

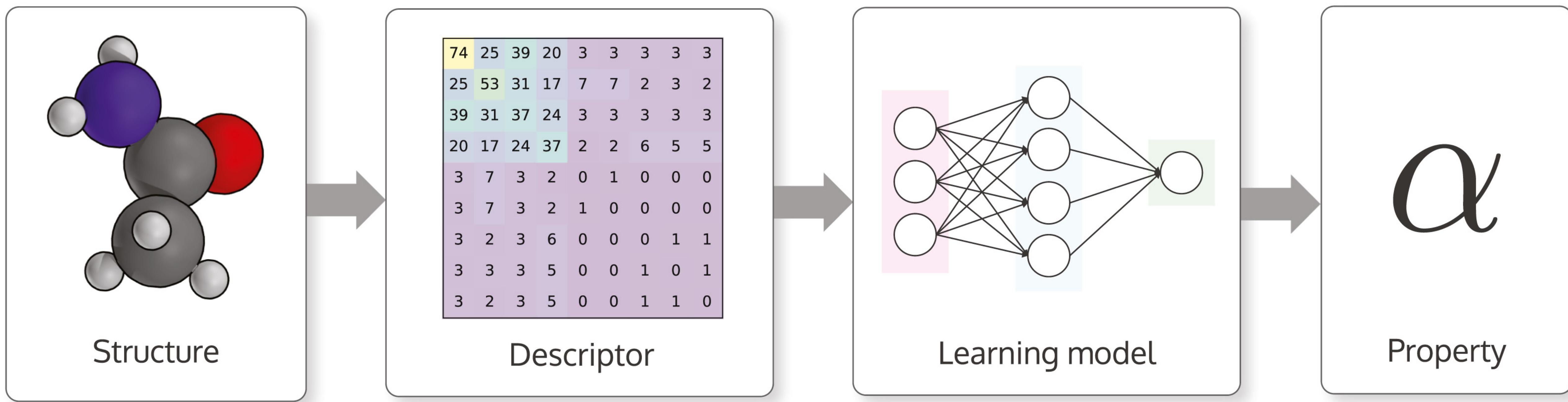
Lecture #6: Atomic structures encoding

Previously on

- Supervised machine learning
- Ridge regression
- Random forest

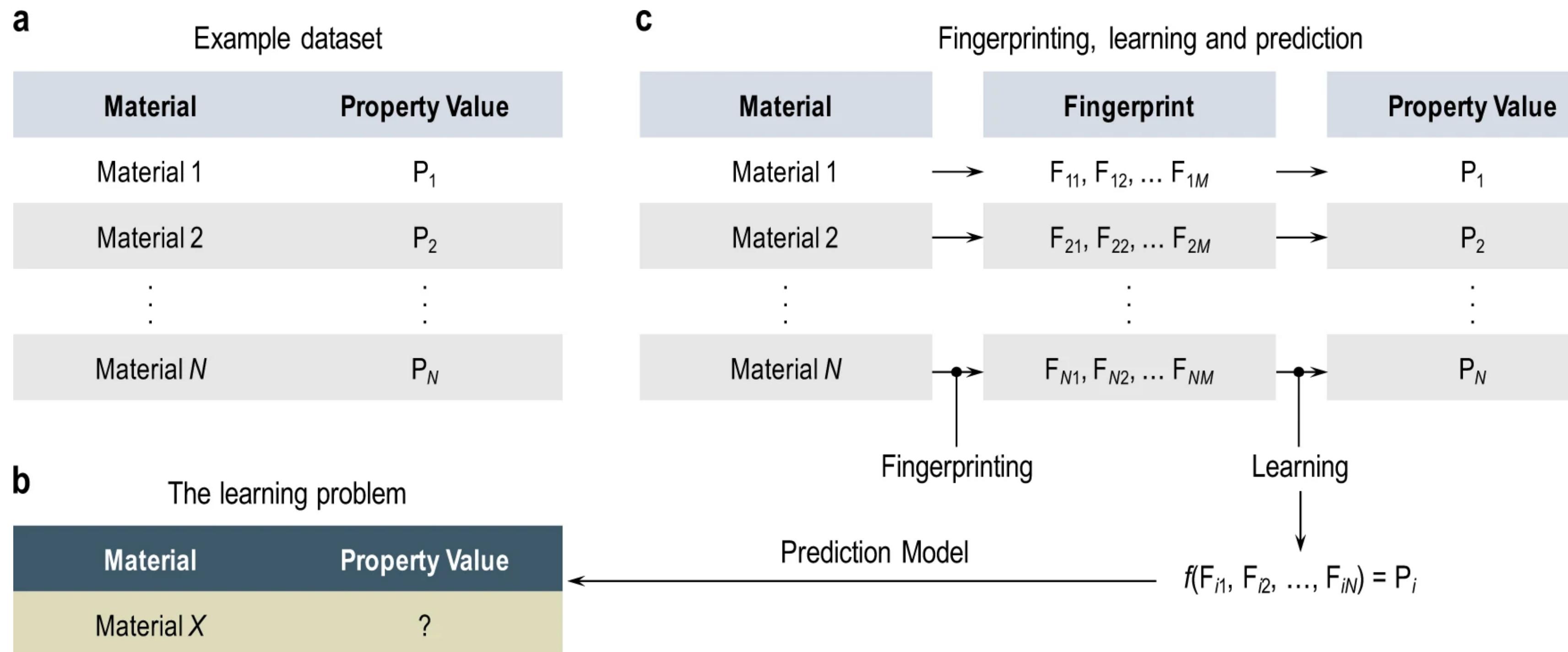
Goals/Agenda

- Hierarchy of crystal/molecular structure descriptors
- Feature importance
- Crystal structure fingerprints



Features

- In classical ML we use features
- to represent materials in a machine-readable format
- Think of it like a fingerprint



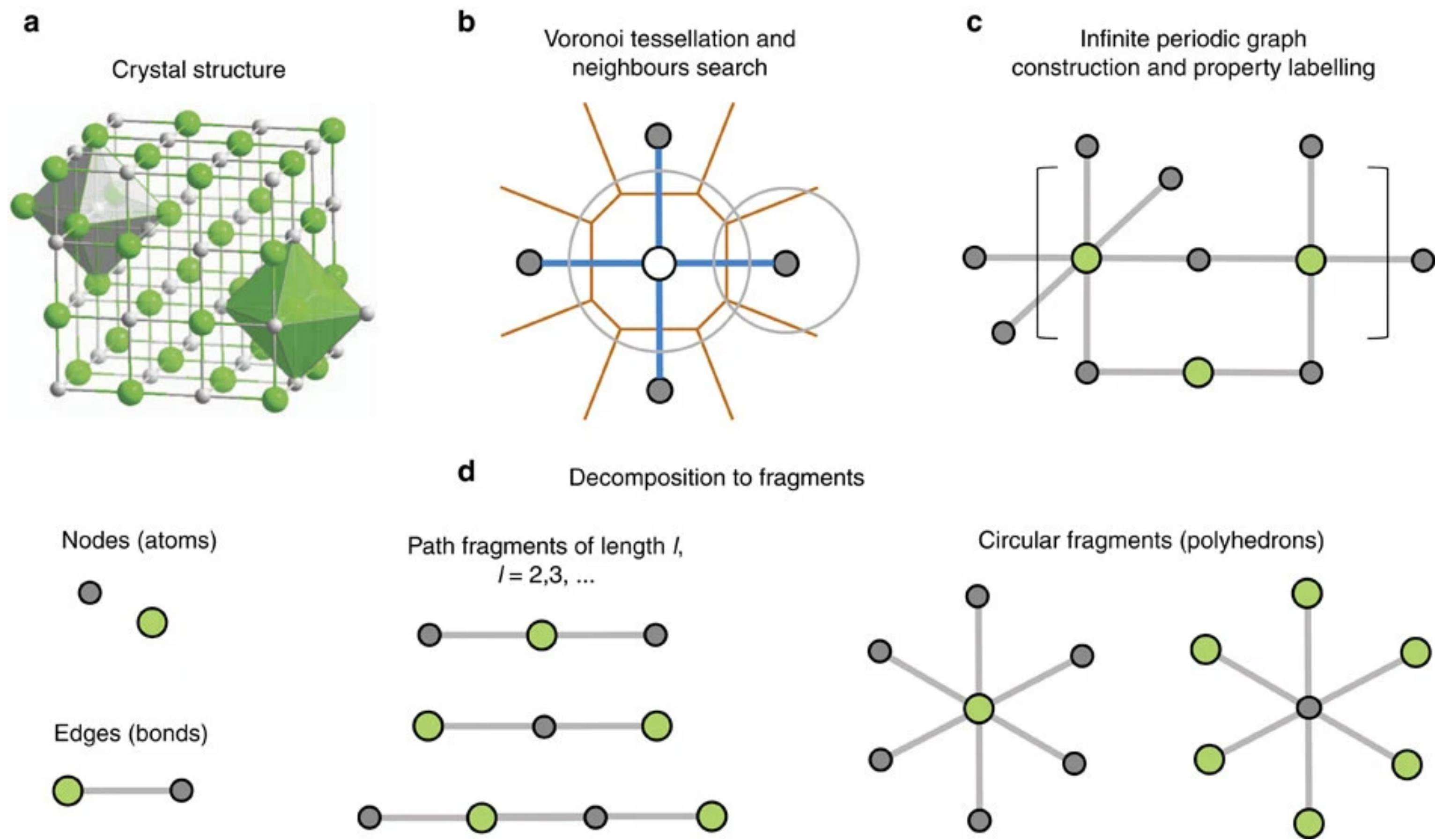
Effective descriptor design is prerequisite to achieving optimal model fit.

