

Representations for Atomistic Machine Learning

Philip Loche
Stuttgart, 2025/10/07

1. Molecular descriptors (cheminformatics)
 2. Atomistic descriptors (condensed matter)
 3. Descriptor based neural network potentials
 4. End-to-end neural network potentials
-

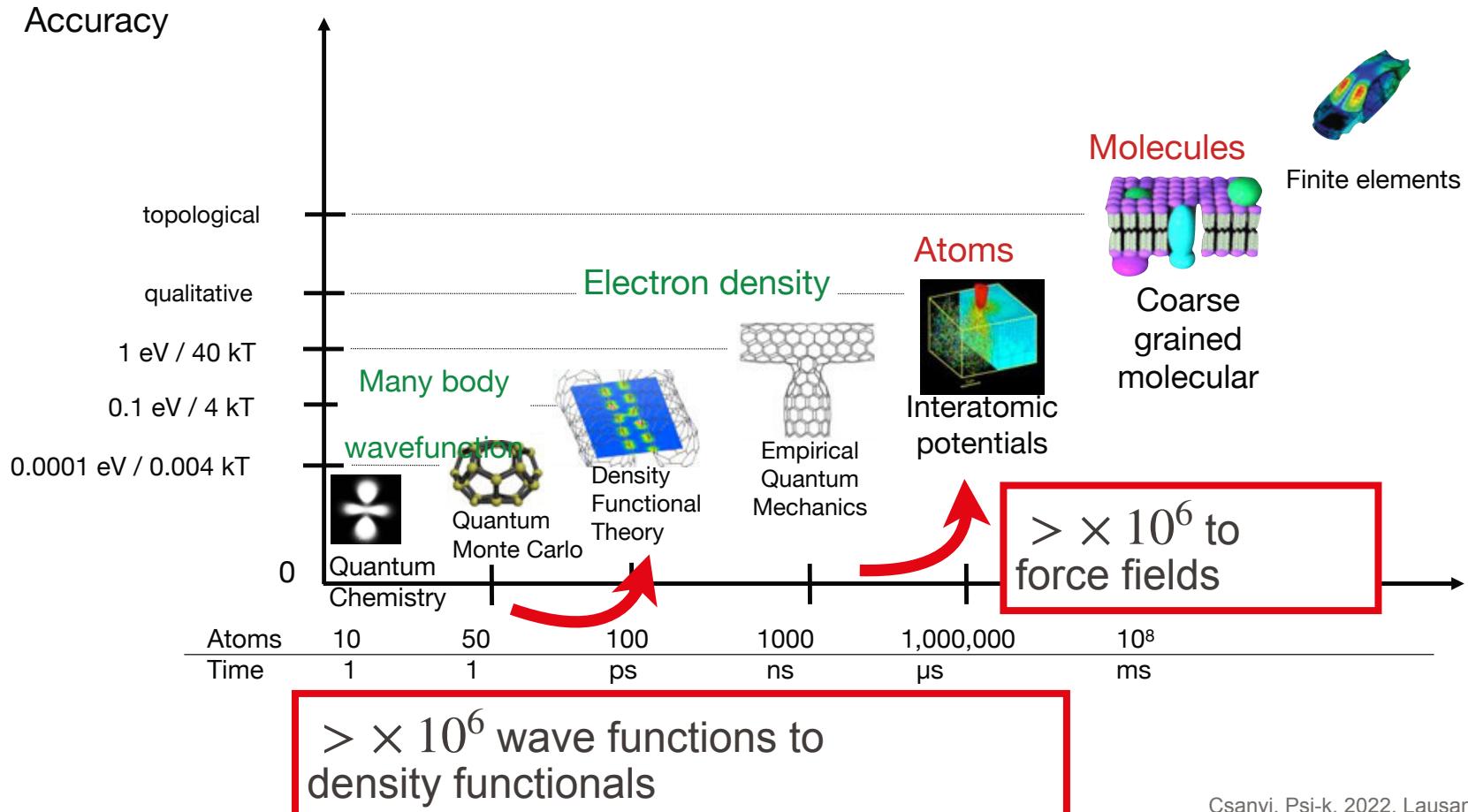
Deliver the vibe of atomistic machine learning!

If you want to dig deeper...

- Musil et. al., “Physics-Inspired Structural Representations for Molecules and Materials”, 10.1021acs.chemrev.1c00021.
- Unke, O. T. et al. “Machine Learning Force Fields”, 10.1021/acs.chemrev.0c01111
- Andrew D. White, Deep learning for molecules & materials, <https://dmol.pub/>

... and try it on your own

- <https://github.com/ceriottm/ml-intro>
- <https://github.com/HFooladi/GNNs-For-Chemists>



DFT: backbone of computational chemistry

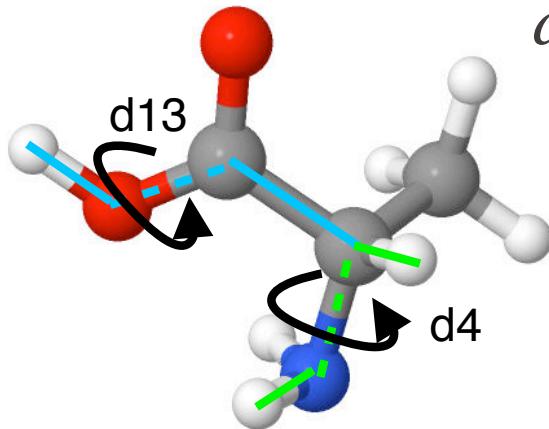
- Quantum mechanics are based on the **Shrödinger's equation** $H\Psi = E\Psi$
 - **Very complex equation** and concepts reduced complexity (and accuracy).
 1. Born-Oppenheimer: Total molecular wavefunction as product of an **electronic and a nuclear wavefunction**.
 2. Kohn-Scham: Base **electron energy** functional $E[\cdot]$ on **density** $\rho(\mathbf{r})$ instead of the wave function Ψ .

DFT is in principle an exact ground state theory.

The **energy functional** needs to be approximated in practice.

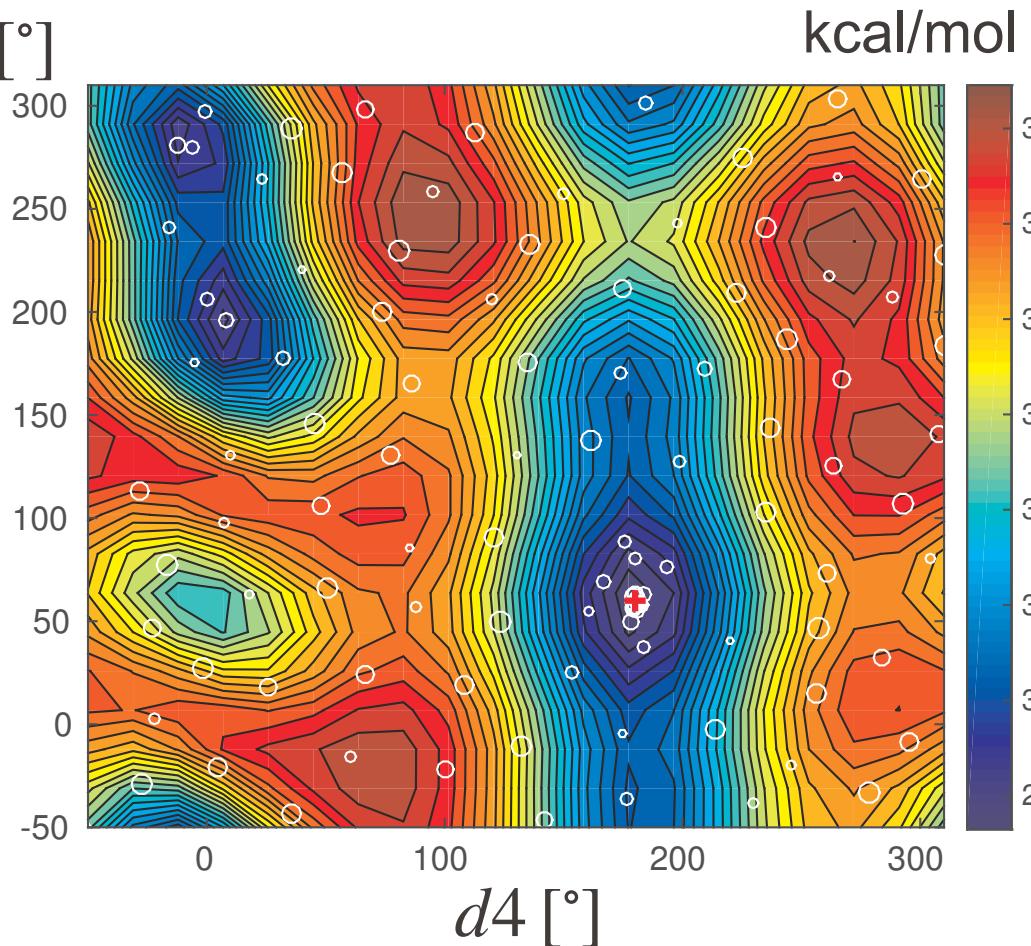


Potential Energy surface (PES)



