# Machine Learning and Quantum Alchemy

Guido Falk von Rudorff, University of Kassel

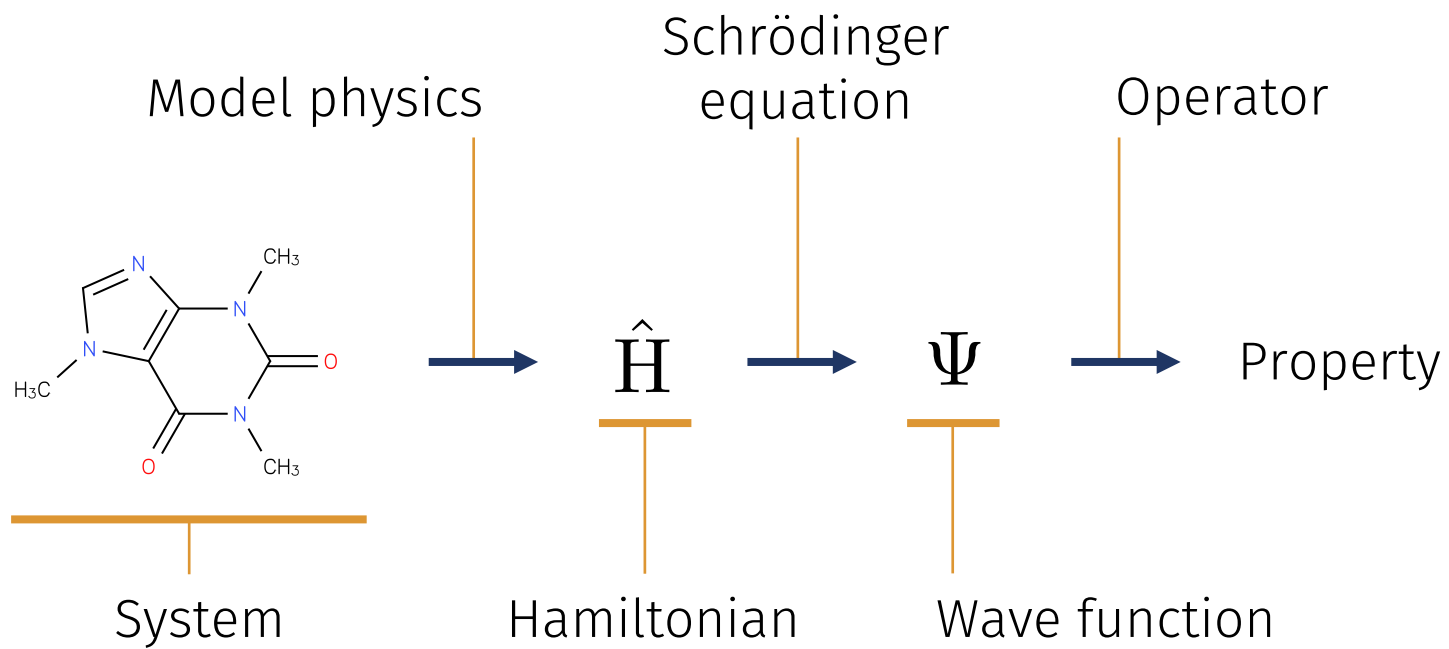✉ vonrudorff@uni-kassel.de     