Geometrical and compositional encoding of atomic structure

Hierarchy of features

Depending on the resolution we have:

- Local descriptors
 - site
- Fragment descriptors
 - bond
 - polyhedron
- Global descriptors
 - chemical family
 - structural type
 - density

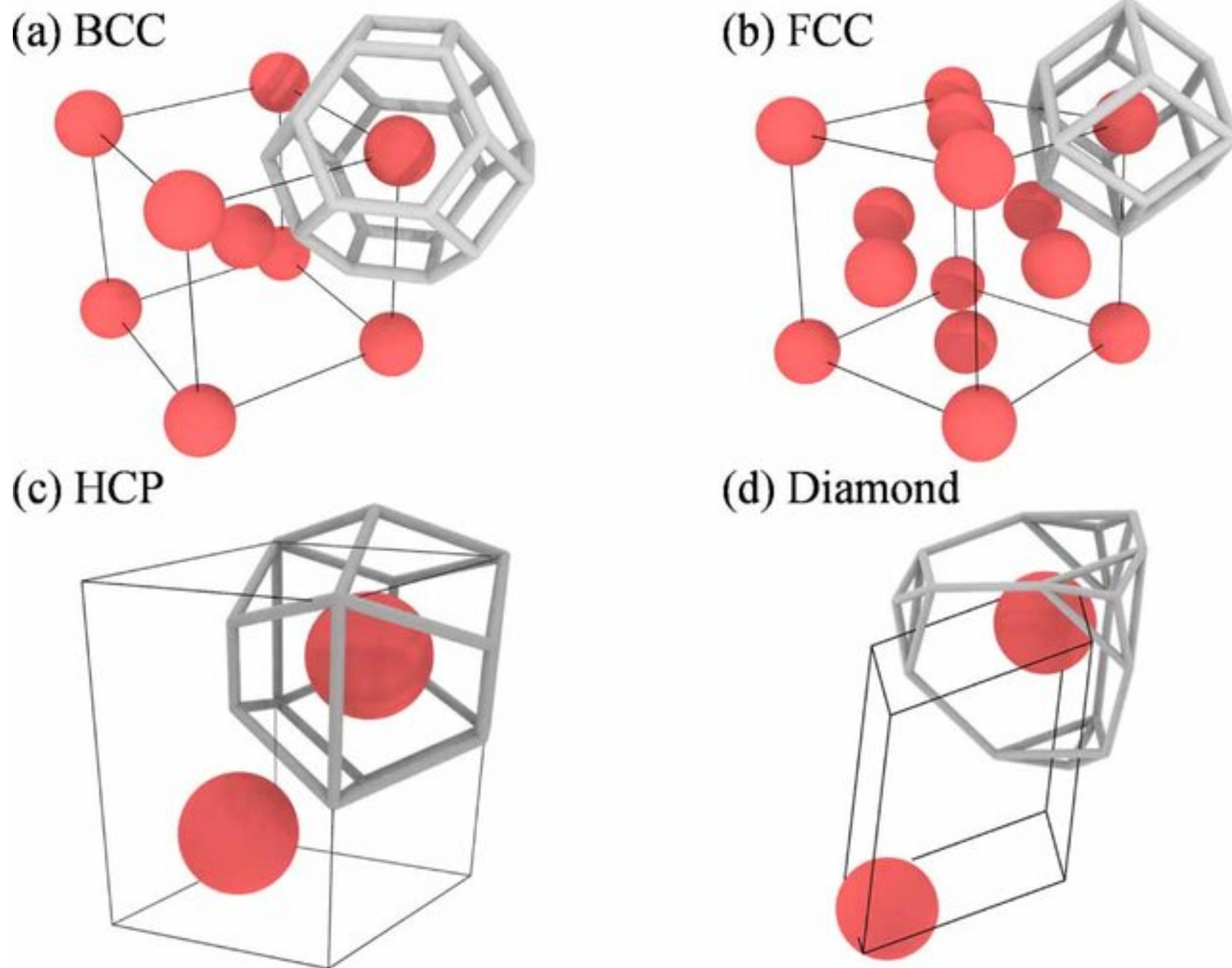


Geometrical (Structural) features

- Atomic packing (e.g. volume per atom, density)
- Voronoi polyhedra features
 - area, volume, face distance, solid angle
- bond distance/angle

Can be calculated for sublattices

Min/max/mean statistics



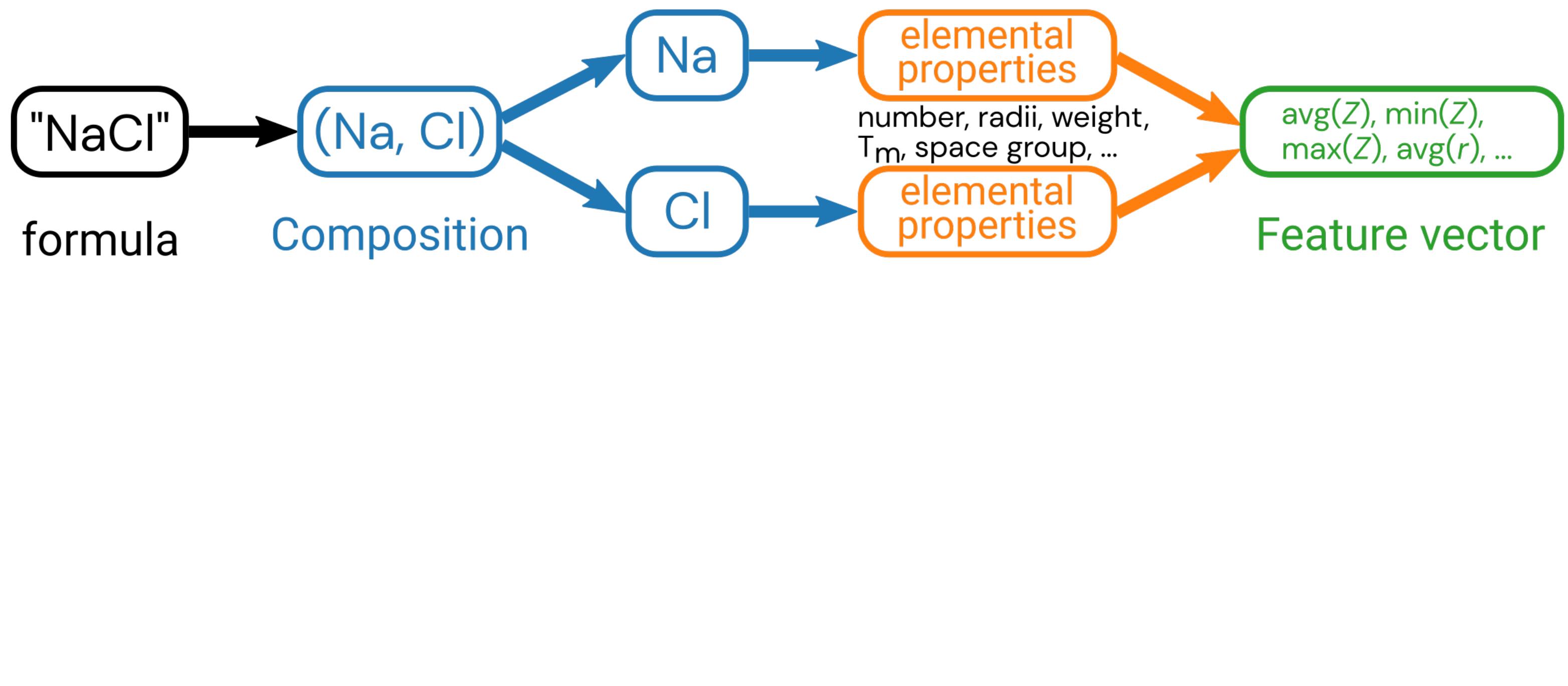
Voronoi polehdra features can be collected by pymatgen's methodology

```
from pymatgen.analysis.local_env import VoronoiNN  
features = VoronoiNN().get_voronoi_polyhedra(structure, site_id)
```

