

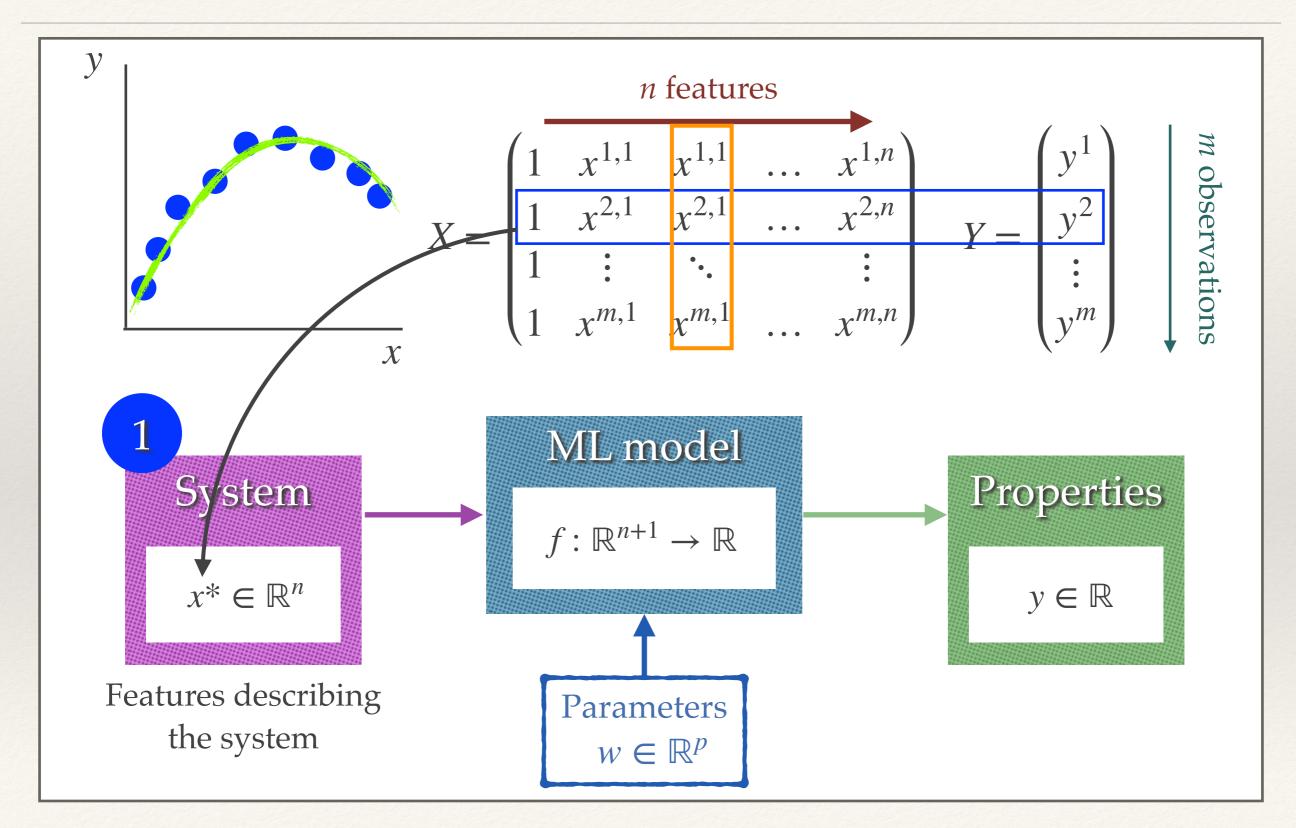
Seyed Mohamad Moosavi, Kevin Maik Jablonka & Berend Smit