🌐 nablachem.org/talks     ⓖ ferchault     𝕏 @ferchault
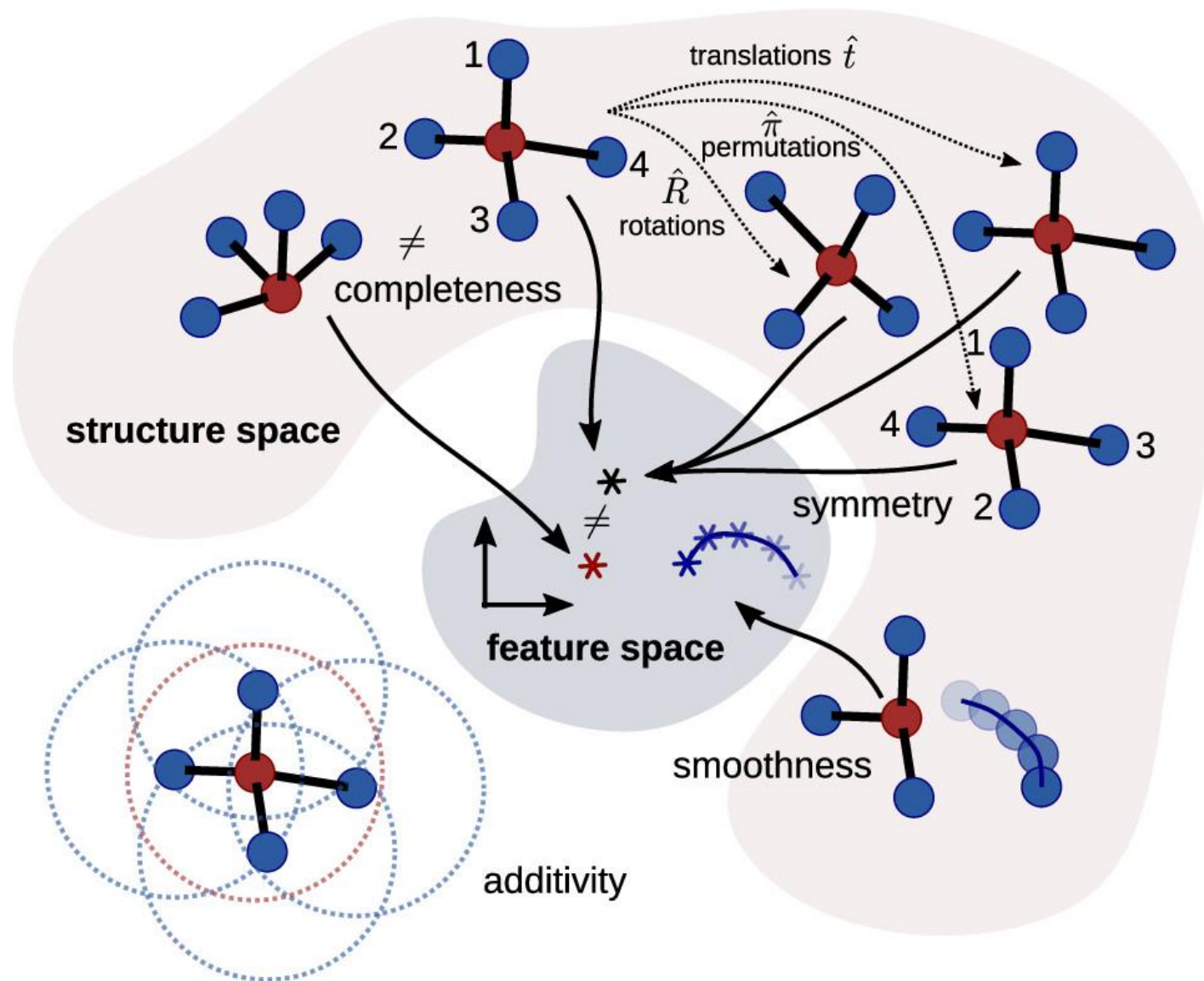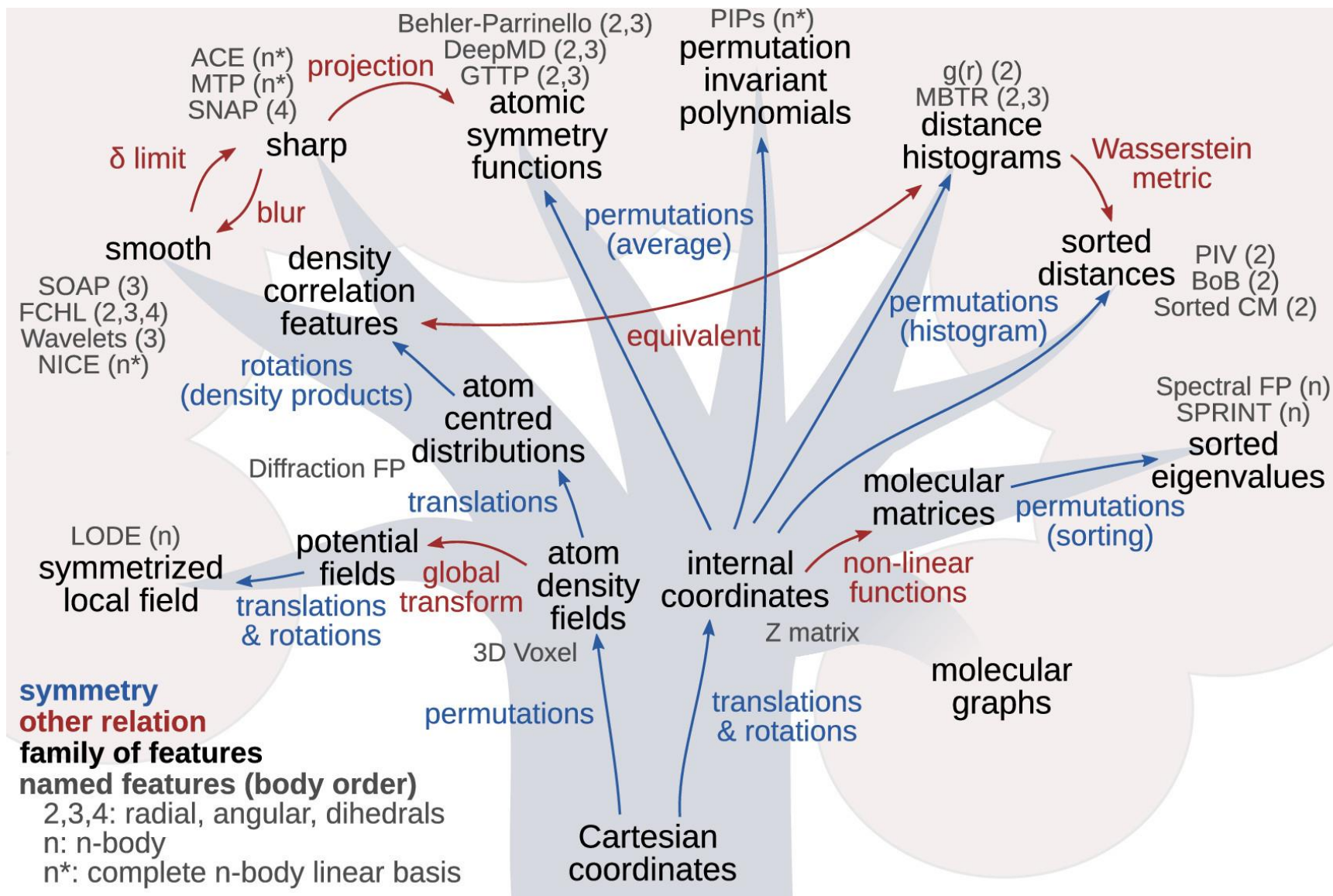
Model physics