Compositional (elemental/atomic) features

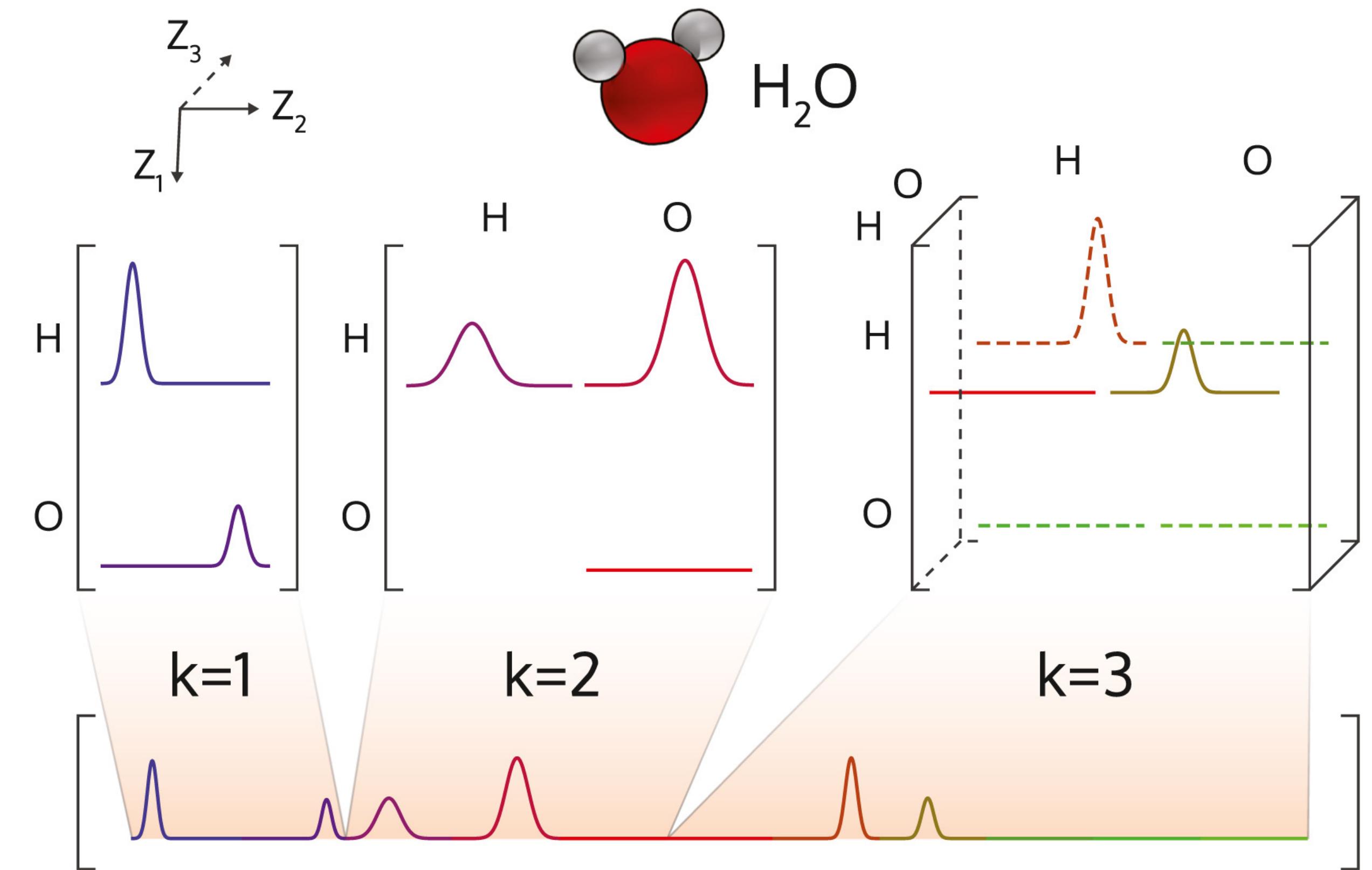
Aggregation over

- Atomic numbers
- Valence electrons
- Covalent radii
- Atomic fraction of each element in a composition
- Stoichiometry related
- Electron affinity
- Ionic radii
- Oxidation states
- Electronegativity
- ...



Atomic structure fingerprints (not taught in the course)

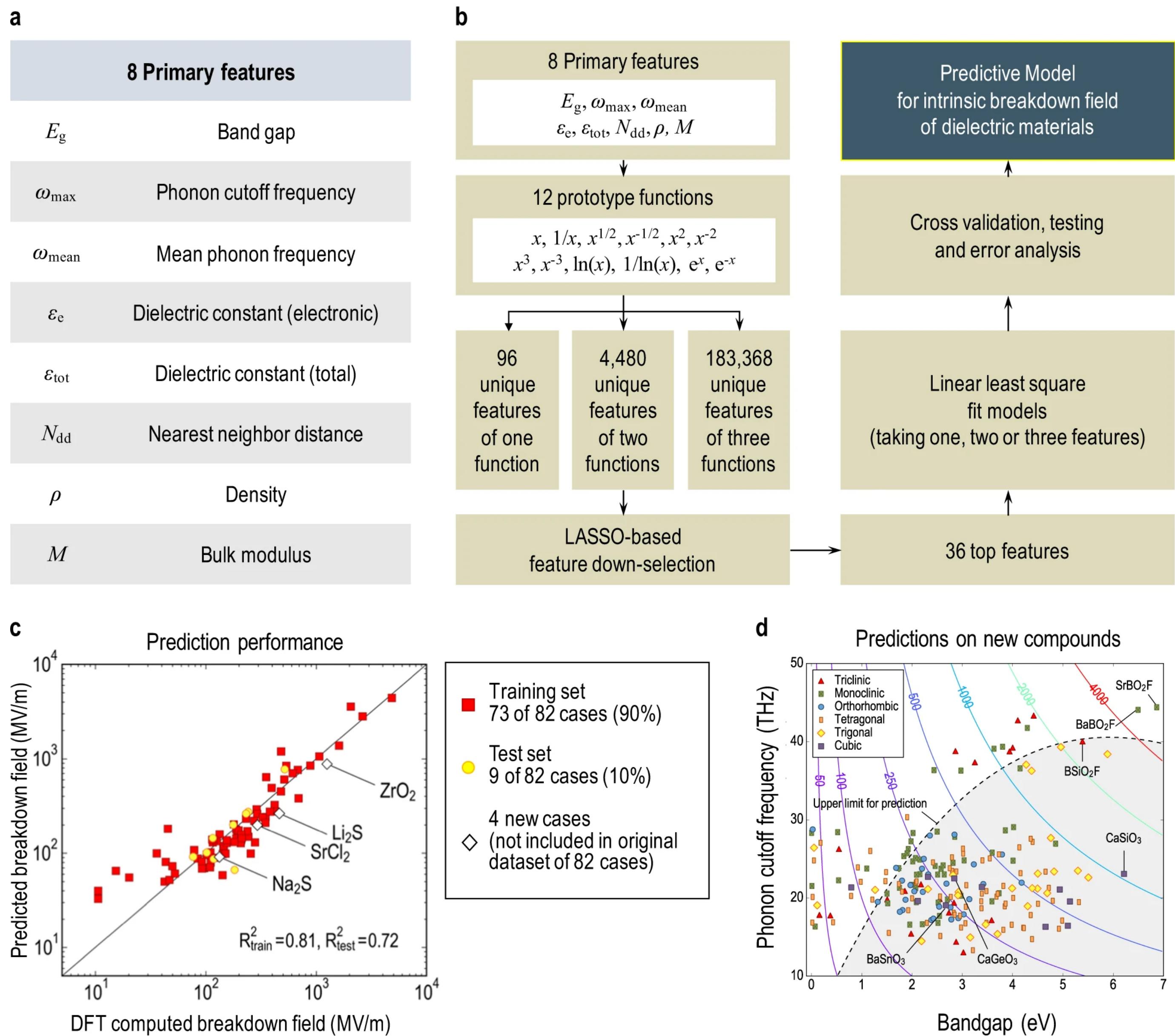
- Atom-centered Symmetry Functions
- Smooth overlap of atomic positions
- Many-body tensor representations
- Pair distribution function
- X-ray diffraction pattern