Basics of machine learning for chemistry and materials science



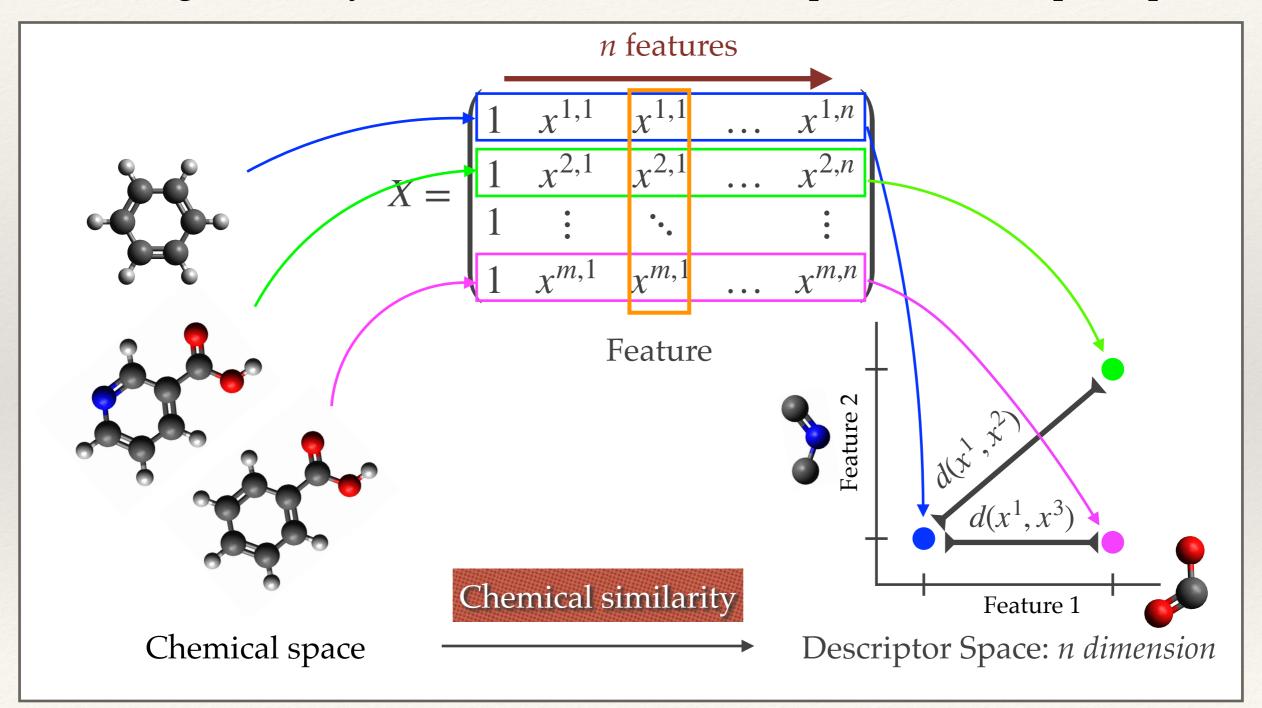


Supervised machine learning



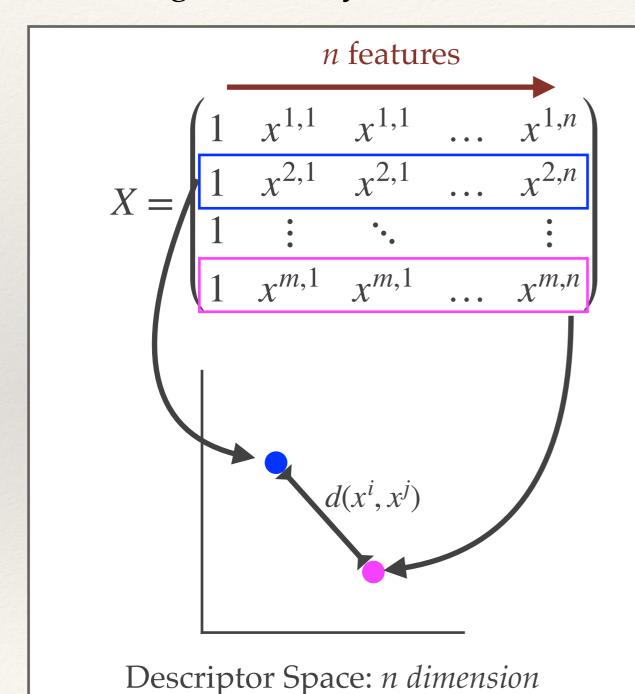
What is featurisation/a descriptor?