Schrödinger
equation

Operator

$\hat{H}$ → $\Psi$ → Property

System

Hamiltonian

Wave function

# Representations
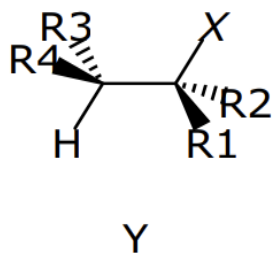
Categorial data vs regression

Solution: „binary", „dummy", „one-hot" encoding

**Encode $n$ categories as vector of length $n-1$ with one category (arbitrary) being the null vector.**

Example:
- A: (0, 0)
- B: (1, 0)
- C: (0, 1)

Chemistry example: (ABBA|CD) -> (0,0,0,0,1,0,0,0,1,0,0,0,0,0,0,0|0,1,0,0,1)



|      | A   | B      | C   | D      | E      |
|------|-----|--------|-----|--------|--------|
| **R$k$** | H   | NO$_2$ | CN  | CH$_3$ | NH$_2$ |
| **X**    | F   | Cl     | Br  |        |        |
| **Y**    | H   | F      | Cl  | Br     |        |

Often: molecules = well-defined bonds

$$O = C = O$$

## Adjacency matrix
- 1 if atoms i and j are bonded
- 0 otherwise

$$\begin{matrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{matrix}$$

## Bond order matrix
- Bond order if atoms i and j are bonded
- 0 otherwise

$$\begin{matrix} 0 & 2 & 0 \\ 2 & 0 & 2 \\ 0 & 2 & 0 \end{matrix}$$

## SMILES
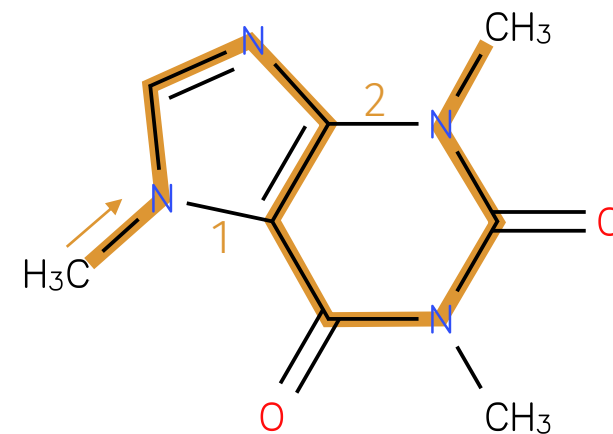- Bonds: Nothing (1), = (2), # (3), $ (4)     O=C=O
- Partial charges                             [Na+]
- Fragments: .                                [Na+].[Cl-]
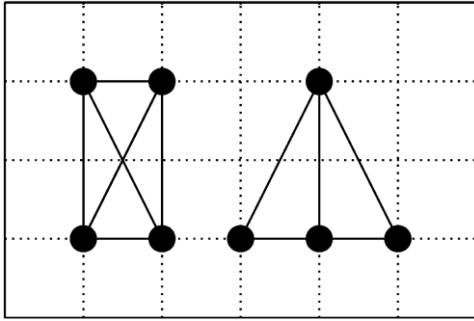- Rings: labels                               C1CCCC1
- Branches: parentheses

CN1C=NC2=C1C(=O)N(C)C(=O)N2C

Many-body descriptions, e.g. all pairwise distances



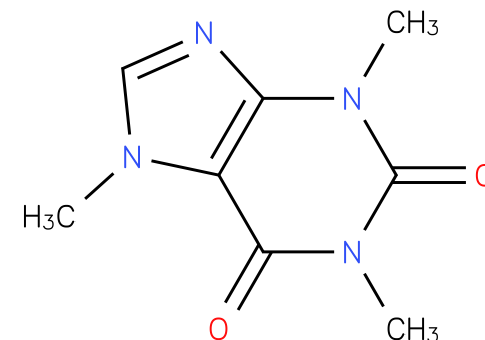Even three- and four-body interactions not unique [1].

