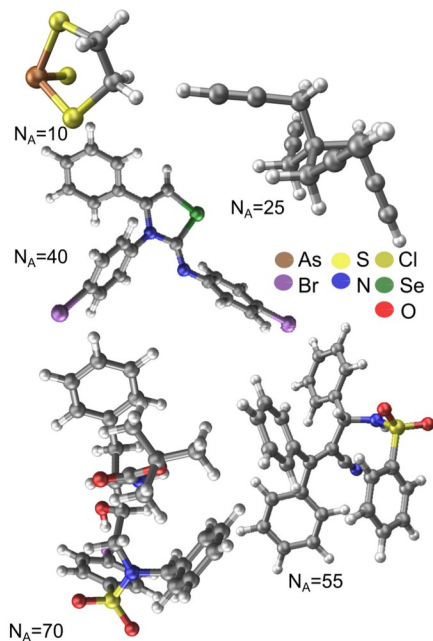
GeoLDM model:
Xu et al. ICML (2023)

Generative deep learning of 3D molecular structures

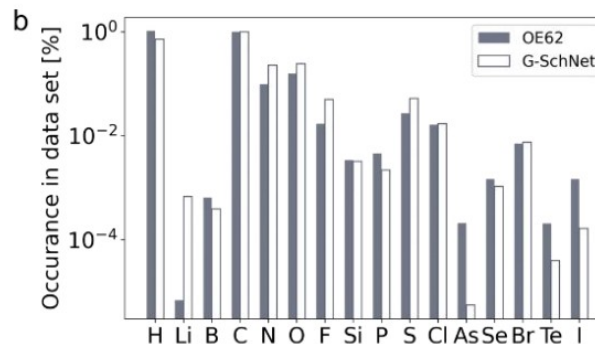


GSchnet

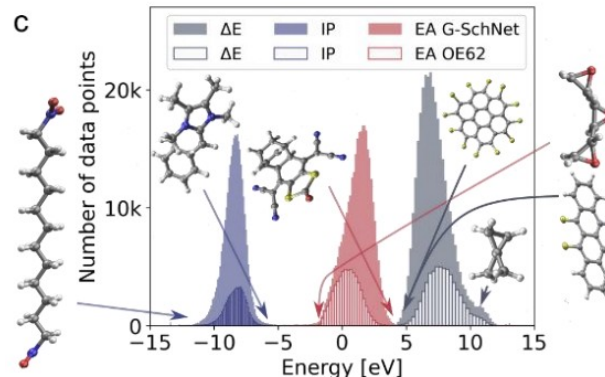
trained on
OE62 database



Generated structures