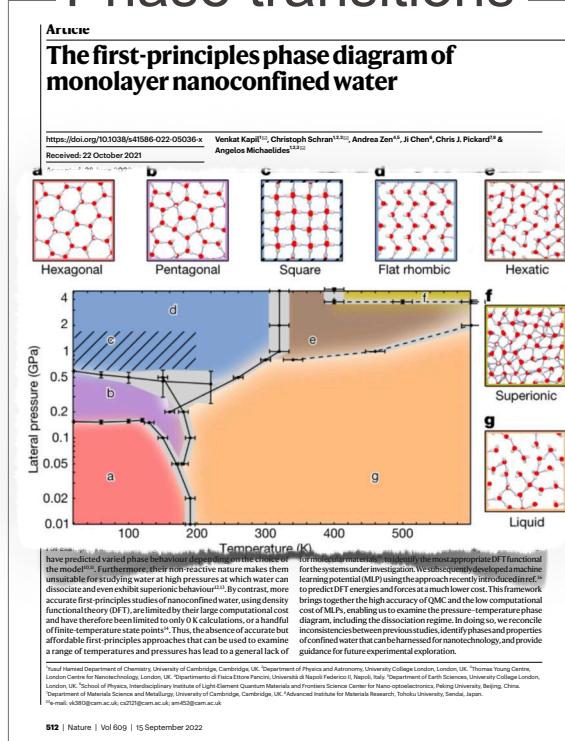
d_{13} [°]

- Total energy in terms of atomic positions
- The minima are the stable isomers/conformers

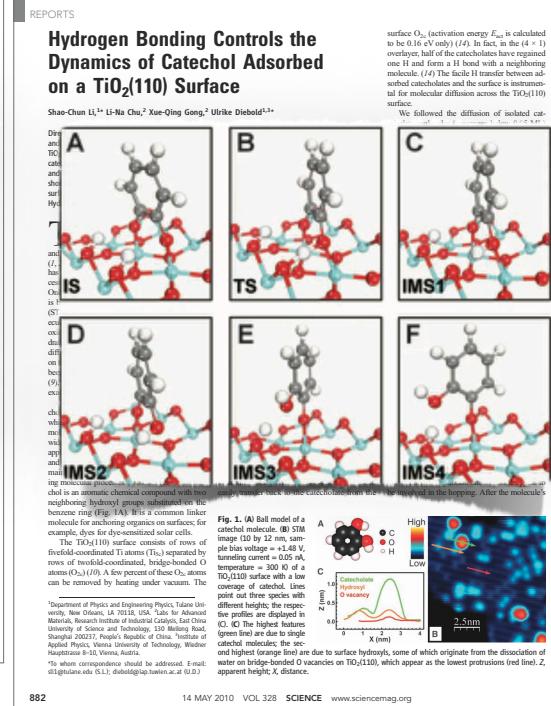


DFT: backbone of computational chemistry

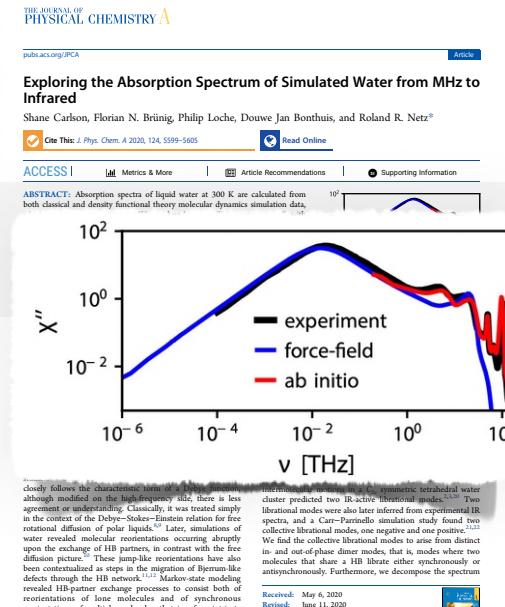
Phase transitions



Chemical reactions



Spectroscopy

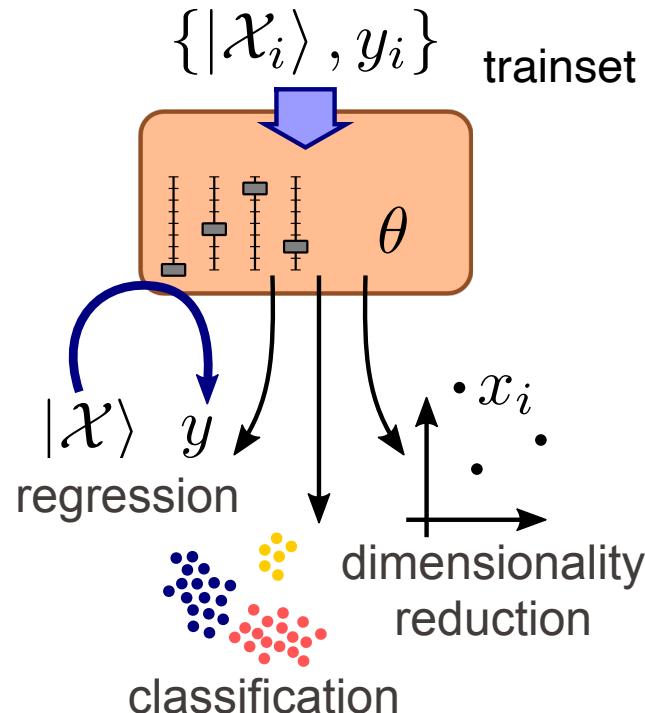


EPFL

Machine learning at the atomic scale

8

Chemical structures (inputs) and their properties (labels) are fed to a learning scheme, tuned by hyper parameters θ , that can then be used to perform different tasks on new data.



What is y_i (target)?

- per-structure property: energy
- per molecule property (dipole moment)
- per atom property (force)

What is X_i (features)?

- also called descriptor, fingerprint or representation
- *How to do good feature engineering?*

Data structure

Data is usually of type:

molecule1, energy1

molecule2, energy2

...

Molecules are then specified in xyz
coordinates:

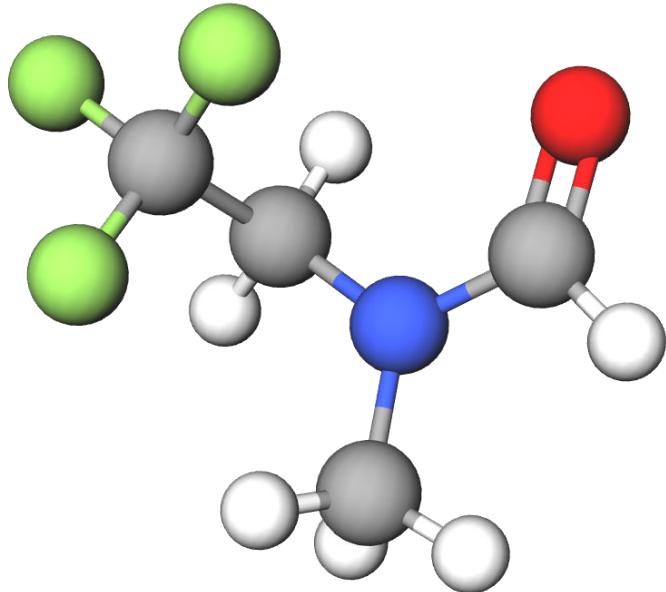
molecule1:

atom1, x1, y1, z1

atom2, x2, y2, z2

....

Additionally, properties (force, dipoles, etc.) might be listed.

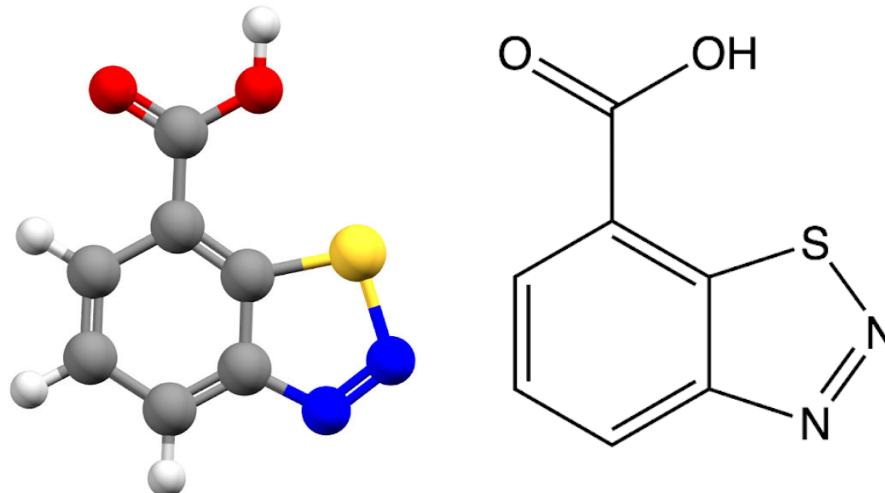


Molecular descriptors (cheminformatics)

Molecular representations

How would you describe this molecule to a machine?

Think about what molecular information could be extracted and how it could be encoded.