Chemistry: **Coulomb Matrix [2]**
- Diagonal:
- Off-diagonal:
- Problem: sorting, uniqueness

$$0.5 Z_i^{2.4}$$
$$Z_i Z_j / \|\mathbf{R}_i - \mathbf{R}_j\|$$

[1] Pozdnyakov et al., *Phys Rev Lett* 2020.  [2] M. Rupp et al., *Phys Rev Lett* 2012.

Collect relevant features of a molecule into a vector.
- Fixed list of functional groups (e.g. Joback method)
- Hash of generated circular atom environments (e.g. ECFP)
- Kinds of generated local environments (e.g. Morgan)
  - Element, # heavy neighbors, # protons, charge, part of ring...
- Fixed checklist of features (e.g. MACCS keys)
  - Subgraphs
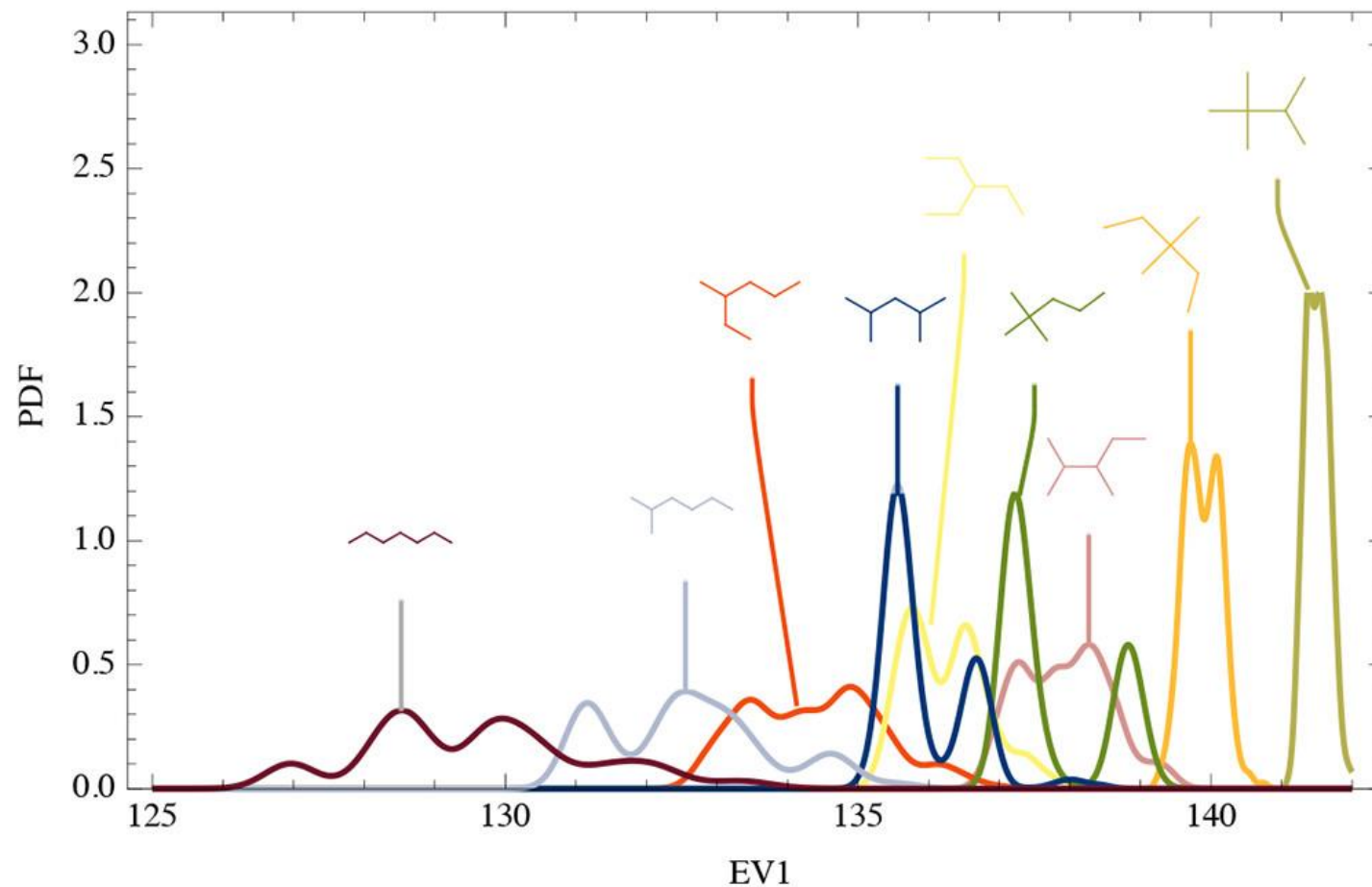