Encoding chemistry into numbers: "chemical space" to descriptor space



What makes a descriptor good?

Encoding chemistry into numbers

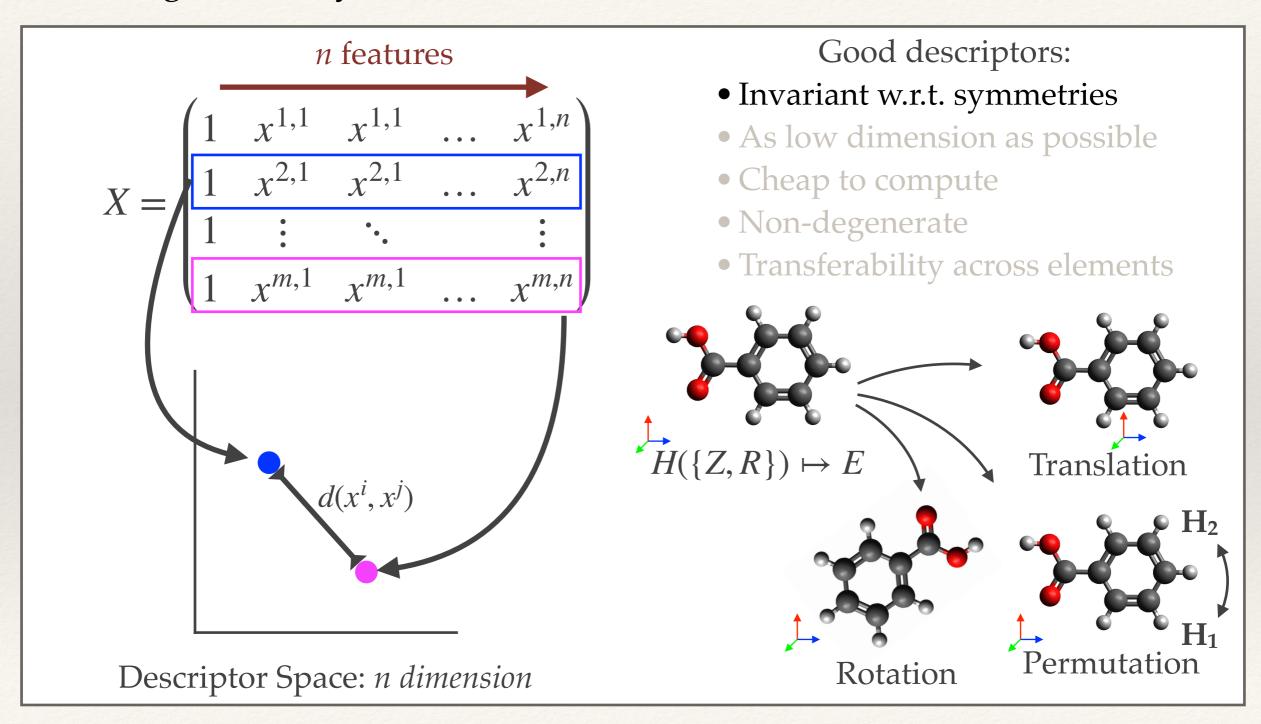


Good descriptors —> obey physics

- Invariant w.r.t. symmetries
- As low dimension as possible
- Cheap to compute
- Non-degenerate
- Transferability across elements

What makes a descriptor good?