Acibenzolar acid
(benzo[d][1,2,3]thiadiazole-7-carboxylic acid)

Descriptors (representation)

A machine-compatible representation of a material,
state, or problem.

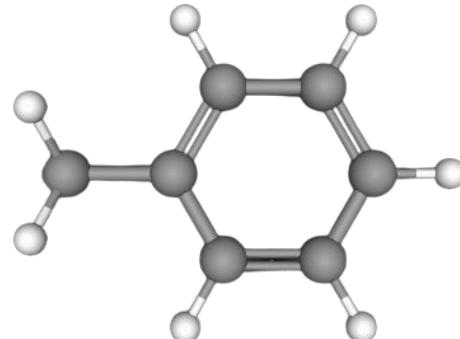
Attribute based

Encodes descriptive
information about the
problem.

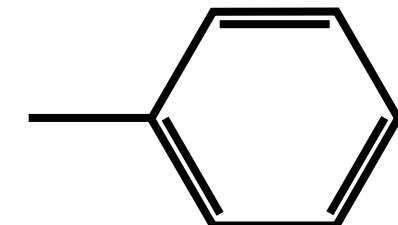
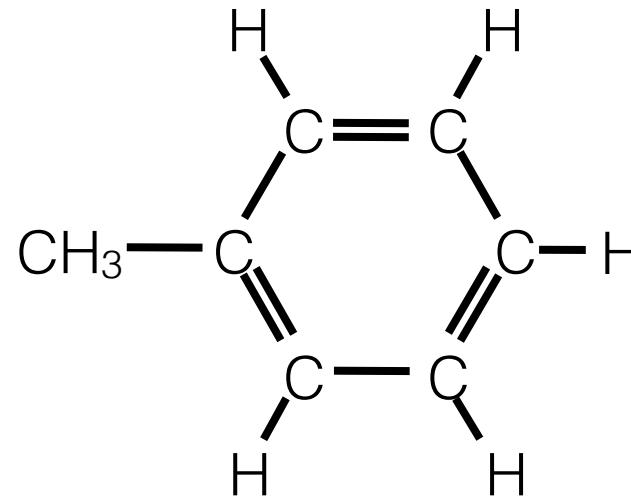
Structure based

Encodes physical or
chemical structure
directly.

Molecules as graphs



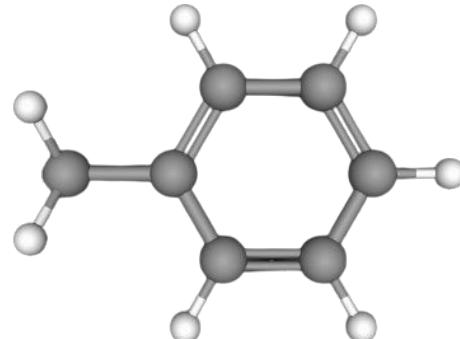
toluene



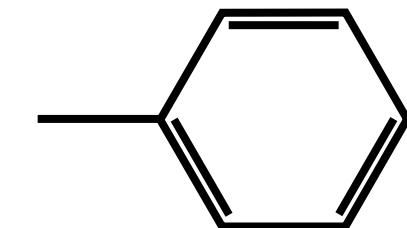
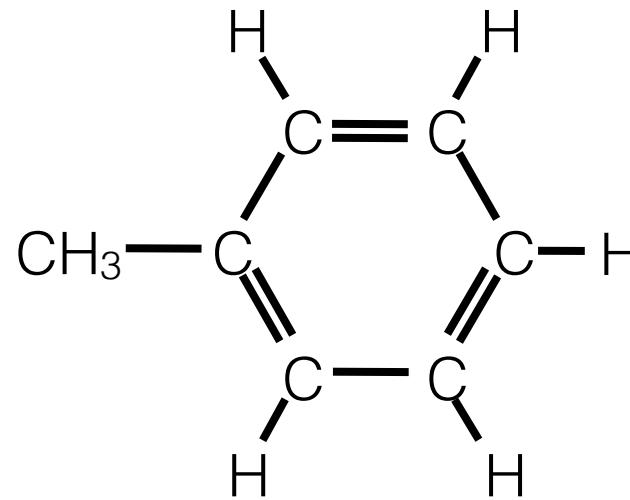
Graph representation of molecules (also called 2D)
Chemists are used to it

But we need something easier for the machines!

Molecules as graphs



toluene



C1=C(C)C=CC=C1

Simplified molecular-input line-entry system (*SMILES*)
text-based molecular encoding

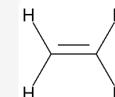
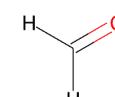
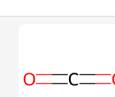
Atoms:

- as atomic symbol: [C]
- implicit fulfilled valency: C -> [CH₄]

Bonds:

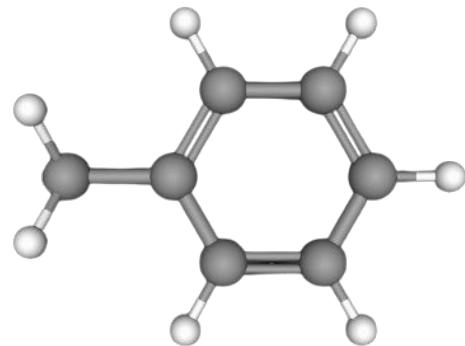
- single: -
- double: =
- triple: #
- Many more rules for branching, rings, charges...

- SMILES tools and generation: CANGEN algorithm, RDkit

Molecule	SMILES	Alternate SMILES	Structure
Ethene	C=C	[CH ₂][CH ₂]	
Formaldehyde	C=O	[CH ₂]=O	
Carbon Dioxide	O=C=O	C(=O)=O	

EPFL Fingerprints: molecules as bit strings

16



bit string: fingerprint
(one hot encoding)



Fingerprinting strategy:

- identify and extract molecular features
- encode them as 1 (present) or 0 (not present)
- or, hash the features
(= convert them to a number in a unique way)

EPFL

Fingerprints: molecules as bit strings

17

Also called dictionary fingerprints

Explicitly define features
(typically with SMILES Arbitrary Target Specification (SMARTS))

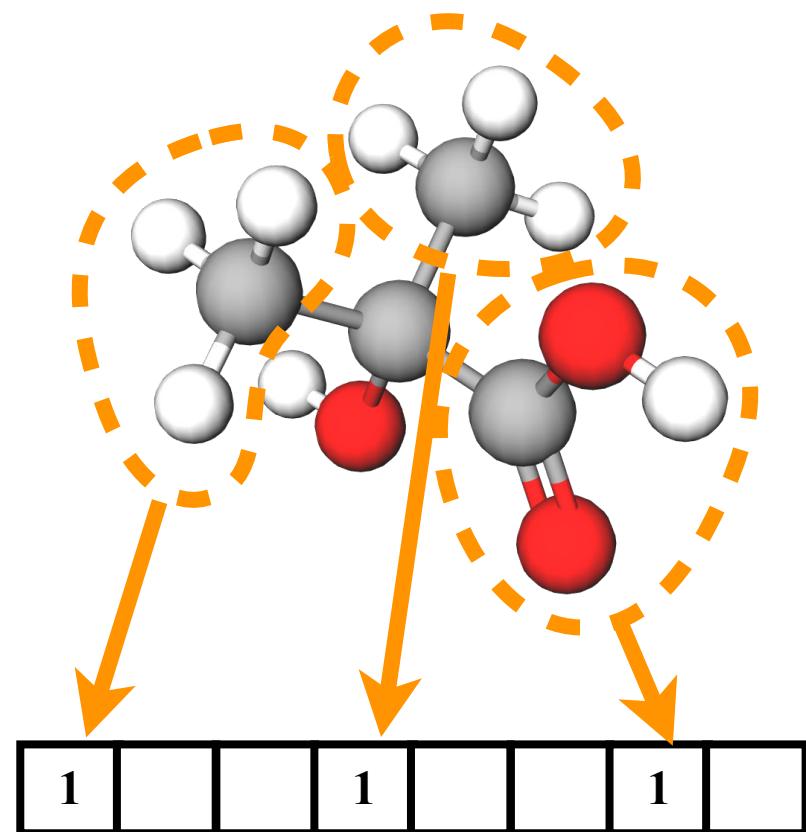
Check if features are present

Examples:

MACCS: 166 features

PubChem: 881 features

Klekota-Roth: 4860 features



EPFL Substructure fingerprints - advantages

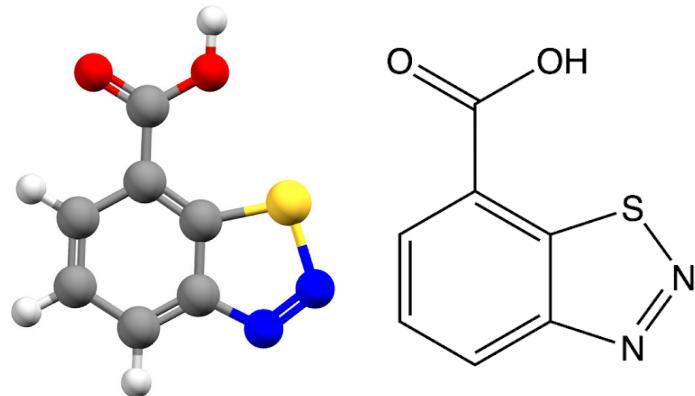
18

1. **Simple:**
easy to implement and use
2. **Compact:**
very small fingerprints, typically 2 to 8k entries
3. **Sasy similarity metrics:**
great for similarity searches
4. **Interpretability**
features encode chemical information directly