**Morgan** (feature: count)
10565946: 2, 348155210: 1, 476388586: 1, 540046244: 1, 553412256: 1, 864942730: 2, 909857231: 1, 1100037548: 1, 1333761024: 1, 1512818157: 1, 1981181107: 1, 2030573601: 1, 2041434490: 1, 2092489639: 3, 2246728737: 3, 2370996728: 1, 2877515035: 1, 2971716993: 1, 2975126068: 2, 3140581776: 1, 3217380708: 4, 3218693969: 1, 3462333187: 1, 3657471097: 3, 3796970912: 1
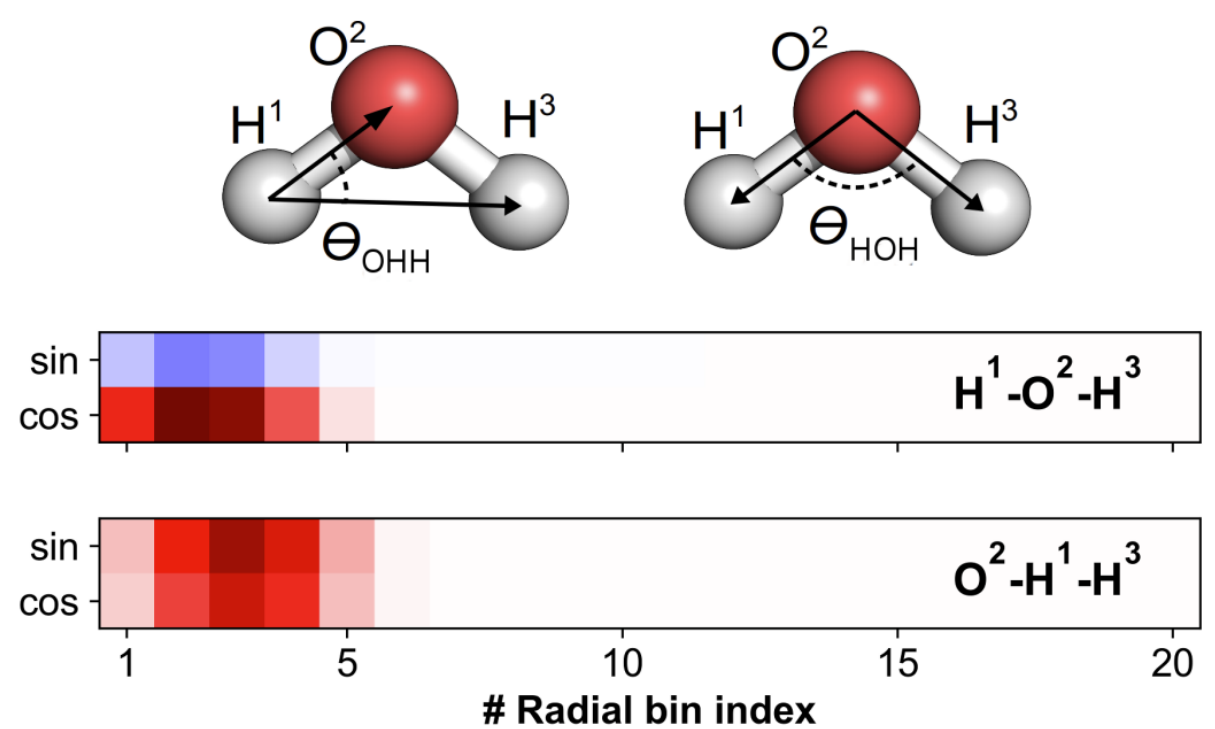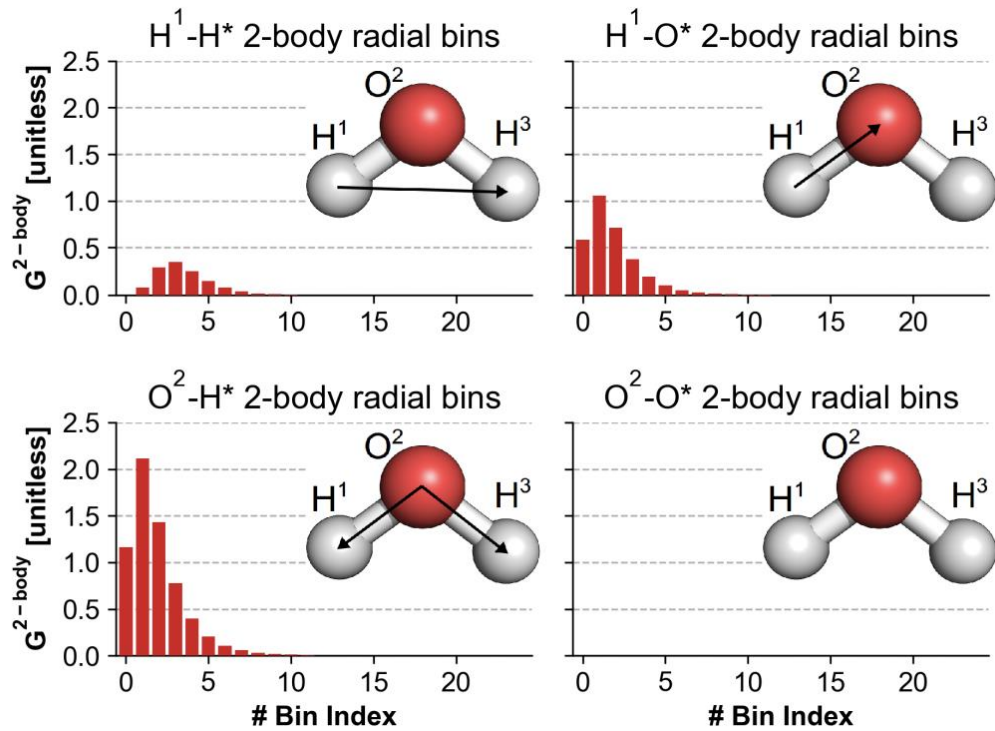
**MACCS keys:**
[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, [0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0]