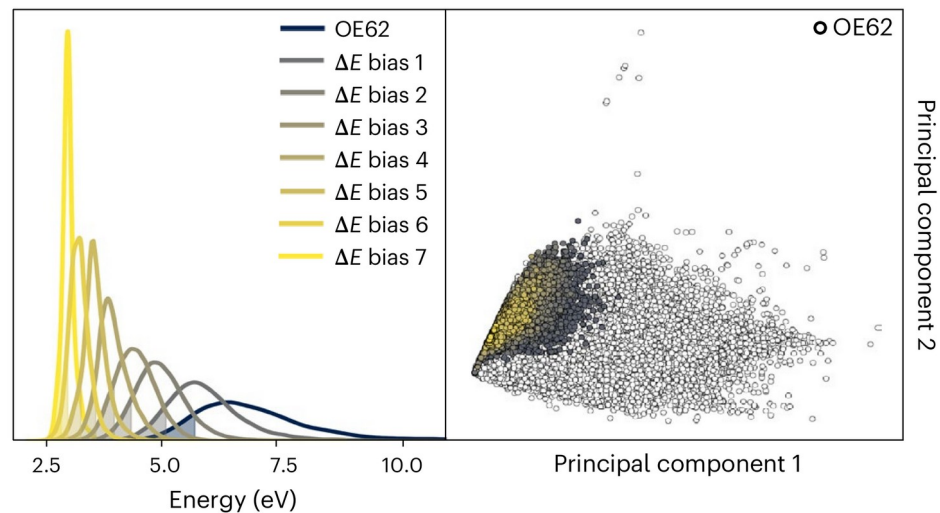
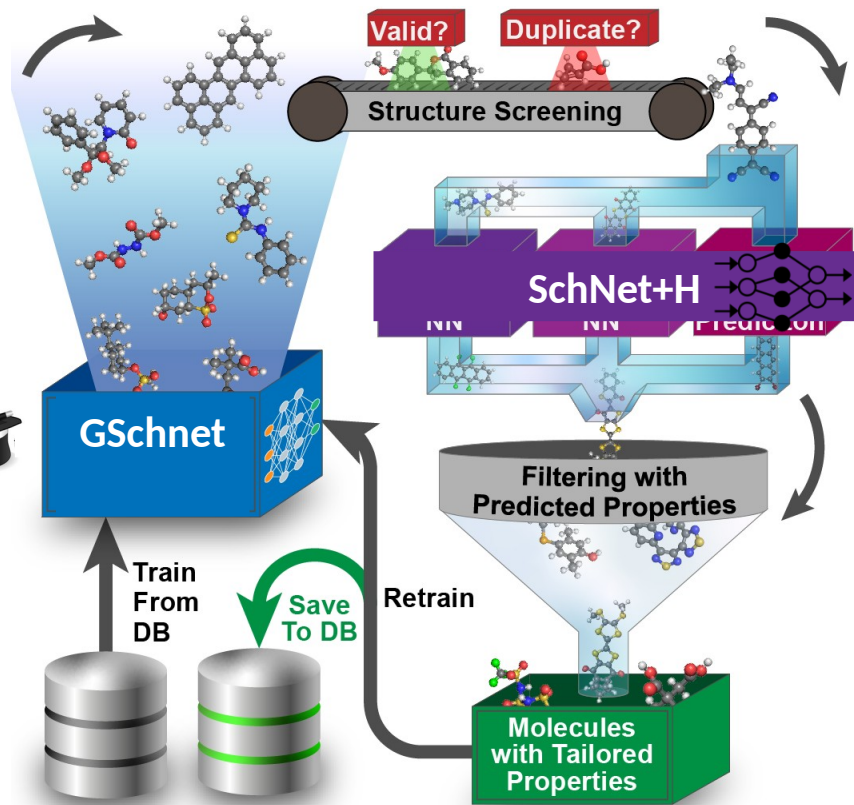


reproduce
elemental
distribution

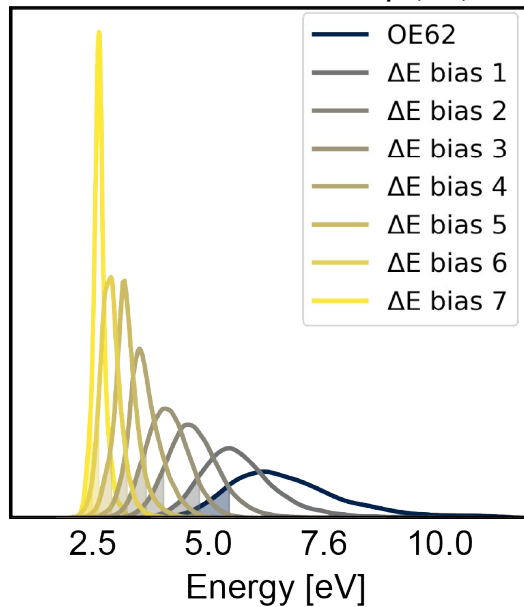


reproduce
property
distribution

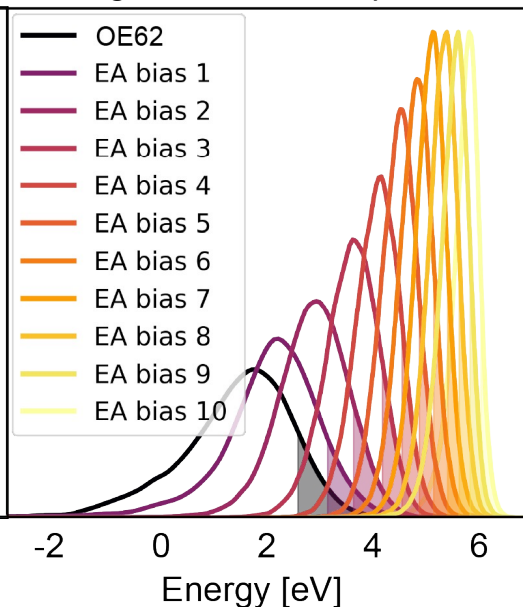
Generative Design of Molecules with Tailored Properties