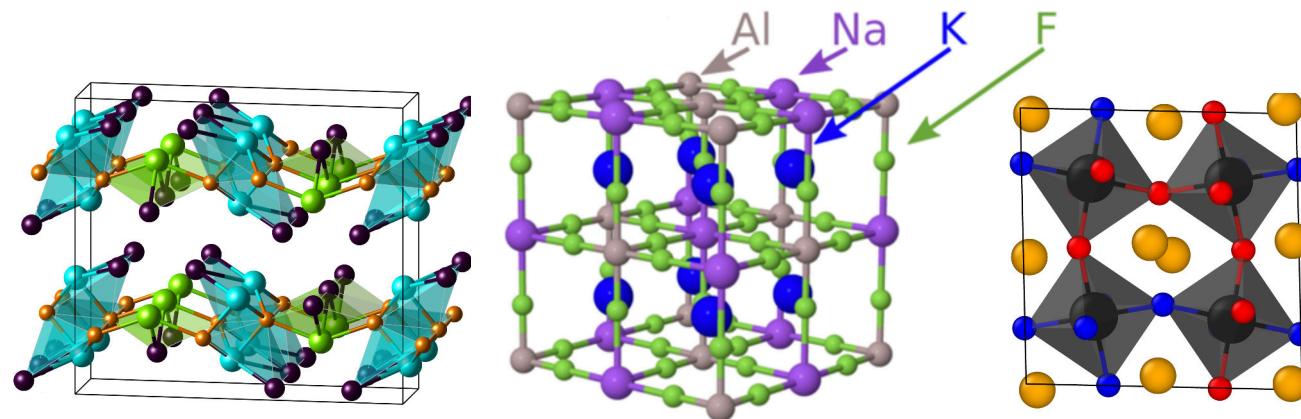
1. **Poor regression accuracy:**
bond lengths and bond angles are not encoded
2. **Lack of transferability:**
fixed feature selection might not translate to new domains
3. **Lack of resolution:**
certain molecules might have the same fingerprint due to limited feature selection
4. **Not conformer sensitive:**
since bond lengths and angle information is absent

Atomistic descriptors (condensed matter)

Representation for solids

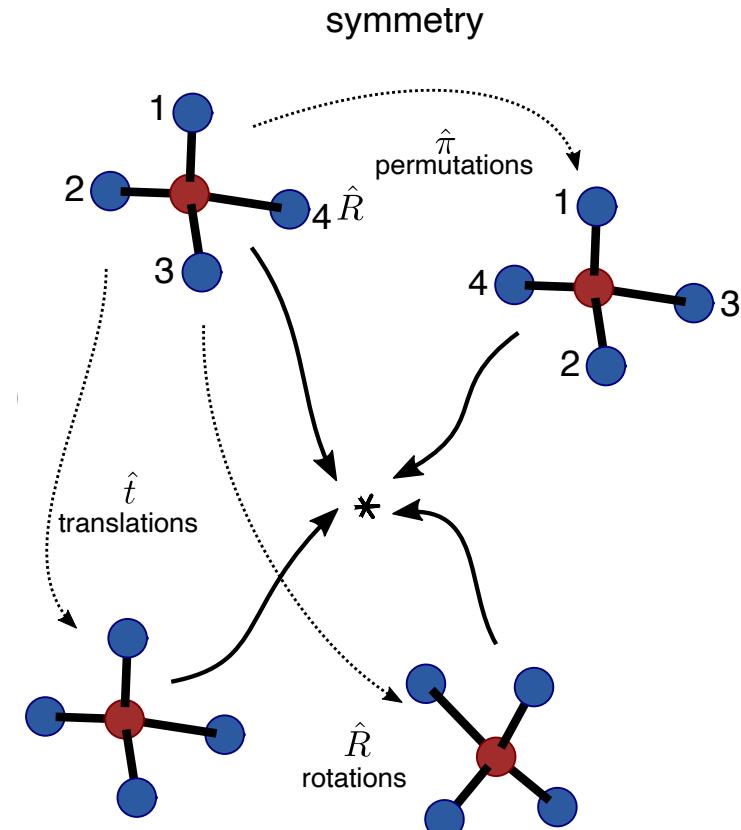


1. Speculate why molecular descriptors are not used to represent solids.
2. What could a descriptor for solids look like; what should it include?



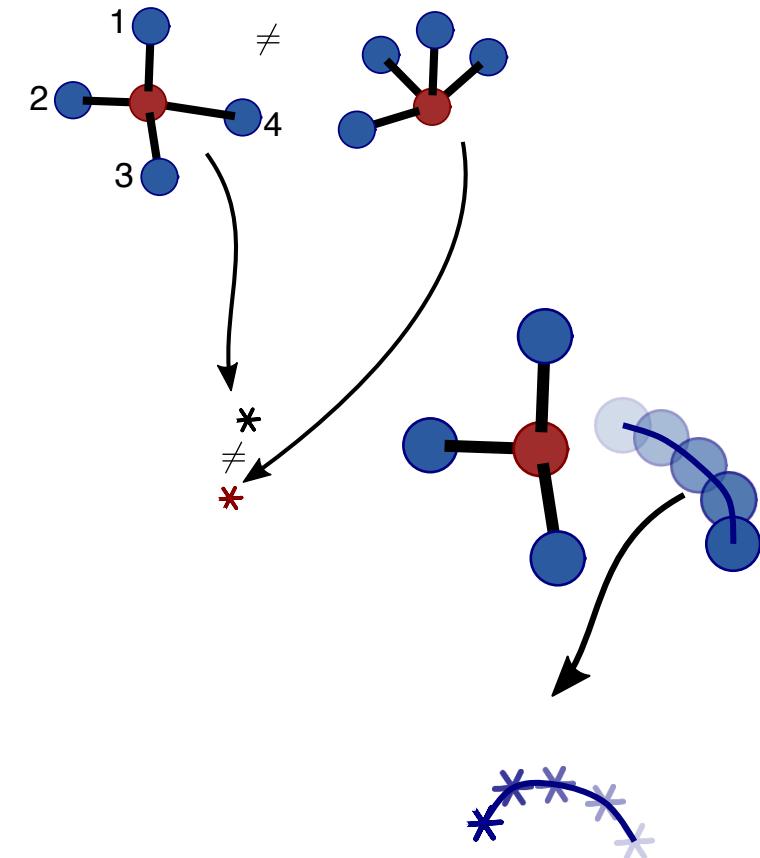
Descriptor requirements

1. Invariant with respect to spatial **translations** of the coordinate system: isometry of space
2. Invariant with respect to **rotations** of the coordinate system: isotropy of space.
3. Invariant with respect to **permutations** (e.g., elements in chemical formula, atomic indices)



Descriptor requirements

4. **Completeness/Unique**: there is a single way to construct a descriptor for the material, problem or state
5. **Continuous**: small changes in the state should translate to small changes in the descriptor
6. Compact and **computationally efficient**