Condense a matrix into its eigenvalues
- No sorting issue, as permutationally invariant
- Smaller: N instead of $N^2$
- Lossy

Smear positions into densities (FCHL / SOAP)
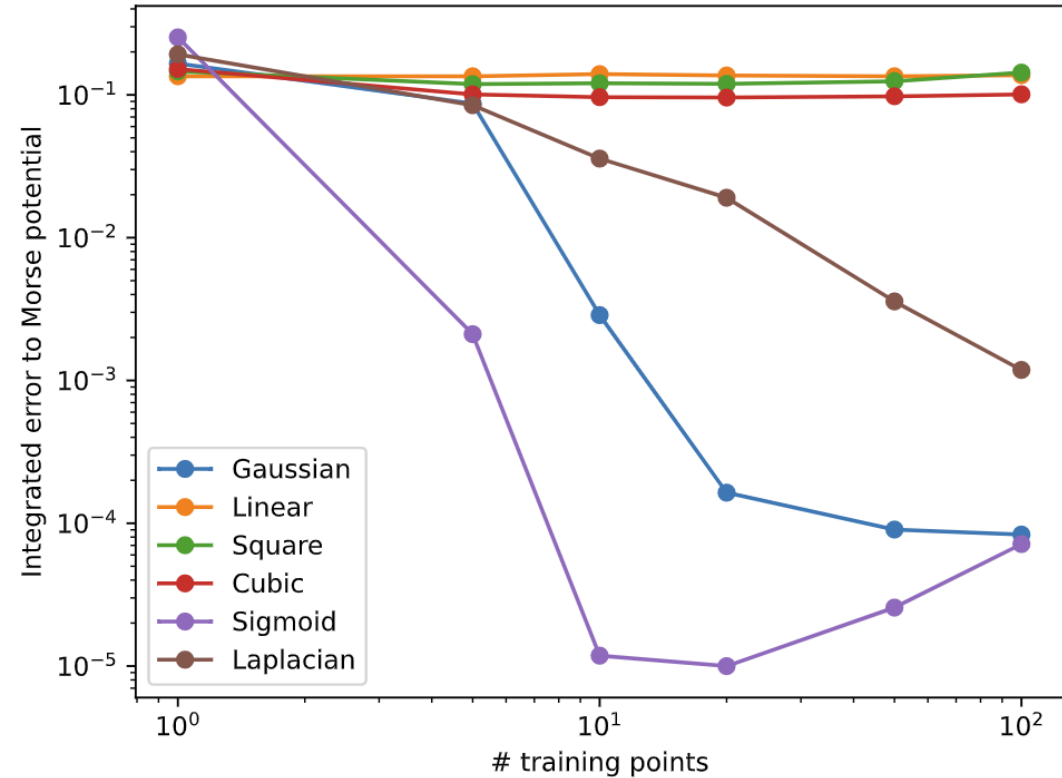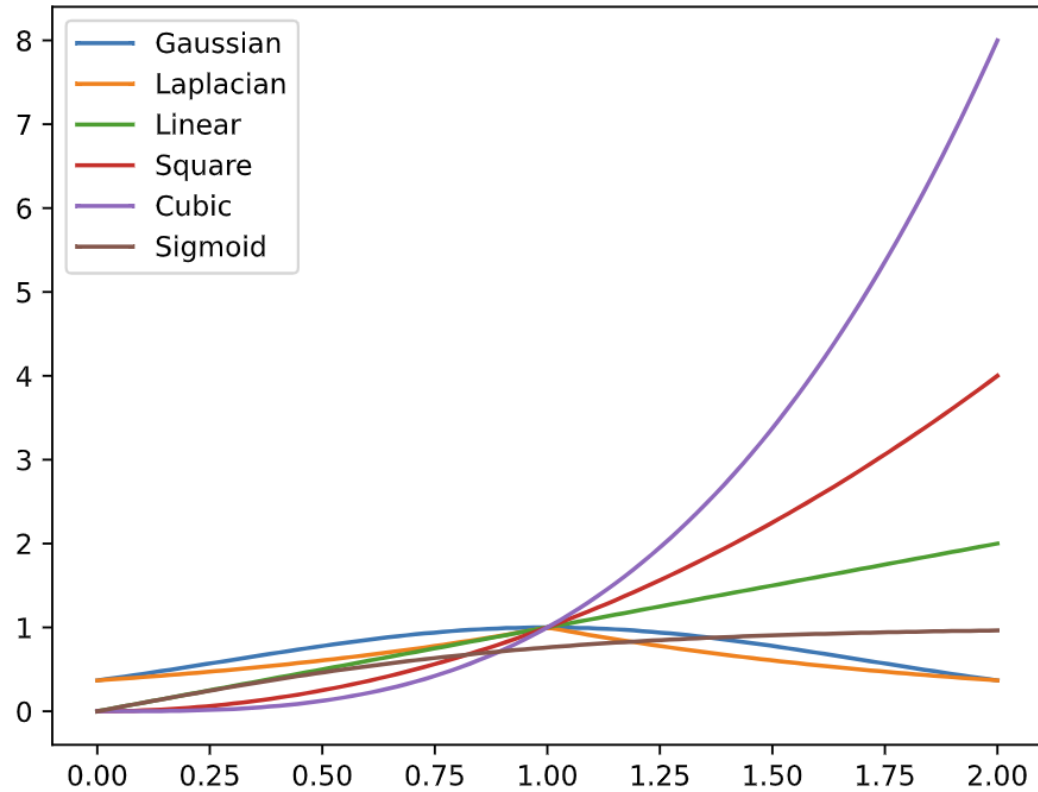
# Kernel-Ridge-Regression