Feature engineering

- Primary descriptors are used to design more complex features
- e.g. by applying set of mathematical operators

See SISSO paper



Pros

- Better performance compared to primary descriptors

Cons

- Increased computational complexity
- Garbage feature needs to be filtered
- Lack of interpretability

Permutation feature importance

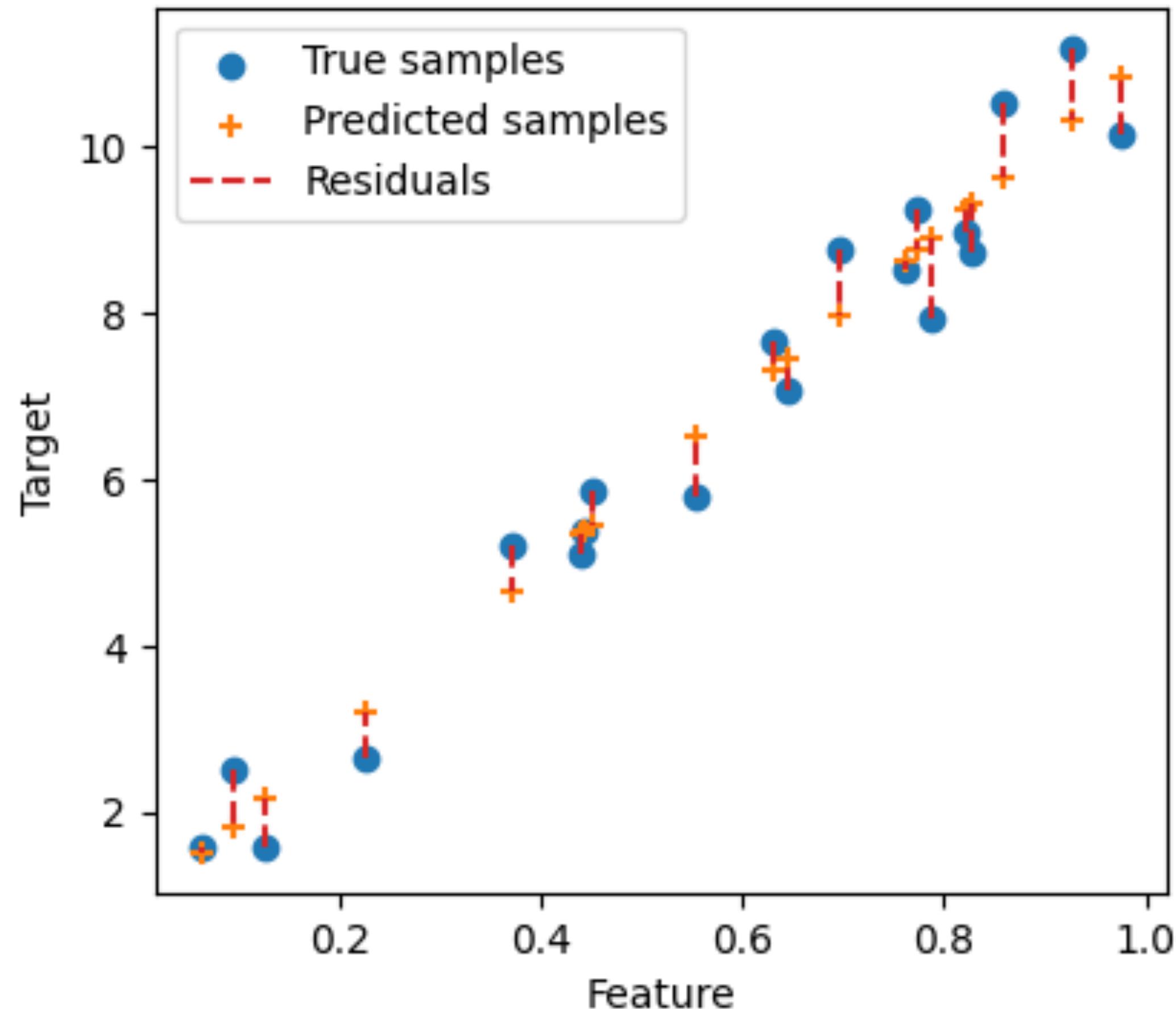
We want to know which features most influence model prediction

Given dataset $\{\mathbf{X}, y\}$

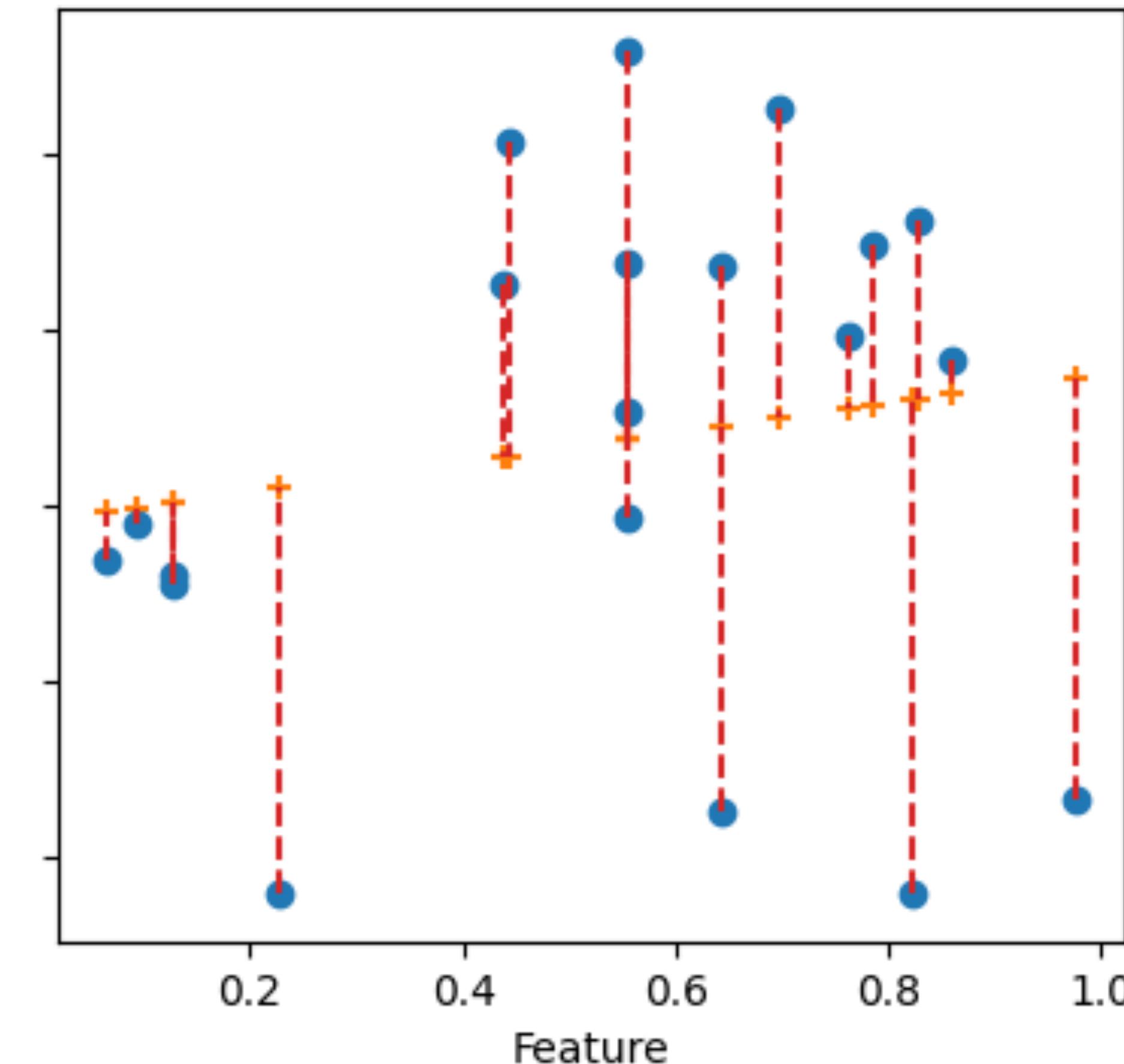
- Fit the model
- Get scores
- Randomly shuffle one of the feature vectors \mathbf{x}_i
- Refit model
- Get scores
- The higher degradation of the model performance the more important the feature