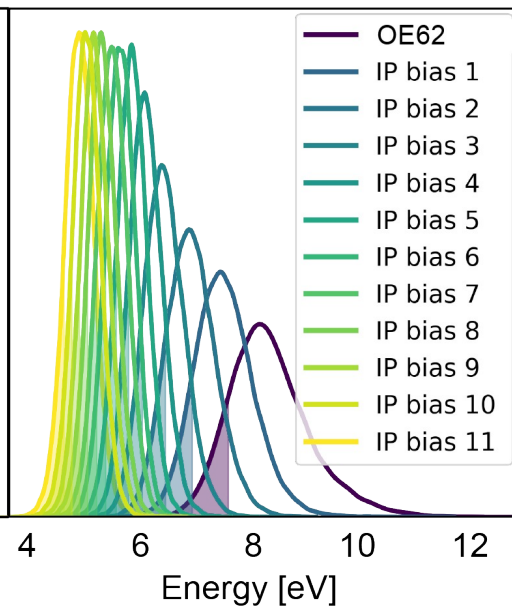
Small Fundamental Gap (ΔE)



High Electron Affinity (EA)

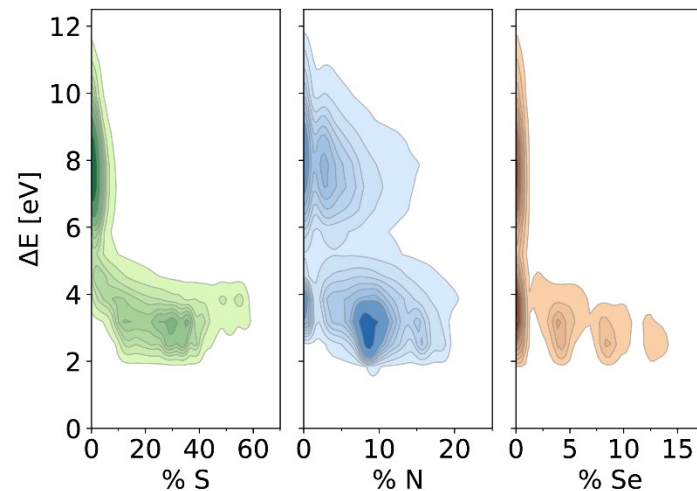
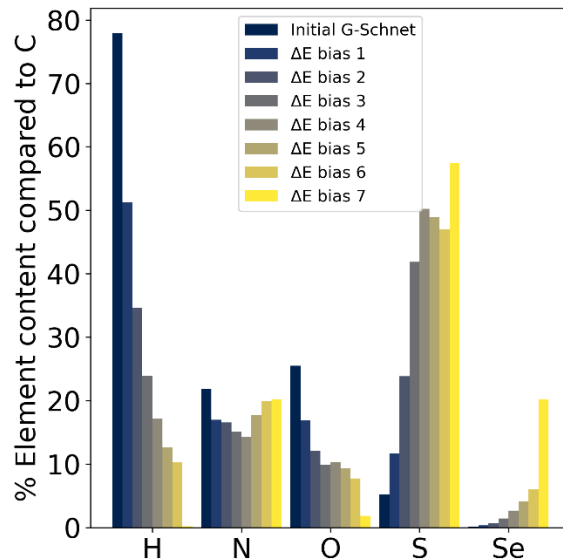
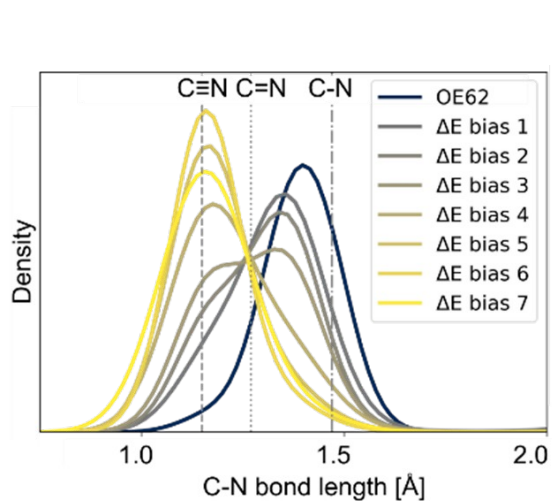
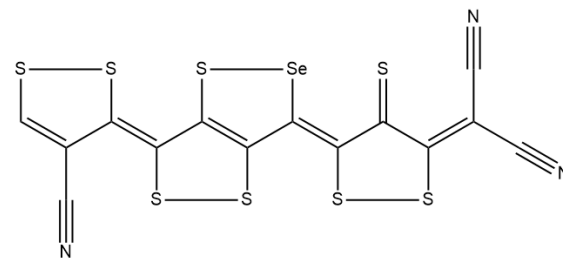