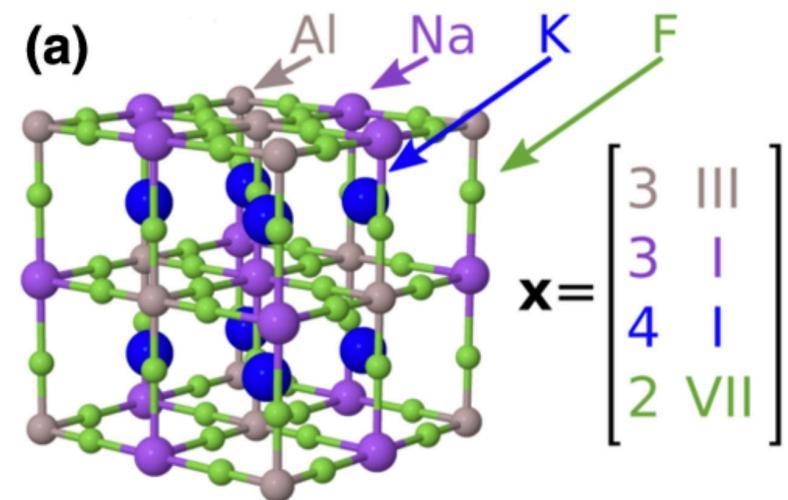


Attribute based descriptors

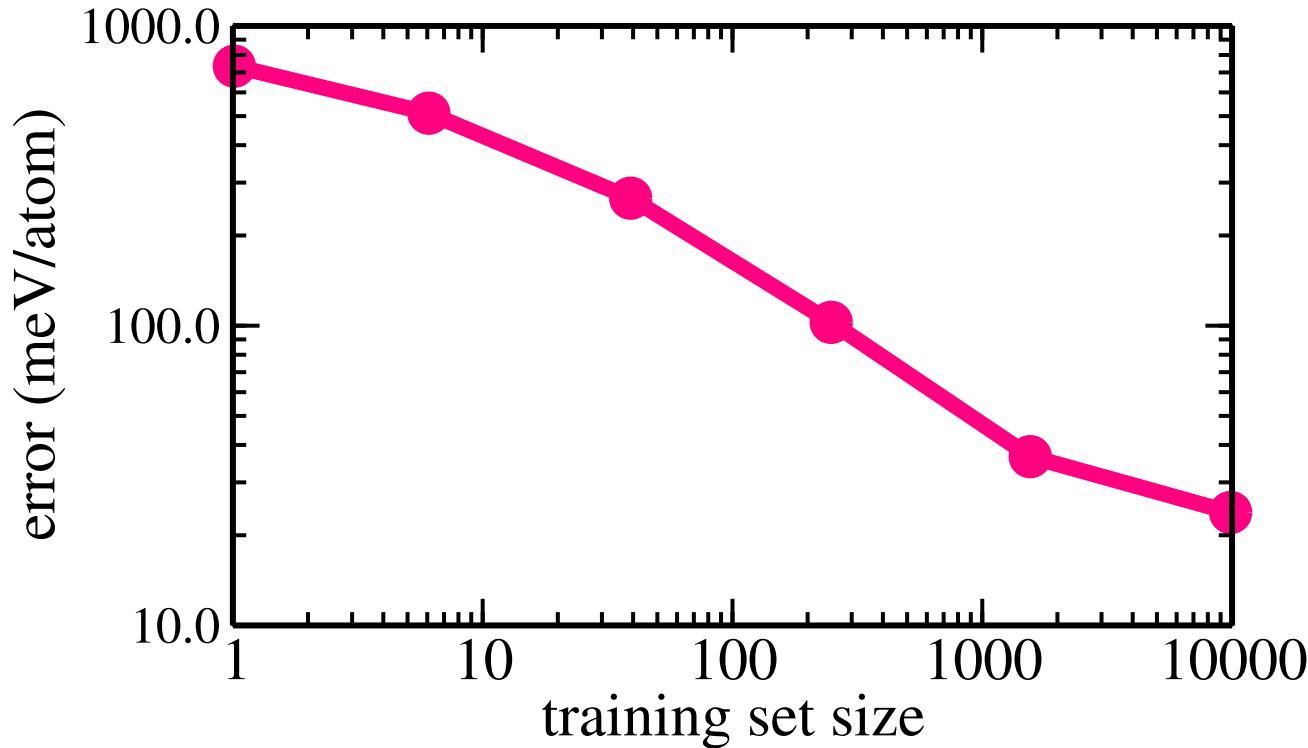
- Easy example: fixed crystal structure
- Elpasolite crystal structure (ABC_2D_6) is very common
- $2 \cdot 10^6$ possible combinations with only main group elements
- Cohesive energy for $\sim 10'000$ materials calculated with DFT

Descriptor

For each lattice position in the structure (here 4), the atom type is encoded by its column and row number in the periodic table



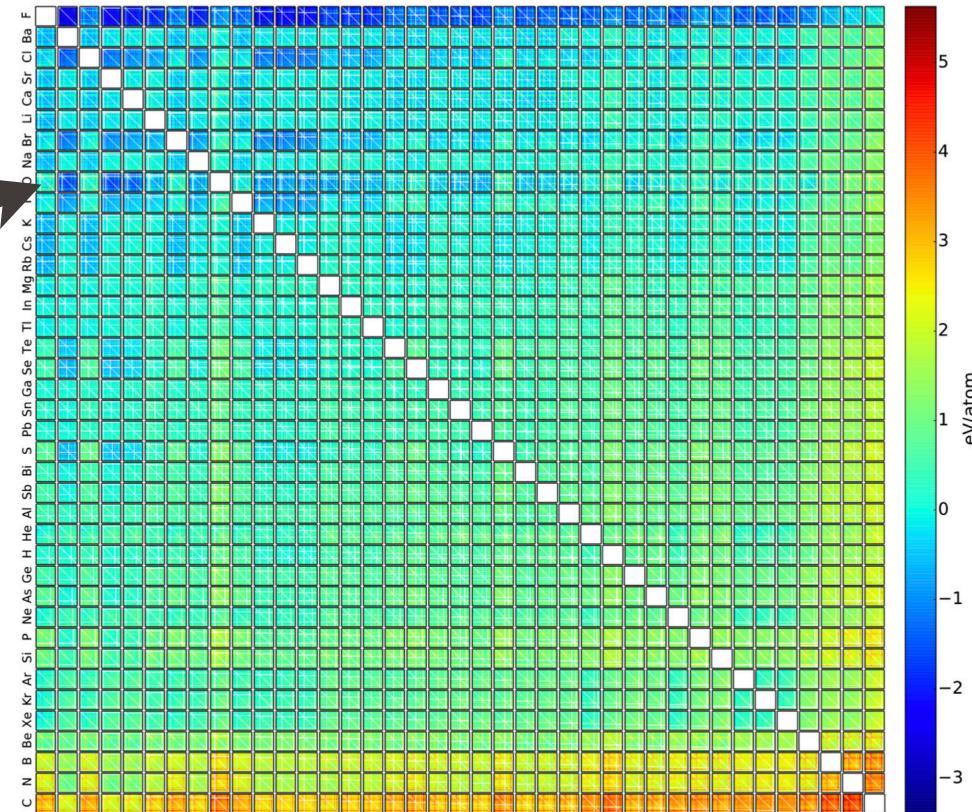
Application: ABC2D6 elpasolite crystals



Inference: Predictions for 2M crystals

Each square is
a 39×39
block of the
3rd and 4th
element

First 2
elements



Problem: How do we present a list of general materials to a machine?

ZnO, Cu₂O₃, YZrO₂, FeCrMnNiCo, CsPbCl₃, BaTiO₂, H₂O etc.

More general solution:

Use the periodic table to expand chemical formulas into descriptor and supplement with data from **material handbooks** and electronic structure theory calculations.

Attribute based descriptors

Atomic Number	Symbol	Name	Atomic Mass
1	H	Hydrogen	1.008
3	Li	Lithium	6.941
4	Be	Beryllium	9.012
11	Na	Sodium	22.990
12	Mg	Magnesium	24.305
19	K	Potassium	39.098
20	Ca	Calcium	40.078
21	Sc	Scandium	44.956
22	Ti	Titanium	47.867
23	V	Vanadium	50.942
24	Cr	Chromium	51.996
25	Mn	Manganese	54.938
26	Fe	Iron	55.845
27	Co	Cobalt	58.933
28	Ni	Nickel	58.693
29	Cu	Copper	63.546
30	Zn	Zinc	65.38
31	Ga	Gallium	69.723
32	Ge	Germanium	72.631
33	As	Arsenic	74.922
34	Se	Selenium	78.972
35	Br	Bromine	79.904
36	Kr	Krypton	83.798
37	Rb	Rubidium	85.468
38	Sr	Strontium	87.62
39	Y	Yttrium	88.906
40	Zr	Zirconium	91.224
41	Nb	Niobium	92.906
42	Mo	Molybdenum	95.95
43	Tc	Technetium	98.907
44	Ru	Ruthenium	101.07
45	Rh	Rhodium	102.906
46	Pd	Palladium	106.42
47	Ag	Silver	107.868
48	Cd	Cadmium	112.411
49	In	Indium	114.818
50	Sn	Tin	118.711
51	Sb	Antimony	121.760
52	Te	Tellurium	127.6
53	I	Iodine	126.904
54	Xe	Xenon	131.294
55	Cs	Cesium	132.905
56	Ba	Barium	137.328
57-71			
72	Hf	Hafnium	178.49
73	Ta	Tantalum	180.948
74	W	Tungsten	183.84
75	Re	Rhenium	186.207
76	Os	Osmium	190.23
77	Ir	Iridium	192.217
78	Pt	Platinum	195.085
79	Au	Gold	196.967
80	Hg	Mercury	200.592
81	Tl	Thallium	204.383
82	Pb	Lead	207.2
83	Bi	Bismuth	208.980
84	Po	Polonium	[208.982]
85	At	Astatine	209.987
86	Rn	Radon	222.018
87	Fr	Francium	223.020
88	Ra	Radium	226.025
89-103			
104	Rf	Rutherfordium	[261]
105	Db	Dubnium	[262]
106	Sg	Seaborgium	[266]
107	Bh	Bohrium	[264]
108	Hs	Hassium	[269]
109	Mt	Meltnerium	[278]
110	Ds	Darmstadtium	[281]
111	Rg	Roentgenium	[280]
112	Cn	Copernicium	[285]
113	Nh	Nihonium	[286]
114	Fl	Flerovium	[289]
115	Mc	Moscovium	[289]
116	Lv	Livermorium	[293]
117	Ts	Tennessee	[294]
118	Og	Oganesson	[294]

57	La	Lanthanum	138.905	58	Ce	Cerium	140.116	59	Pr	Praseodymium	140.908	60	Nd	Neodymium	144.242	61	Pm	Promethium	144.913	62	Sm	Samarium	150.36	63	Eu	Europium	151.964	64	Gd	Gadolinium	157.25	65	Tb	Terbium	158.925	66	Dy	Dysprosium	162.500	67	Ho	Holmium	164.930	68	Er	Erbium	167.259	69	Tm	Thulium	168.934	70	Yb	Ytterbium	173.055	71	Lu	Lutetium	174.967
89	Ac	Actinium	227.028	90	Th	Thorium	232.038	91	Pa	Protactinium	231.036	92	U	Uranium	238.029	93	Np	Neptunium	237.048	94	Pu	Plutonium	244.064	95	Am	Americium	243.061	96	Cm	Curium	247.070	97	Bk	Berkelium	247.070	98	Cf	Californium	251.080	99	Es	Einsteinium	[254]	100	Fm	Fermium	257.095	101	Mendelevium	258.1	102	No	Nobelium	259.101	103	Lr	Lawrencium	[262]	



EPFL Attribute-based descriptors: MAGPIE

29

Elemental Property Stats (115):

6 statistics: mean, variance, max, min, range, mode; of 22 elemental properties: Z, row, column, radius, electronegativity, s, p, d... valence electrons, unfilled states, magnetic moment....