Effect of permuting a predictive feature

Linear model without feature permutation
Mean Absolute Error: 0.51



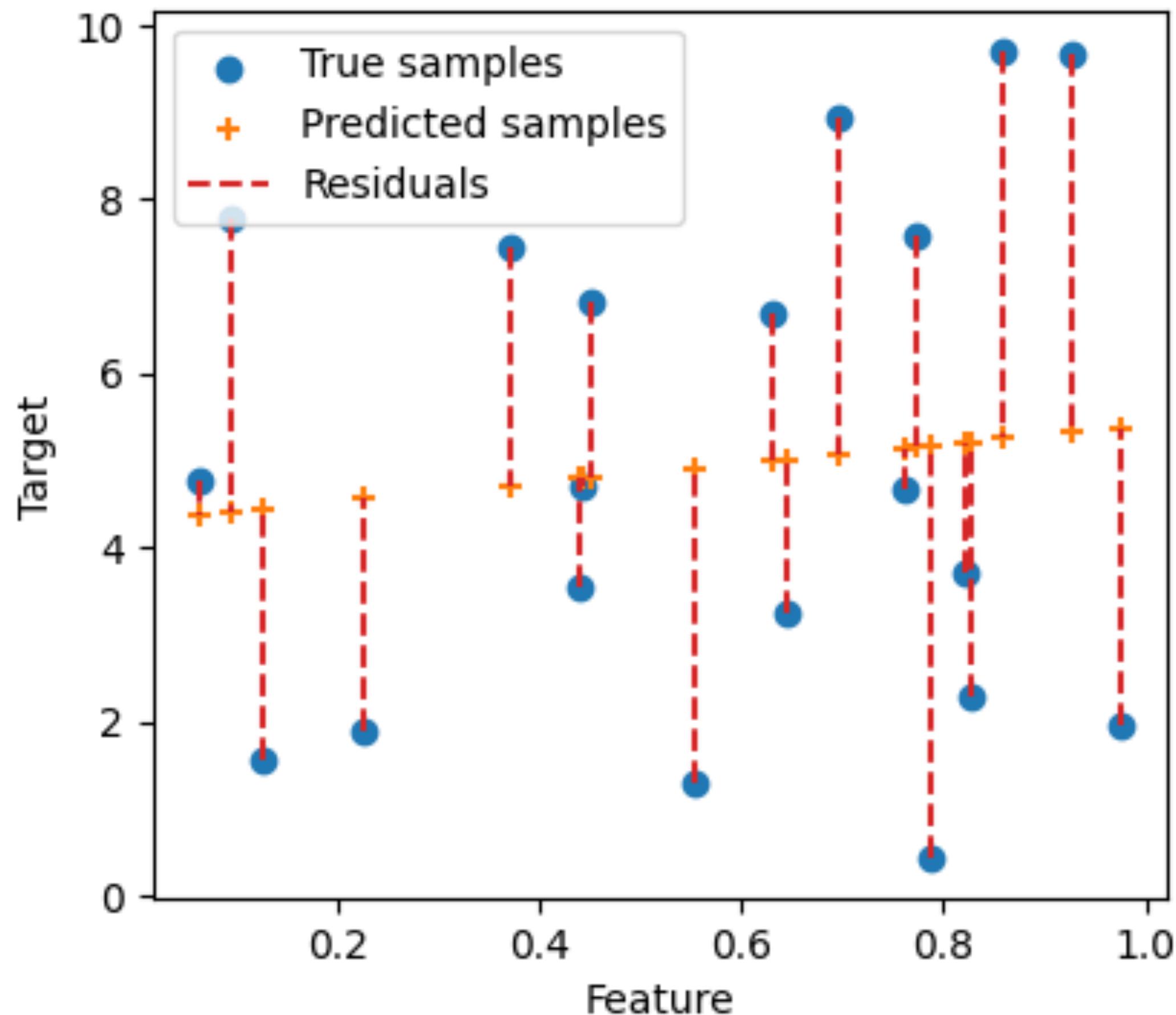
Linear model with feature permutation
Mean Absolute Error: 2.28



Effect of permuting a non-predictive feature

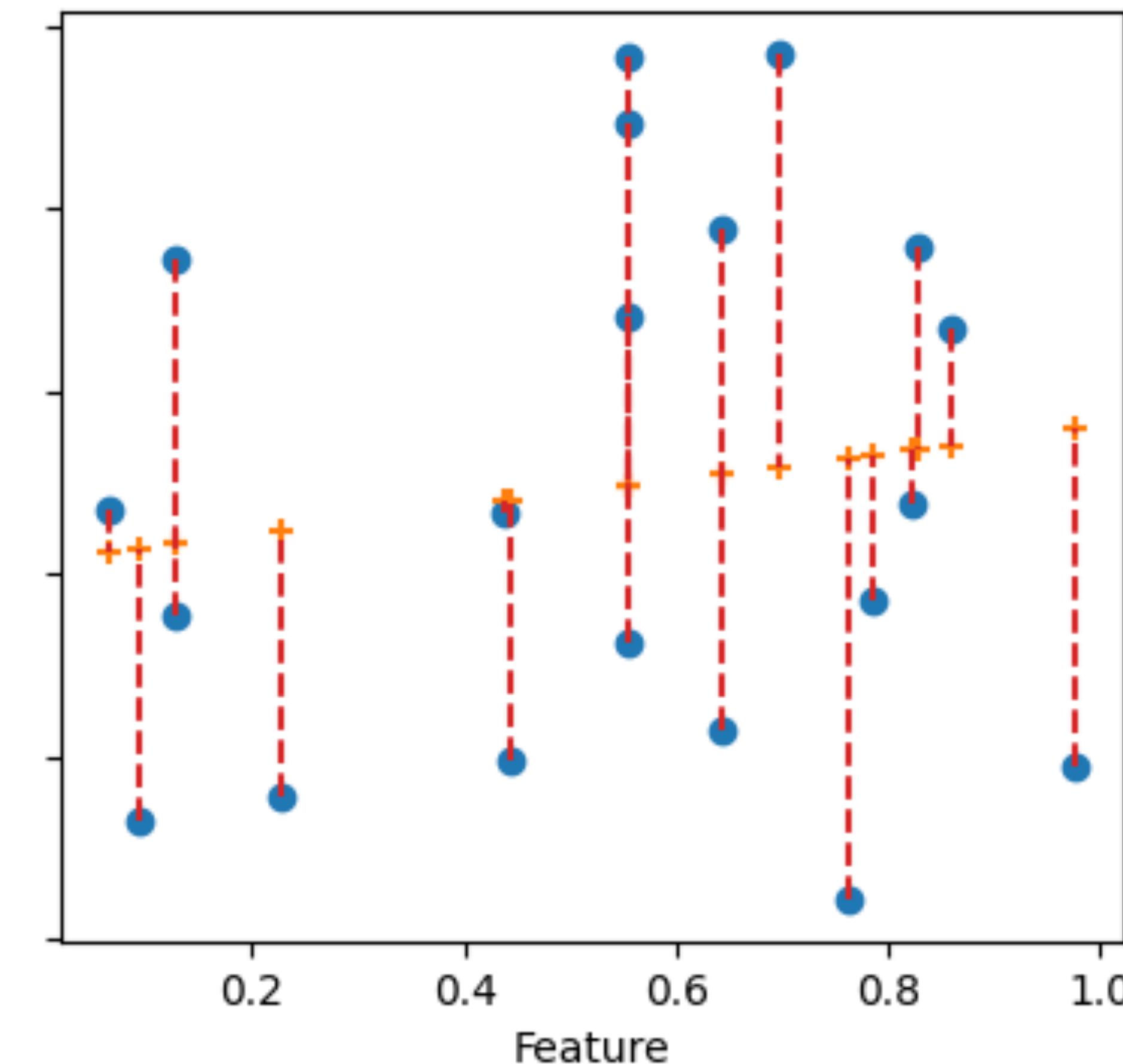
Linear model without feature permutation