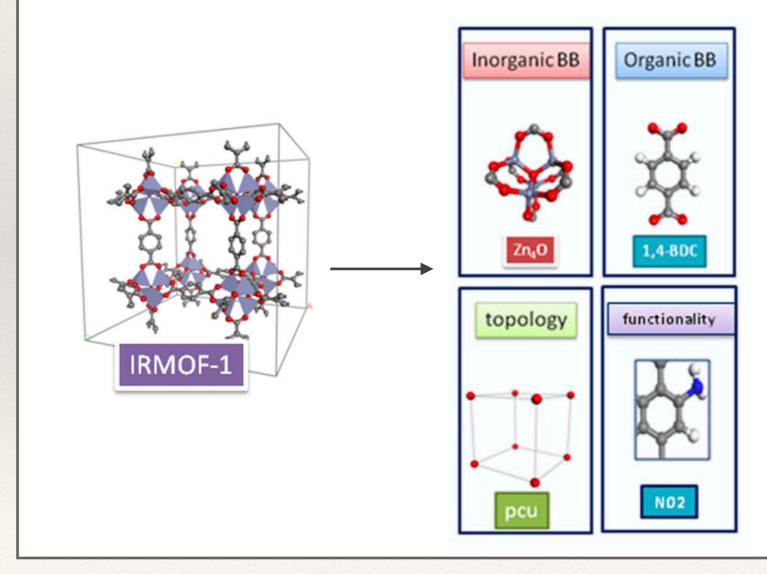
Encoding chemistry into numbers



Adhoc descriptors or properties:

Based on chemical intuition
In principle, can work but often not generalisable

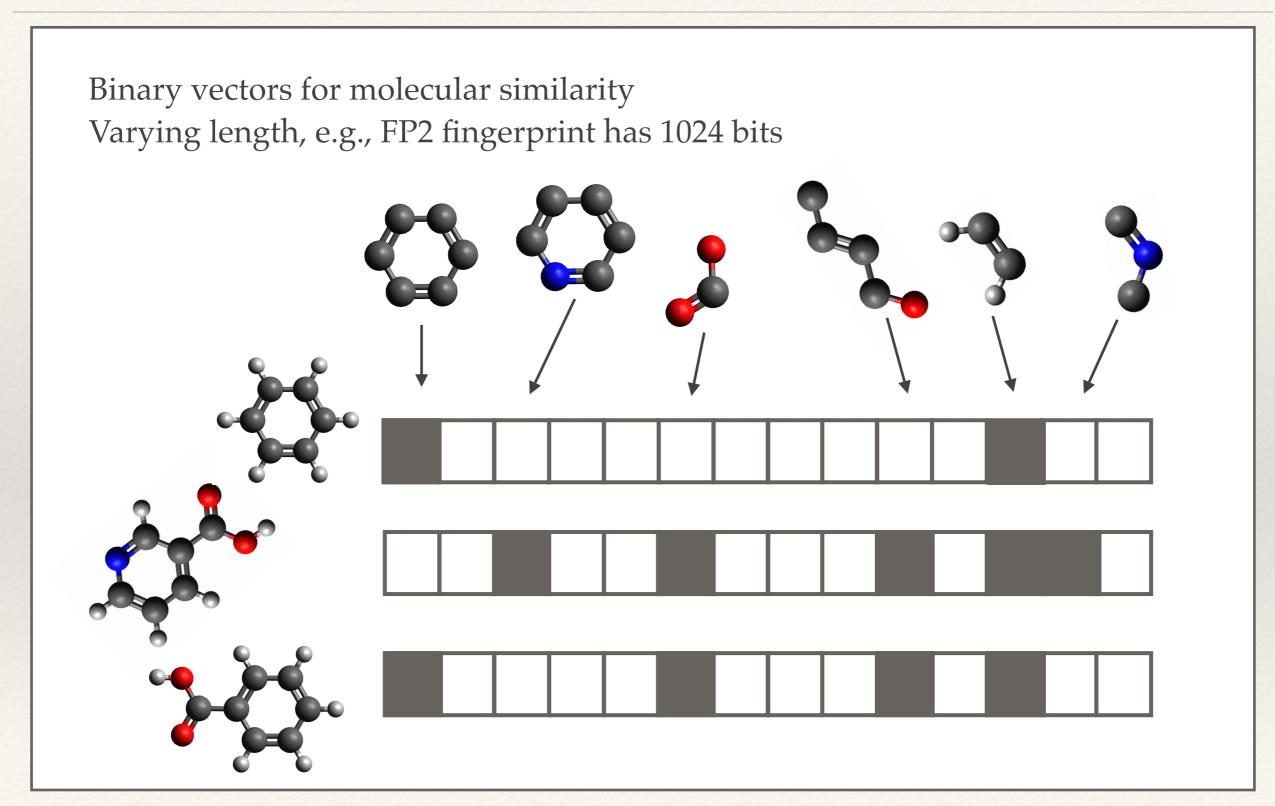
One hot featurisation

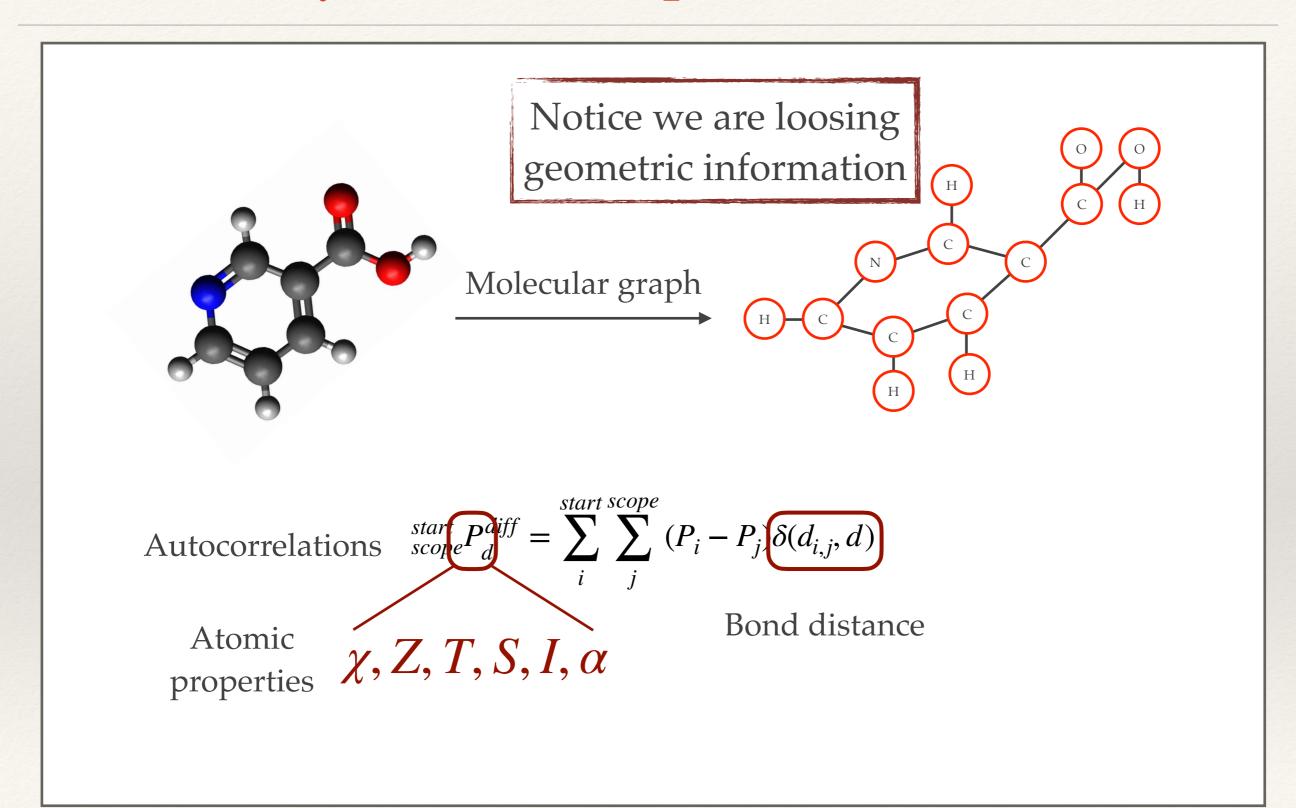


Computed properties:

- Atom identity
- Maximum positive charge
- Minimum negative charge
- etc.