Stoichiometry:

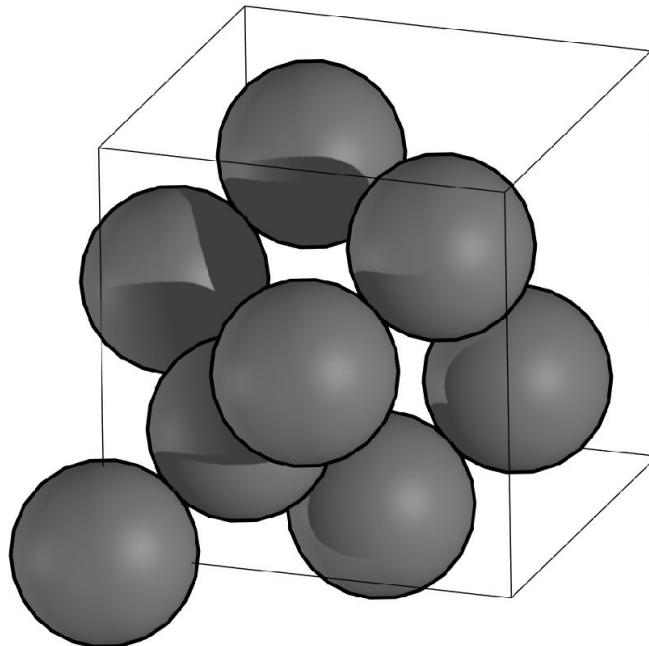
no. of components (fractions of elements)

Electronic-structure based:

fractions of s, p, d... electrons in the material

1. **Simple:**
easy to implement and use
2. **Compact:**
very, very small fingerprints
3. **Invariant:**
invariances to rotation and translation built in
4. **Interpretability**
features encode information on chemical elements and certain properties directly

1. **poor regression accuracy:**
bond lengths and bond angles are not encoded
2. **no structural information:**
polymorphs of the same material cannot be resolved
3. **lack of transferability:**
fixed feature selection might not translate to new domains

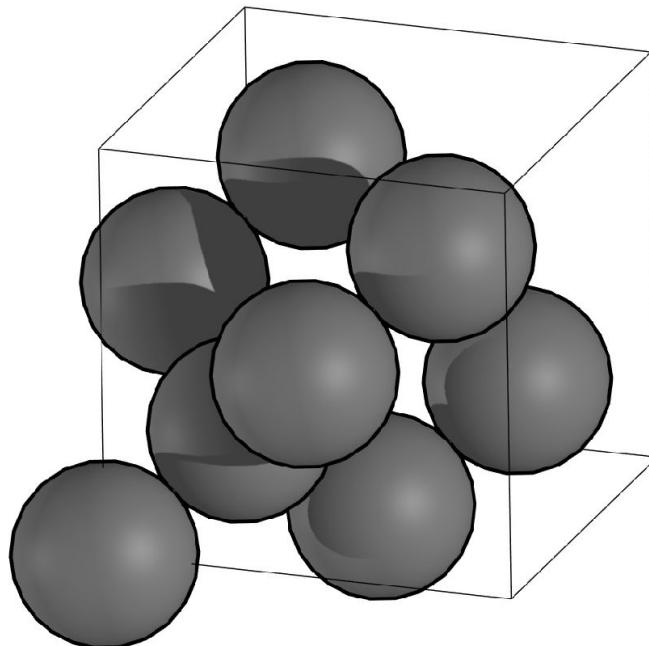


diamond

$$\mathbf{M}_{IJ} = \begin{cases} 0.5Z_I^{2.4} & \text{for } I = J \\ \frac{Z_I Z_J}{|\mathbf{R}_I - \mathbf{R}_J|} & \text{for } I \neq J \end{cases}$$

↑
atomic positions

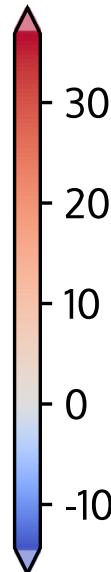
atomic charges

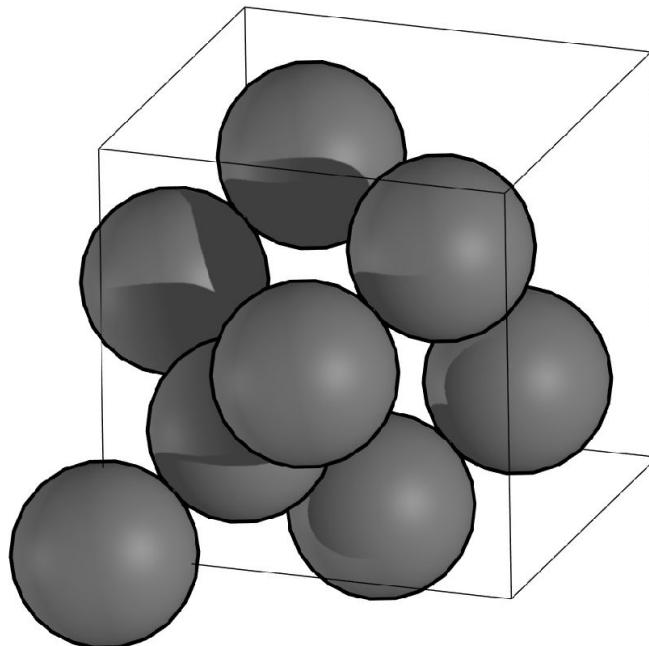


diamond

not periodic

36.9	23.3	14.3	9.3	14.3	9.3	14.3	9.3
23.3	36.9	23.3	14.3	23.3	14.3	23.3	14.3
14.3	23.3	36.9	23.3	14.3	12.2	14.3	12.2
9.3	14.3	23.3	36.9	12.2	14.3	12.2	14.3
14.3	23.3	14.3	12.2	36.9	23.3	14.3	12.2
9.3	14.3	12.2	14.3	23.3	36.9	12.2	14.3
14.3	23.3	14.3	12.2	14.3	12.2	36.9	23.3
9.3	14.3	12.2	14.3	12.2	14.3	23.3	36.9





diamond

Electrostatic energy between atoms i and j :

$$\phi_{ij} = \sum_{\mathbf{n}} \frac{Z_i Z_j}{|\mathbf{R}_i - \mathbf{R}_j| + \mathbf{n}}$$

Infinite lattice sum:

$$\mathbf{n} = h\mathbf{a} + k\mathbf{b} + l\mathbf{c}$$

You will learn tomorrow
how to evaluate!

Simplification: Sine matrix

$$M_{ij}^{\text{sine}} = \begin{cases} 0.5Z_i^{2.4} & \forall i = j \\ \phi_{ij} & \forall i \neq j \end{cases}$$

$$\phi_{ij} = Z_i Z_j | \mathbf{B} \cdot \sum_{k=\{x,y,z\}} \hat{\mathbf{e}}_k \sin^2 (\pi \mathbf{B}^{-1} \cdot (\mathbf{R}_i - \mathbf{R}_j)) |^{-1}$$

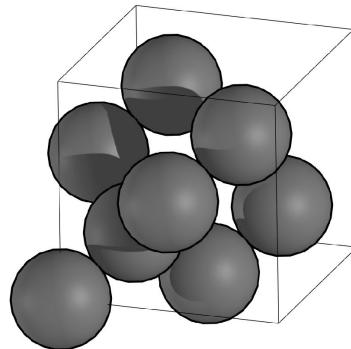


matrix formed by
lattice vectors



Cartesian unit vectors

Comparison of the three matrices



Coulomb matrix

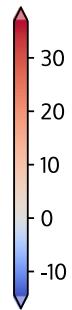
36.9	23.3	14.3	9.3	14.3	9.3	14.3	9.3	9.3
23.3	36.9	23.3	14.3	23.3	14.3	23.3	14.3	14.3
14.3	23.3	36.9	23.3	14.3	12.2	14.3	12.2	12.2
9.3	14.3	23.3	36.9	12.2	14.3	12.2	14.3	14.3
14.3	23.3	14.3	12.2	36.9	23.3	14.3	12.2	12.2
9.3	14.3	12.2	14.3	23.3	36.9	12.2	14.3	9.3
14.3	23.3	14.3	12.2	14.3	12.2	36.9	23.3	9.3
9.3	14.3	12.2	14.3	12.2	14.3	23.3	36.9	36.9