Low Ionisation Potential (IP)

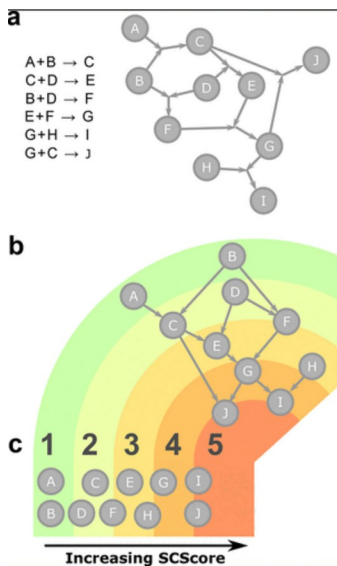


Bonding Descriptor Trends

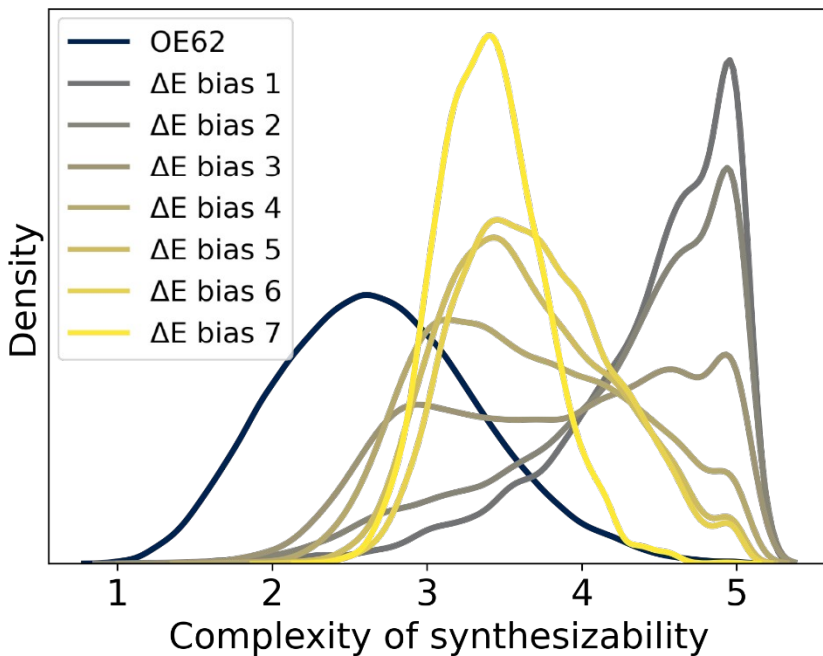
Small Fundamental Gap (ΔE) Biasing



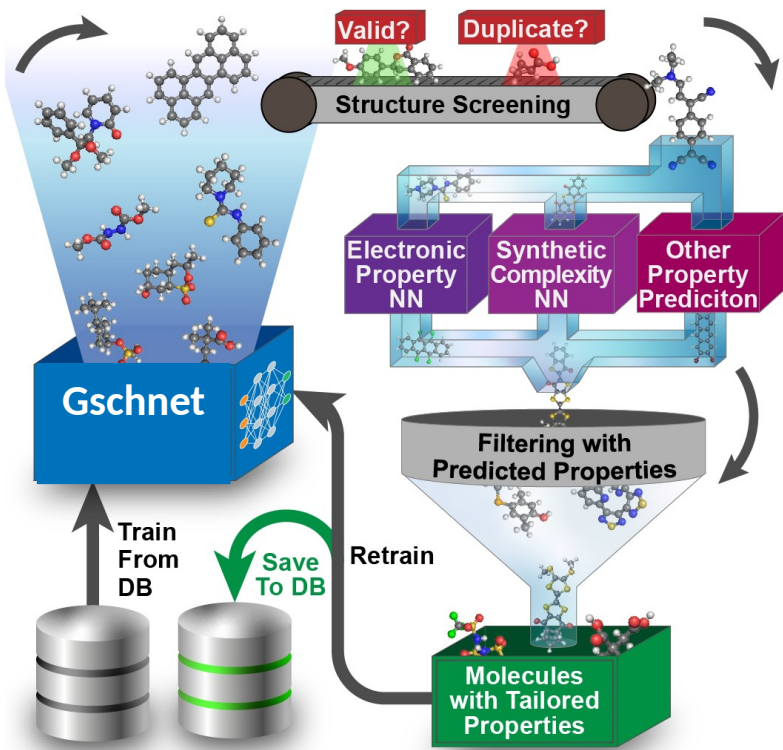
Synthetic Viability