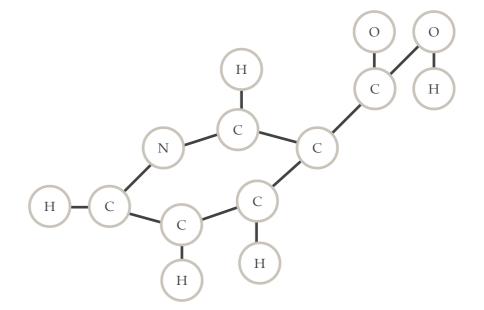
Fragment based descriptors: fingerprints





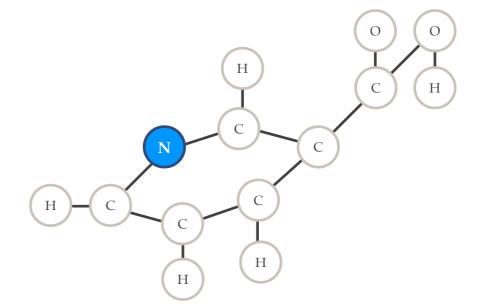
Autocorrelations
$$\underset{scope}{start} P_d^{diff} = \sum_{i}^{start} \sum_{j}^{scope} (P_i - P_j) \delta(d_{i,j}, d)$$

$$\sum_{all}^{[N]} \chi_1^{diff} = \sum_{i}^{[N]} \sum_{j}^{all} (\chi_i - \chi_j) \delta(d_{i,j}, 1)$$



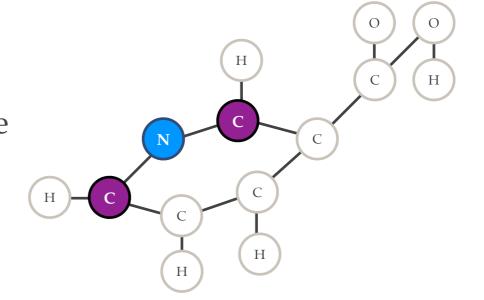
Autocorrelations
$$\underset{scope}{start} P_d^{diff} = \sum_{i}^{start} \sum_{j}^{scope} (P_i - P_j) \delta(d_{i,j}, d)$$

$$\sum_{all}^{[N]} \chi_1^{diff} = \sum_{i}^{[N]} \sum_{j}^{all} (\chi_i - \chi_j) \delta(d_{i,j}, 1)$$



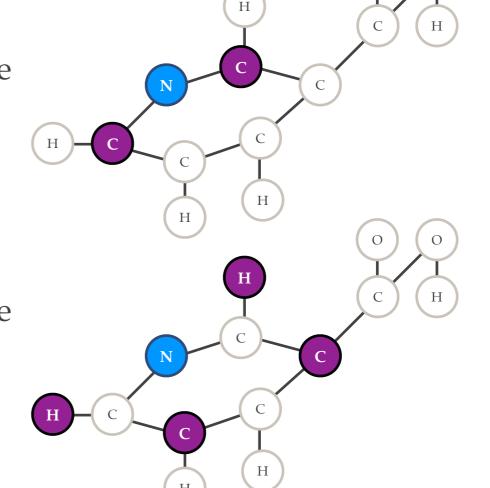
Autocorrelations
$$\underset{scope}{start} P_d^{diff} = \sum_{i}^{start} \sum_{j}^{scope} (P_i - P_j) \delta(d_{i,j}, d)$$

$$\sum_{all}^{[N]} \chi_1^{diff} = \sum_{i}^{[N]} \sum_{j}^{all} (\chi_i - \chi_j) \delta(d_{i,j}, 1)$$



Autocorrelations
$$\underset{scope}{start} P_d^{diff} = \sum_{i}^{start} \sum_{j}^{scope} (P_i - P_j) \delta(d_{i,j}, d)$$