## Procedure
- Get *i* data points with scalar property (label)     $\{q_i\}$
  - E.g. atomisation energy
- Calculate all representations     $\{\mathbf{M}_i\}$
  - typically ~1k
- Find distance and kernel matrices     $\mathbf{D}, \mathbf{K}$
  - Symmetric
- Train model for predictions     $\{\tilde{q}_i\}$

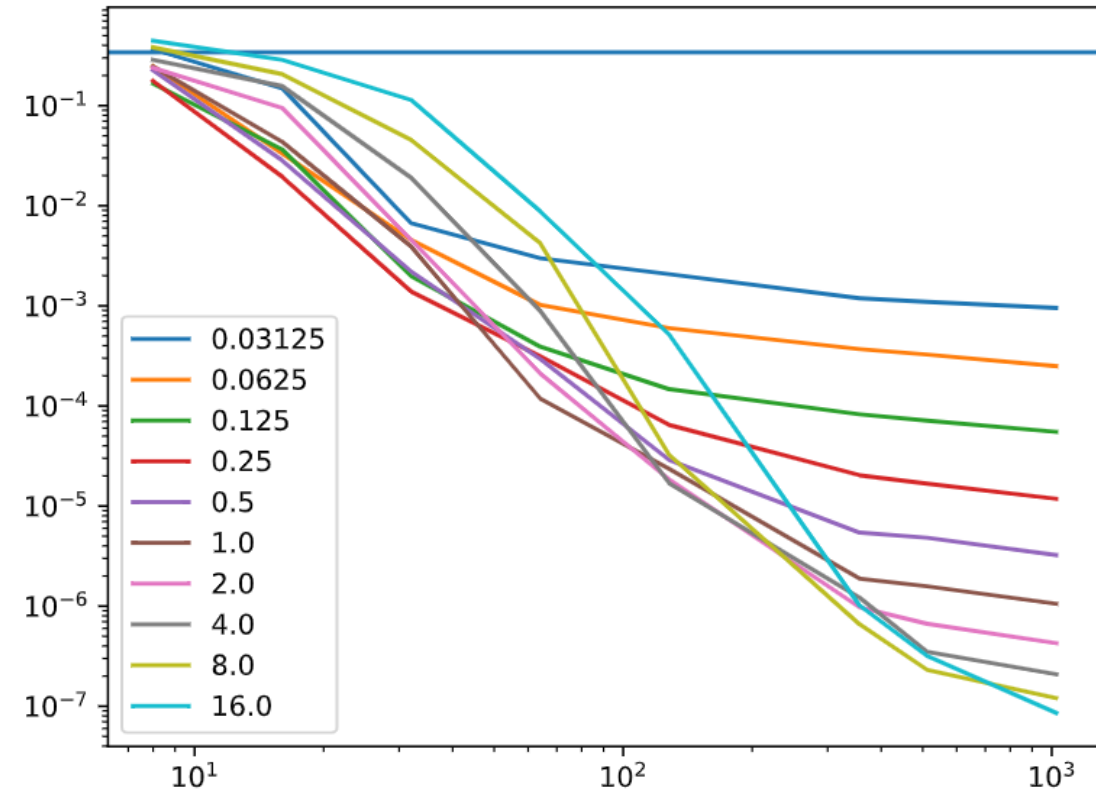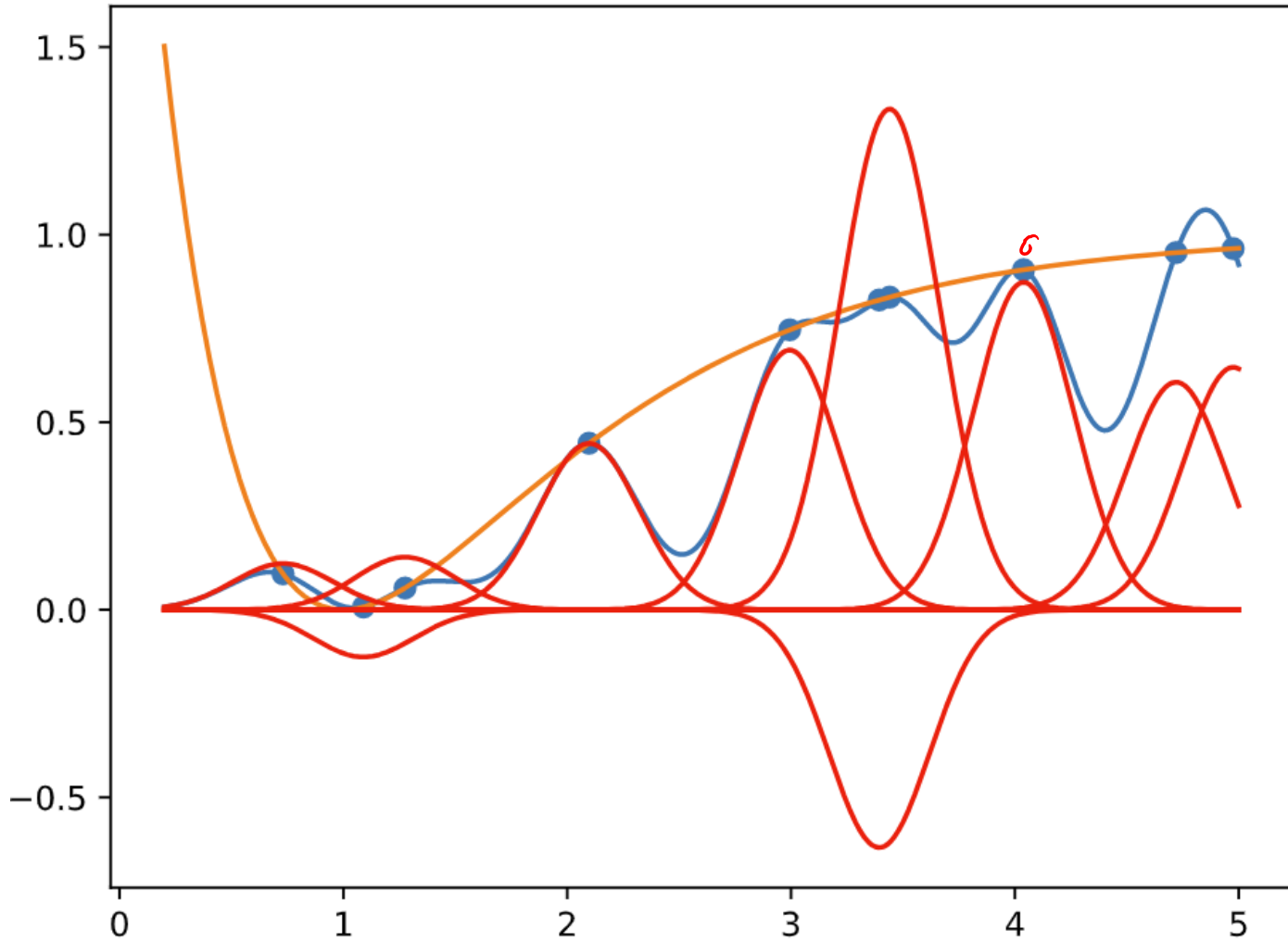$$\arg\min_{\alpha} \sum_i (q_i - \tilde{q}_i)^2 + \lambda \sum_{ij} \alpha_i \alpha_j k_{ij}$$

$$\Rightarrow \alpha = (\mathbf{K} + \lambda\mathbf{I})^{-1}\, y \qquad\qquad \tilde{q}(\mathbf{M}) = \sum_i \alpha_i k(\mathbf{M}, \mathbf{M}_i)$$

Gaussian kernel:   $k(x, y) = \exp(-\gamma\|x - y\|^2)$

$$\iint g(x) K(x, y) g(y) \, dx \, dy \geq 0$$

Pros
- One-Step learning
- Explainable model
- Easy ablatation studies (="what if certain data was not known")
- Efficient cross-validation

Cons
- Need to have features
- Memory requirements can be challenging
- Plenty of hyperparameters
- Instable for large data sets

How to fix
Make them overcomplete
Consider nearest neighbors only
More compute
Regularize

# Summary Kernel Ridge Regression

- Explainable method

- Requires kernel function

- Can be learned without optimization

- Requires hyperparameter scans

- Highly susceptible to kernel function choice

- Regularization required especially for noisy data

- Easy to implement, hard to implement efficiently