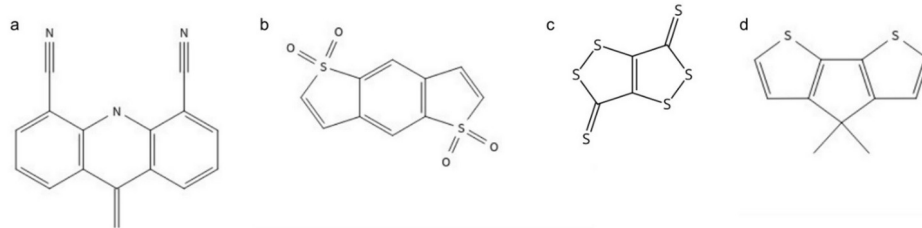
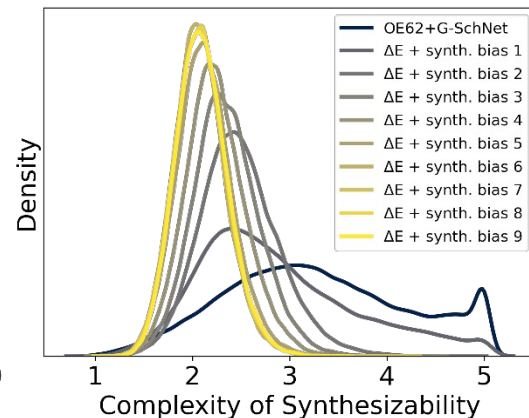
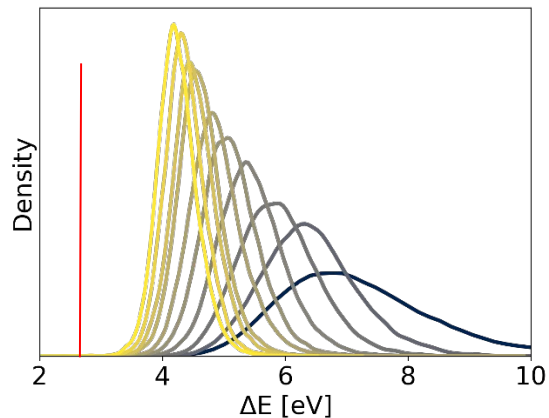
- High selenium content leads to molecules that are difficult to synthesise.
- We quantify this with the SCScore metric.¹