1 bond distance
$$\sum_{i=1}^{[N]} \chi_1^{diff} = \sum_{i=1}^{[N]} \sum_{j=1}^{all} (\chi_i - \chi_j) \delta(d_{i,j}, 1)$$



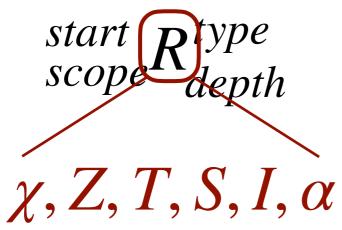
$$\sum_{\substack{[N]\\all}} \chi_2^{diff} = \sum_{i}^{[N]} \sum_{j}^{all} (\chi_i - \chi_j) \delta(d_{i,j}, 2)$$

RACs for MOFs, hands on session in the afternoon



$$\frac{start}{scope} P_d^{diff} = \sum_{i}^{start} \sum_{j}^{scope} (P_i - P_j) \delta(d_{i,j}, d)$$

$$_{scope}^{start}P_{d}^{prod} = \sum_{i}^{start}\sum_{j}^{scope}(P_{i} \times P_{j})\delta(d_{i,j},d)$$



Crystal Graph Separating Linkers Full Linker **Functional** Linker Metal Connecting Center Groups

https://github.com/hjkgrp/molSimplify

Encoding geometry: Coulomb matrix

Inspired by how quantum mechanics works:

$$H(\{Z,R\}) \xrightarrow{\Psi} E \qquad \Leftrightarrow \qquad \{Z,R\} \xrightarrow{\mathrm{ML}} E$$

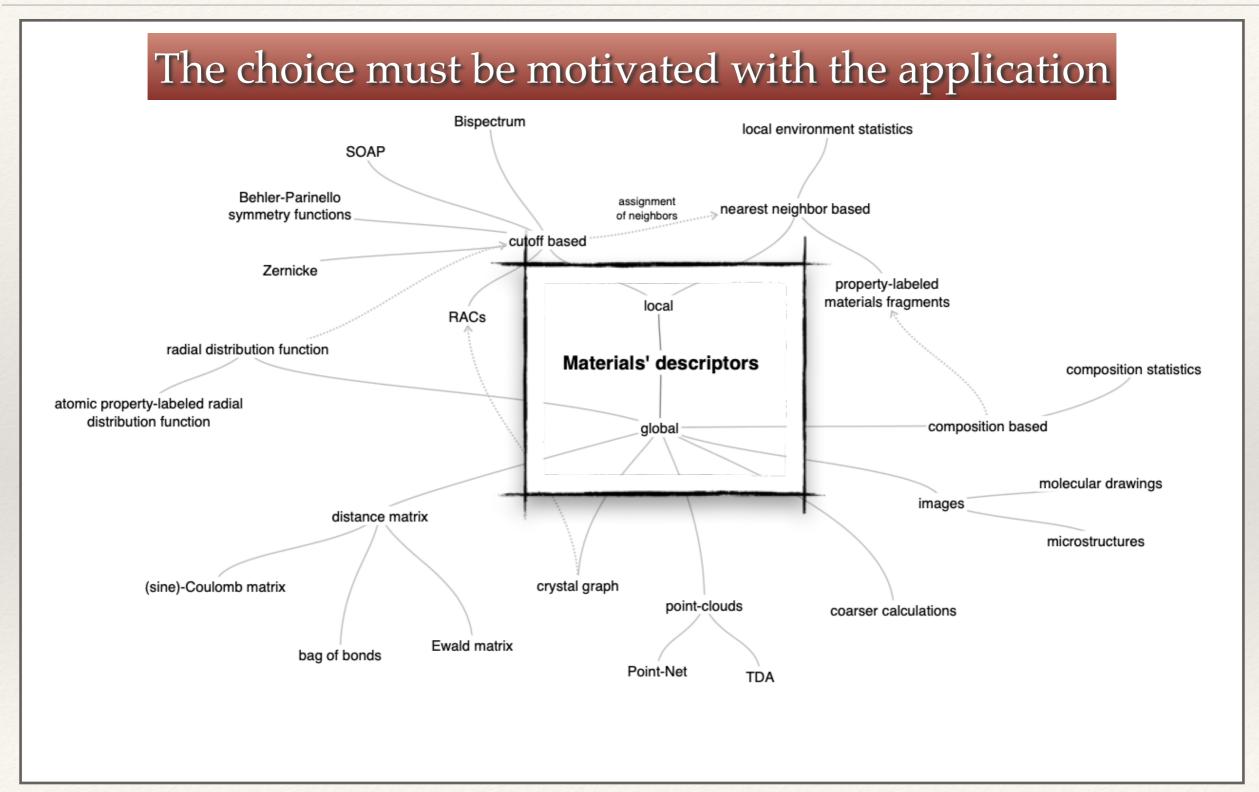
$$M_{ij} = \begin{cases} 0.5 \, Z_i^{2.4} & i = j \\ \frac{Z_i \, Z_j}{(|\mathbf{r}_i - \mathbf{r}_j|)} & i \neq j \end{cases}$$