Mean Absolute Error: 2.53



Linear model with feature permutation

Mean Absolute Error: 2.49



Pros

- Works for any model
- Yields variance of the feature importance

Cons

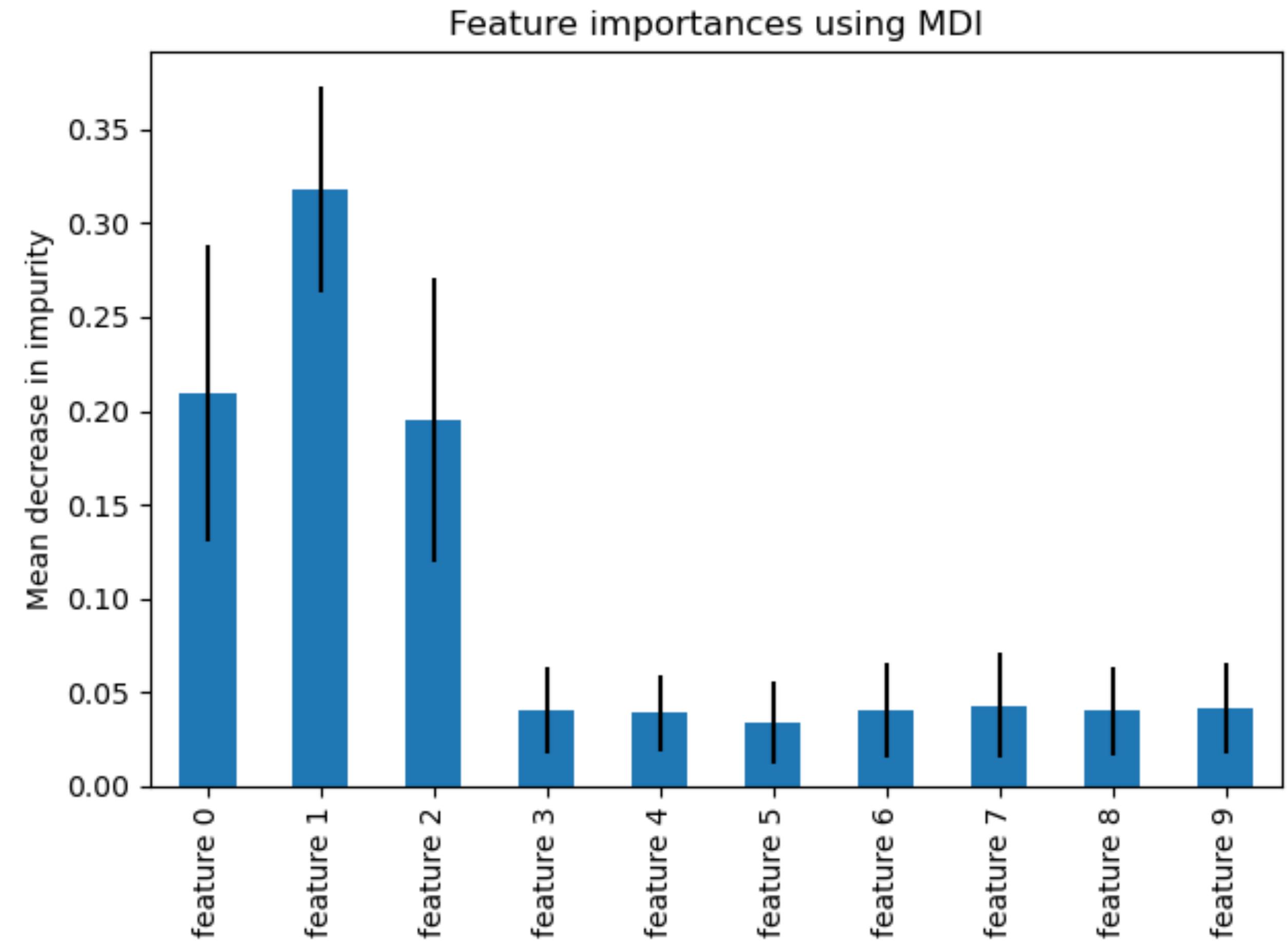
- Computational complexity
- Requires refitting the model for each feature
- Wrong output for correlated features

Intuition behind the impurity-based feature importance

Decision tree:

- Selects the best feature for splitting at each node
- Best split is measured by gain in purity (Gini impurity or Entropy)
- The higher gain the higher importance of the feature

In the case of Random Forest the feature importances are "computed as the mean and standard deviation of accumulation of the impurity decrease within each tree."



Pros

- You have this information after fitting the model
- Yields variance of the feature importance