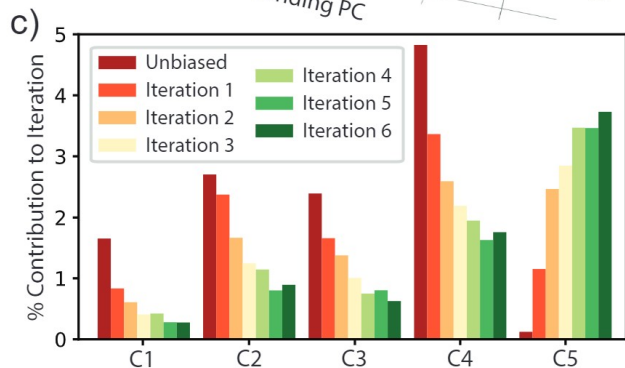
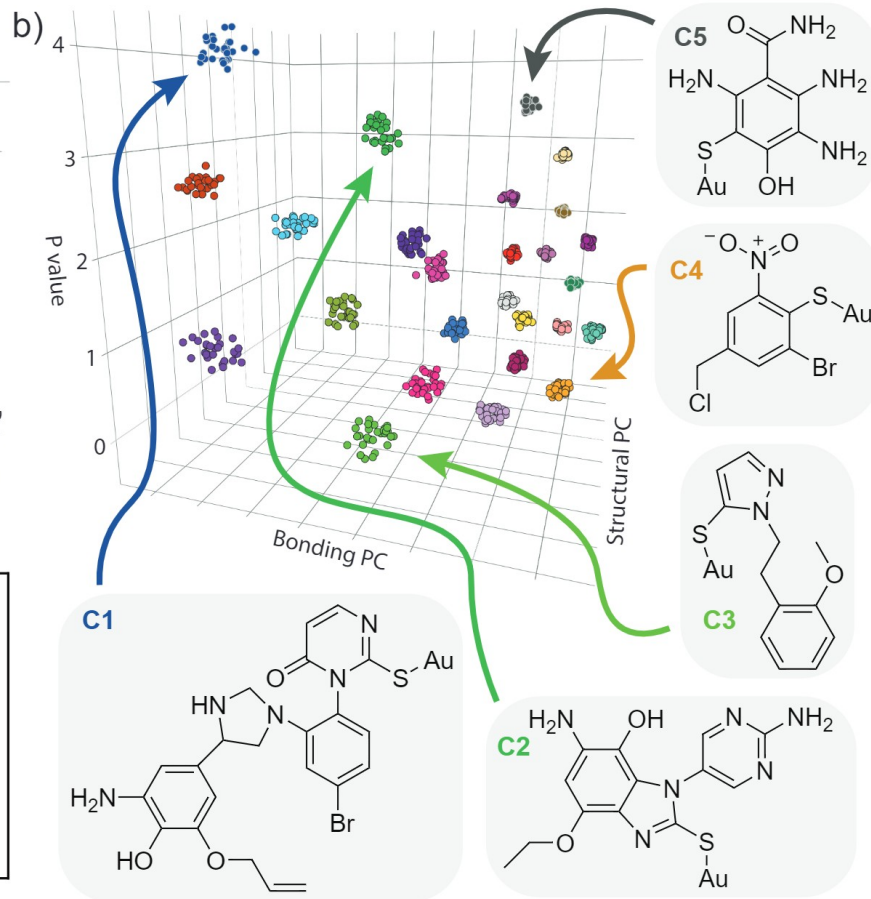
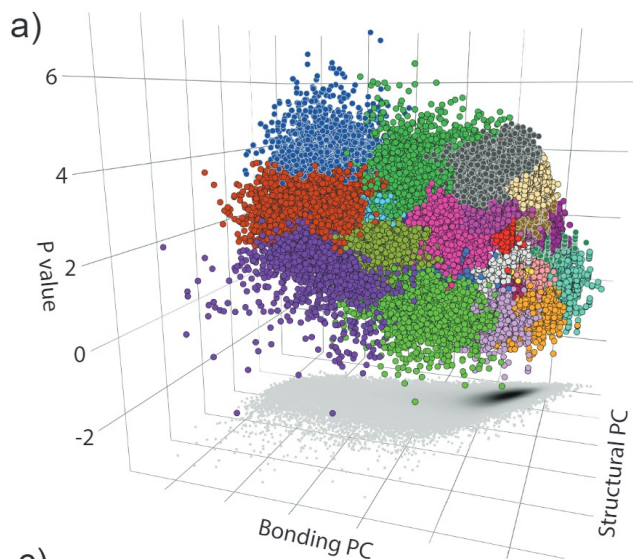
Generative Design of Molecules with Multi-Property Optimisation



Combine electronic screening and synthetic complexity screening



Interpreting the Data: Clustering



Thank you!

Go and use ML methods for your research!

BUT PLEASE

- Remember: learning \neq understanding
- Embrace reproducibility (clear workflows, write tutorials)
- Embrace openness (publish your models, data & scripts!)