Ewald sum matrix

-14.3	-2.0	-5.9	-2.0	-5.9	-2.0	-5.9	-2.0
-2.0	-14.3	-2.0	-5.9	-2.0	-5.9	-2.0	-5.9
-5.9	-2.0	-14.3	-2.0	-5.9	-2.0	-5.9	-2.0
-2.0	-5.9	-2.0	-14.3	-2.0	-5.9	-2.0	-5.9
-5.9	-2.0	-5.9	-2.0	-14.3	-2.0	-5.9	-2.0
-2.0	-5.9	-2.0	-5.9	-2.0	-14.3	-2.0	-5.9
-5.9	-2.0	-5.9	-2.0	-5.9	-2.0	-14.3	-2.0
-2.0	-5.9	-2.0	-5.9	-2.0	-5.9	-2.0	-14.3

Sine matrix

36.9	11.7	7.1	11.7	7.1	11.7	7.1	11.7	11.7
11.7	36.9	11.7	7.1	11.7	7.1	11.7	7.1	11.7
7.1	11.7	36.9	11.7	7.1	11.7	7.1	11.7	11.7
11.7	7.1	11.7	36.9	11.7	7.1	11.7	7.1	11.7
7.1	11.7	7.1	11.7	36.9	11.7	7.1	11.7	11.7
11.7	7.1	11.7	7.1	11.7	36.9	11.7	7.1	7.1
7.1	11.7	7.1	11.7	7.1	11.7	36.9	11.7	7.1
11.7	7.1	11.7	7.1	11.7	7.1	11.7	36.9	11.7
7.1	11.7	7.1	11.7	7.1	11.7	7.1	11.7	36.9



not periodic



periodic



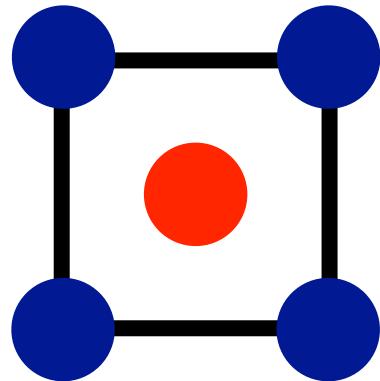
EPFL Ewald and sine matrices - advantages

37

1. **Simple:**
easy to implement and use
2. **Compact:**
small descriptor
3. **Invariant:**
invariances to rotation and translation built in
4. **Quantitative structure information**
bond lengths now encoded

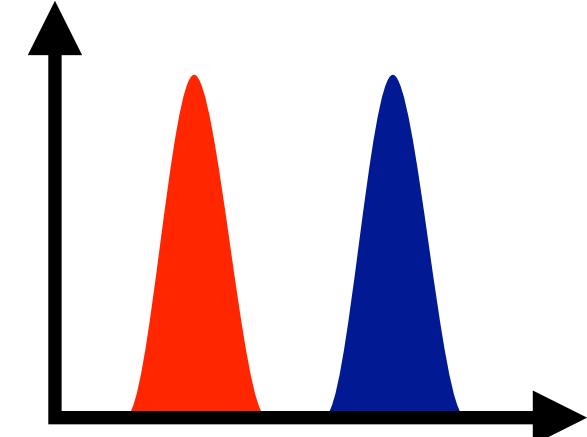
1. **Moderate regression accuracy:**
structural information still limited
2. **Not continuous:**
small changes in the structure can switch the rows and columns
3. **Lack of interpretability:**
features correspond to inverse distances, but are becoming hard to interpret
4. **Not transferable**
Descriptor size changes with number of atoms

From global to local



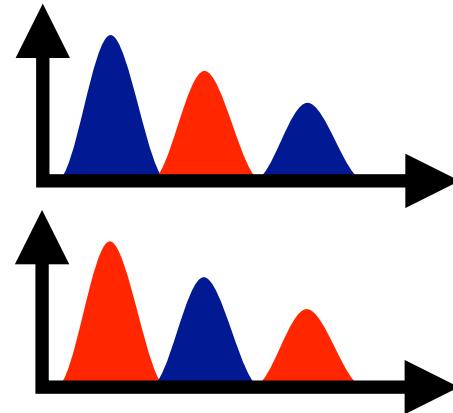
global

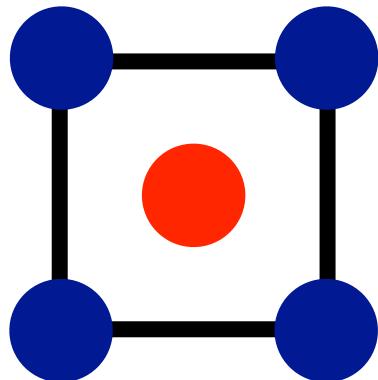
One descriptor for the
whole structure



local

One descriptor for
the **every atom**





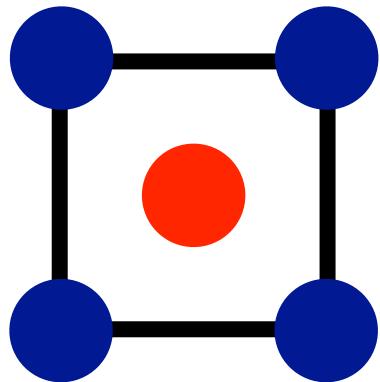
Locality principle:
many properties are local
or approximately additive

The total energy in DFT is not additive!

But we treat it as such approximately:

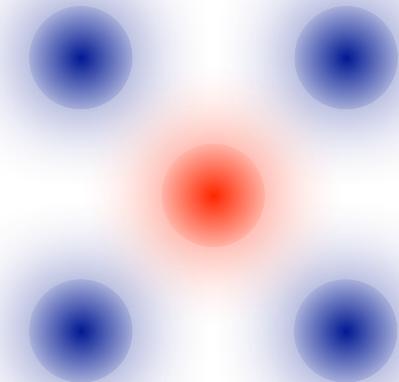
$$E = \sum_i^{N_{atoms}} E_i$$

Smooth overlap of atomic orbitals (SOAP)



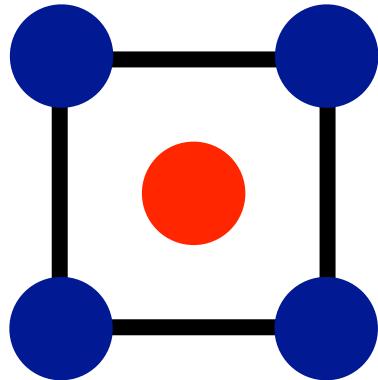
Gaussian broadening
atomic positions into a
density field.

“Density trick”

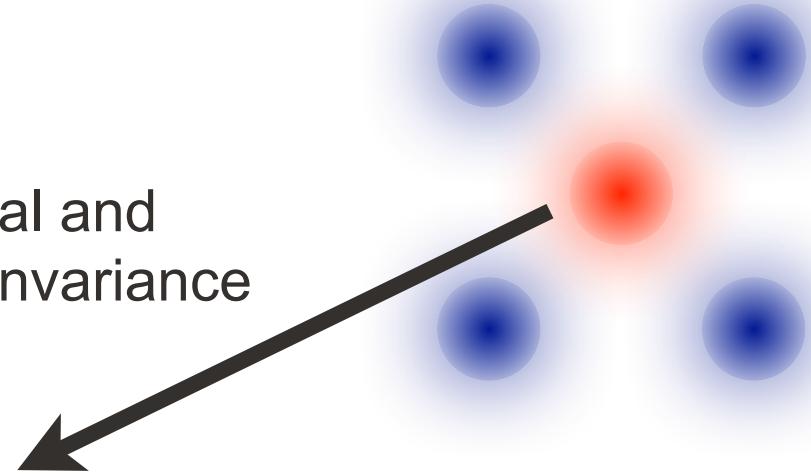


$$\rho_i(\mathbf{r}) = \sum_j e^{-\frac{1}{2\sigma^2} |\mathbf{r} - \mathbf{r}_{ij}|^2}$$
$$\mathbf{r}_{ij} = \mathbf{R}_i - \mathbf{R}_j$$

Smooth overlap of atomic orbitals (SOAP)