Cons

- Wrong output for high-cardinality features

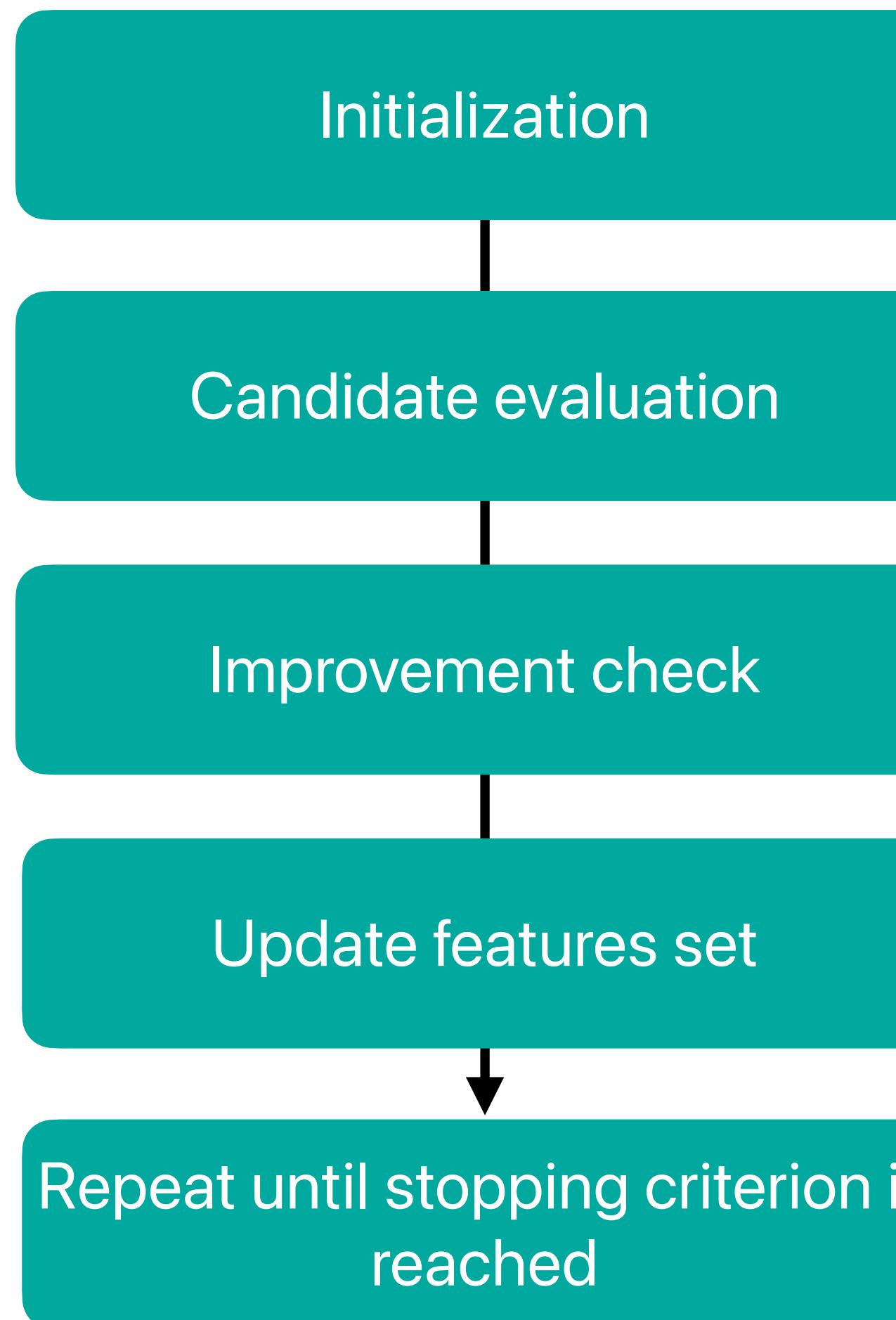
Feature selection

- Drop features with low variance
- Backward/Forward feature selection
- LASSO
- Tree-based feature selection

Why?

- Remove redundant features
- Develop robust model
- The less features the less compute is needed
- Feature collection = time = money
 - save your resources

Forward feature selection algorithm



```
Initialize: S =  $\emptyset$ , C = all features, best_score =  $-\infty$ , stall = 0

While C  $\neq \emptyset$  AND stall < patience:
    best_f = None, current_best =  $-\infty$ 
    For each f in C:
        score = performance(S  $\cup$  {f})
        If score > current_best:
            current_best = score, best_f = f

    If current_best > best_score + tolerance:
        S = S  $\cup$  {best_f}, C = C \ {best_f}
        best_score = current_best, stall = 0
    Else:
        stall += 1
```

How to collect features

Design your own featurizer

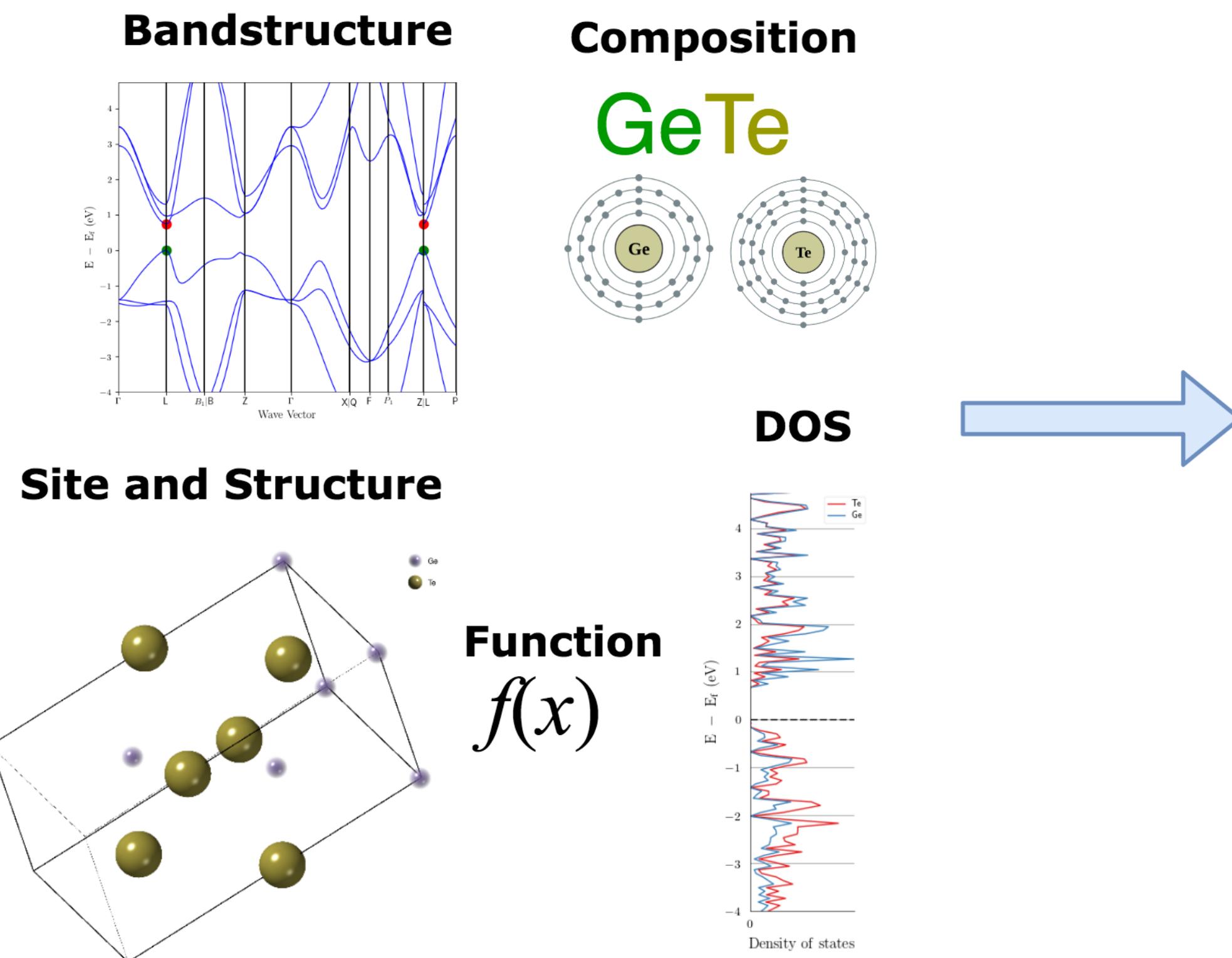
- You are an expert in your field
- Consider essential descriptors of your target
- Use pymatgen/ase/etc.

Ready-to-use Featurizers

- [Matminer](#) Python Library
- [DScribe](#) Python library

Gross level features

- Generate gross level features non-specific to your task
- Select the best candidates



pandas
 $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$

Structure		Composition		Function	
Feature 1	Feature 2	Feature 3	Feature 4	Feature 1, 1/x	Feature 1, logx
0.1	1e-14	.003	223	10	-1
4.2	1.2e-12	.002	14	.238	.62
1.1	1e-6	.0031	101	.91	.041

The most important properties of an ideal descriptor*:

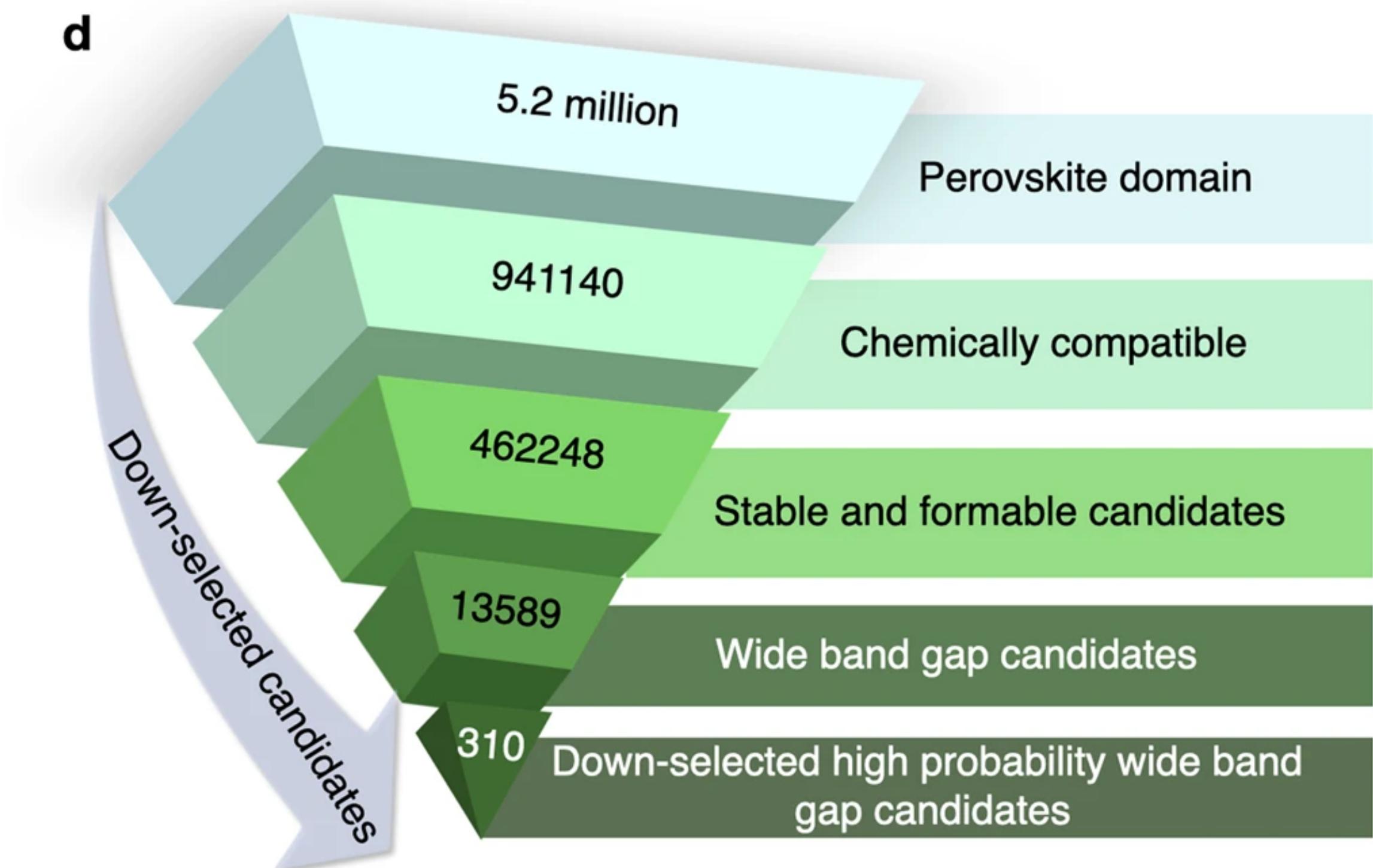
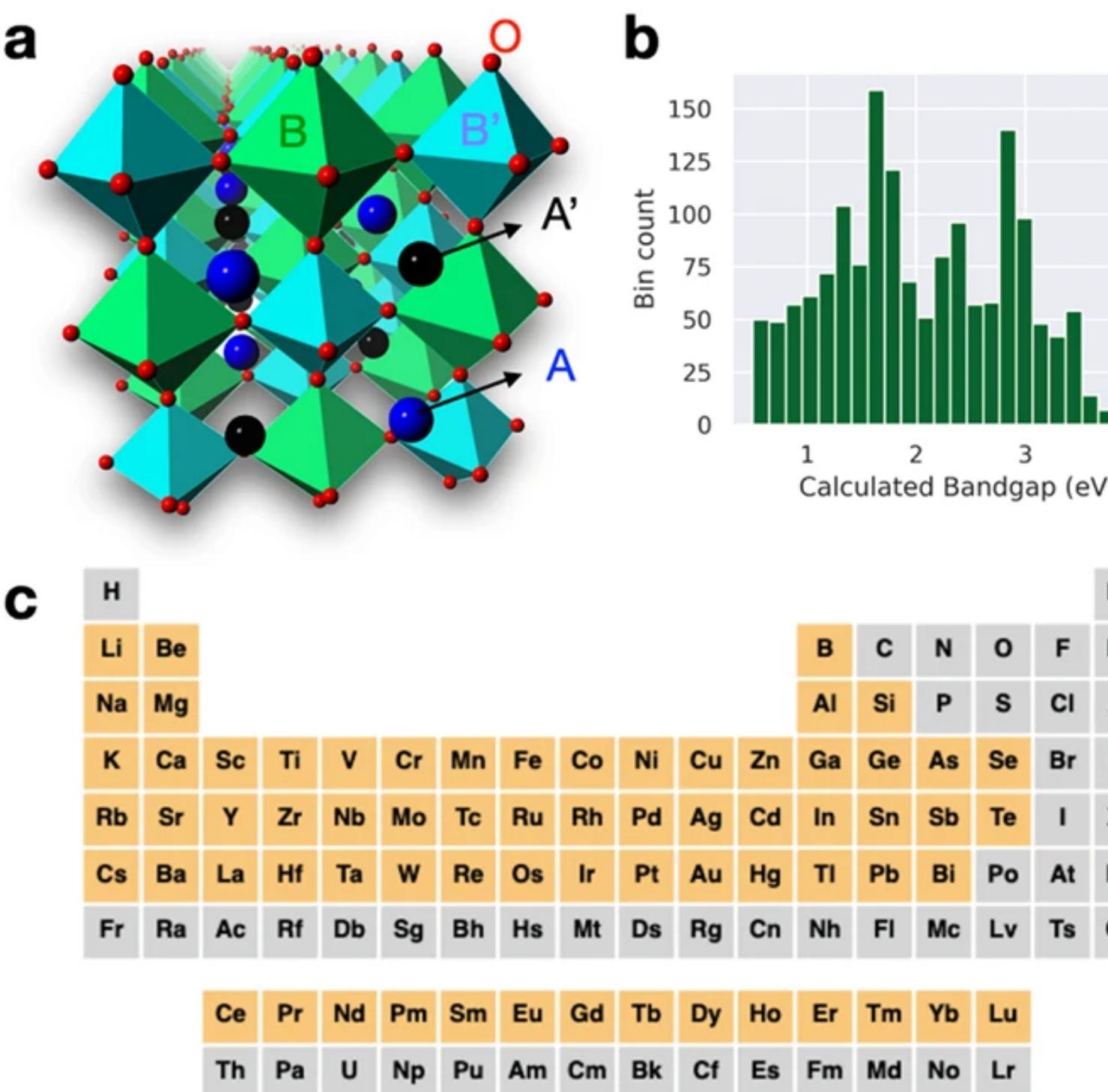
- Invariant with respect translation of the coordinate system
- Invariant with respect to rotation of the coordinate system
- Invariant with respect to permutation of atomic indices: changing the enumeration of atoms does not affect the target
- Unique: single way to construct a descriptor and the descriptor itself corresponds to a single property
- Continuous: small changes in the atomic structure -> small changes in the descriptor
- Compact
- Computationally cheap

Which features to consider?

- What accuracy do I need?
- At which scale your property is defined?
 - Site (i.e. defect formation energy)
 - Bond (i.e. migration barrier)
 - Structure (i.e. hardness)
- Size of the chemical space?
 - Several elements or the whole periodic table
- Diversity of crystal structures?
 - Fixed structural type? Not fixed?
- Do we know (expect) any relationship?
- How many samples do I have? 100? 100,000?

You may find the right feature design by answering these questions

Example



Considered features

- Specific structural type - double perovskites $A_xA'_{2-x}B_yB'_{2-y}X_6$
- Wide chemical space

“...the intention was to use the simplest possible inputs and minimize computational overhead necessary to compile inputs while achieving high prediction accuracies by effectively incorporating notions of chemical similarity across different chemistries”

Abbreviation	Feature
<i>Elemental</i>	
HOMO	Highest Occupied Molecular Orbital (eV)
LUMO	Lowest Unoccupied Molecular Orbital (eV)
IE	Ionization energy (kJ/mol)
X	Pauling Electronegativity
Z radius	Zunger's Pseudopotential radius (a.u.)
EA	Electron affinity (kJ/mol)
<i>Geometric</i>	
t	Tolerance factor
μ	Octahedral factor
$\bar{\mu}$	Mismatch factor

Features importance for E_g prediction