Similarity is defined as:

The difference in eigenvalues of Ms between two systems

$$d(x^{i}, x^{j}) = d(\epsilon^{i}, \epsilon^{j}) = \sqrt{\sum_{I} |\epsilon^{i} - \epsilon^{j}|}$$

And there exist many more!



Encoding local environments: symmetry functions

Chemical Locality Assumption: decomposing property into local environments

property(descriptor) =
$$\sum_{i}^{\text{atoms}} \text{models}_{i}(\text{descriptor}_{i})$$

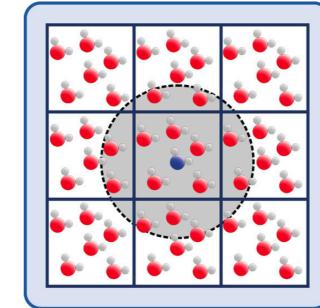
Energy can be decomposed into atomic contributions

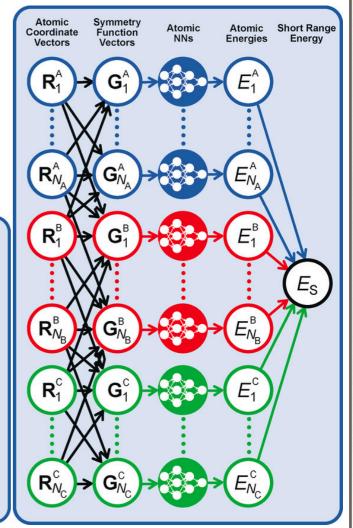
- —> this approach is used to describe PES
- —> Scalable to large systems
- —> differentiability of descriptors is essential

$$E_{S} = \sum_{\nu=1}^{N_{atoms,\nu}} \sum_{\mu=1}^{N_{elem}} E_{\mu}^{\nu}$$

$$f_{cut}(r_{ij}) = \begin{cases} \frac{1}{2} \left[\cos \left(\pi \frac{r_{ij}}{r_{c}} \right) + 1 \right] & \text{for } r_{ij} \leq r_{cut} \\ 0 & \text{for } r_{ij} > r_{cut} \end{cases}$$

$$G_{i}^{2} = \sum_{j} \exp \left[-\eta_{i}(r_{ij} - r_{si})^{2} \right] f_{cut}(r_{ij})$$





Complexity and richness

Higher Learning Capability

Symmetry functions

3D structure

Coulomb matrix

Connectivity-based RACs

Fragment-based Fingerprints

Ad hoc descriptors

Summary of featurisation

- * The aim is to map chemical space to numbers, such that:
 - Chemical similarity is preserved
 - Physics obeyed
- Many kinds of representation exist
 - * Global vs. Local
 - Richness and complexity
- * Should choose representation based on the application