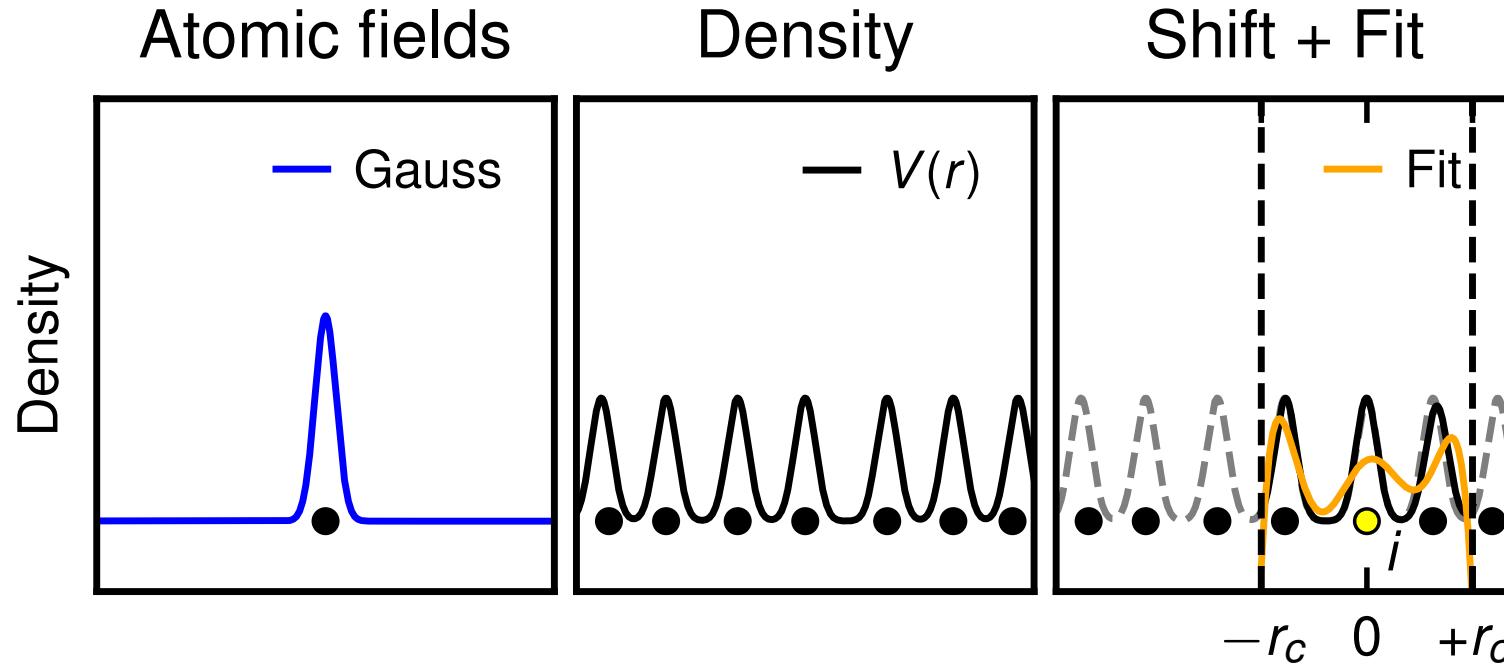
translational and
rotational invariance



Expand in on a basis spherical harmonics:

$$\rho_i(\mathbf{r}) = \sum_{nlm} c_{i,nlm} R_n(r) Y_{lm}(\hat{\mathbf{r}})$$

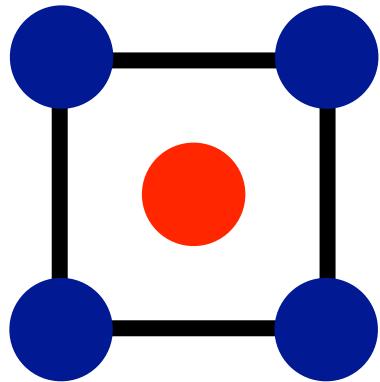
Smooth overlap of atomic orbitals (SOAP)



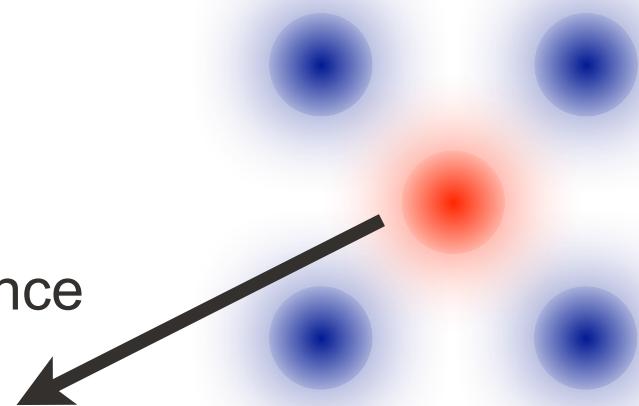
Project/fit the density on a basis built from radial & angular functions.

Take the coefficients of that projection as pairwise representation
also called spherical expansion $\rho(r)$.

Smooth overlap of atomic orbitals (SOAP)



translational and
rotational invariance



Combine spherical expansions and take zeroth order term
of for power spectrum:

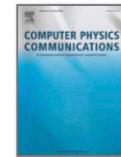
final SOAP descriptor →
$$p_{nn'l}^{Z_1, Z_2} = \pi \sqrt{\frac{8}{2l+1}} \sum_m \left(c_{nlm}^{Z_1} \right)^* c_{n'l'm}^{Z_2}$$

In practice, implementation of atomic descriptors
is a bit painful... 😞



Computer Physics Communications

Volume 247, February 2020, 106949



DSScribe: Library of descriptors for machine learning in materials science ☆

Lauri Himanen ^a✉, Marc O.J. Jäger ^a, Eiaki V. N.
David Z. Gao ^{b, c}, Patrick Rinke ^{a, f}, Adam S. Foste

✉ Show more

<https://doi.org/10.1016/j.cpc.2019.106949>

Under a Creative Commons license

metatensor and metatomic: foundational libraries for interoperable atomistic machine learning

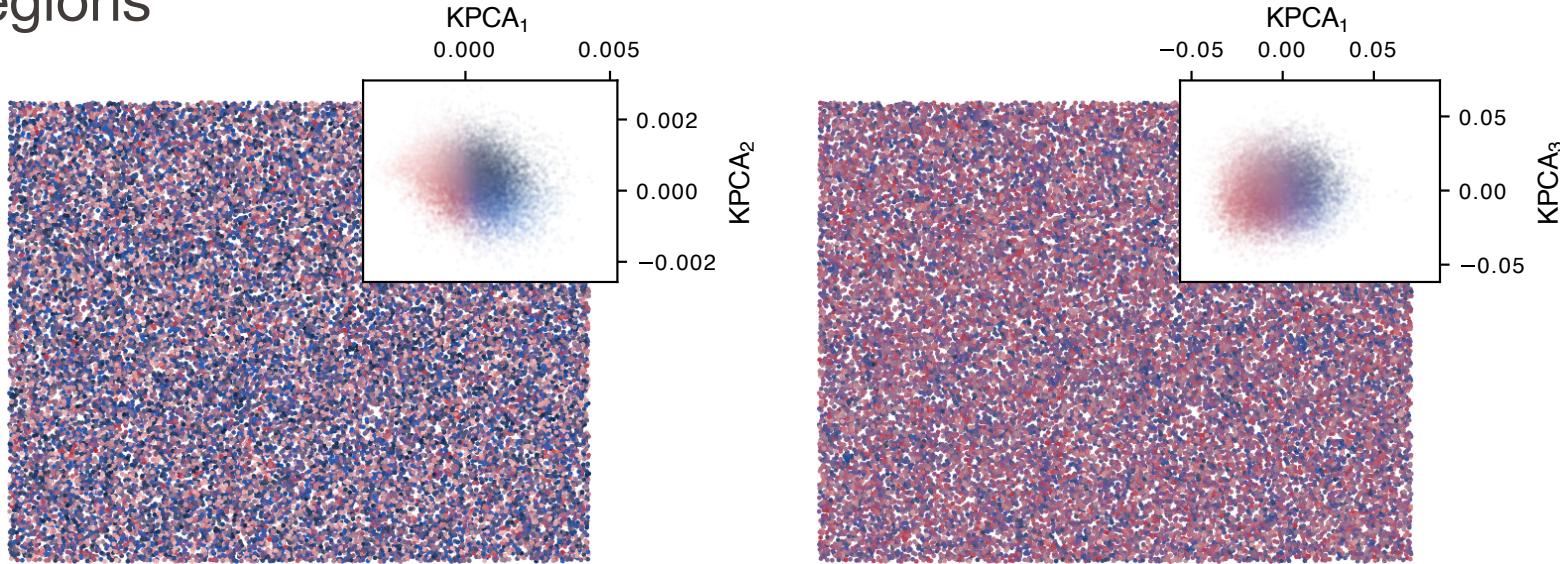
Filippo Bigi,^[a] Joseph W. Abbott,^[a] Philip Loche, Arslan Mazitov, Davide Tisi, Marcel F. Langer, Alexander Goscinski, Paolo Pegolo, Sanggyu Chong, Rohit Goswami, Sofiia Chorna, Matthias Kellner, Michele Ceriotti, and Guillaume Fraux

*Laboratory of Computational Science and Modeling, Institute of Materials,
École Polytechnique Fédérale de Lausanne, 1015 Lausanne, Switzerland*

Incorporation of machine learning (ML) techniques into atomic-scale modeling has proven to be an extremely effective strategy to improve the accuracy and reduce the computational cost of simulations. It also entails conceptual and practical challenges, as it involves combining very different mathematical

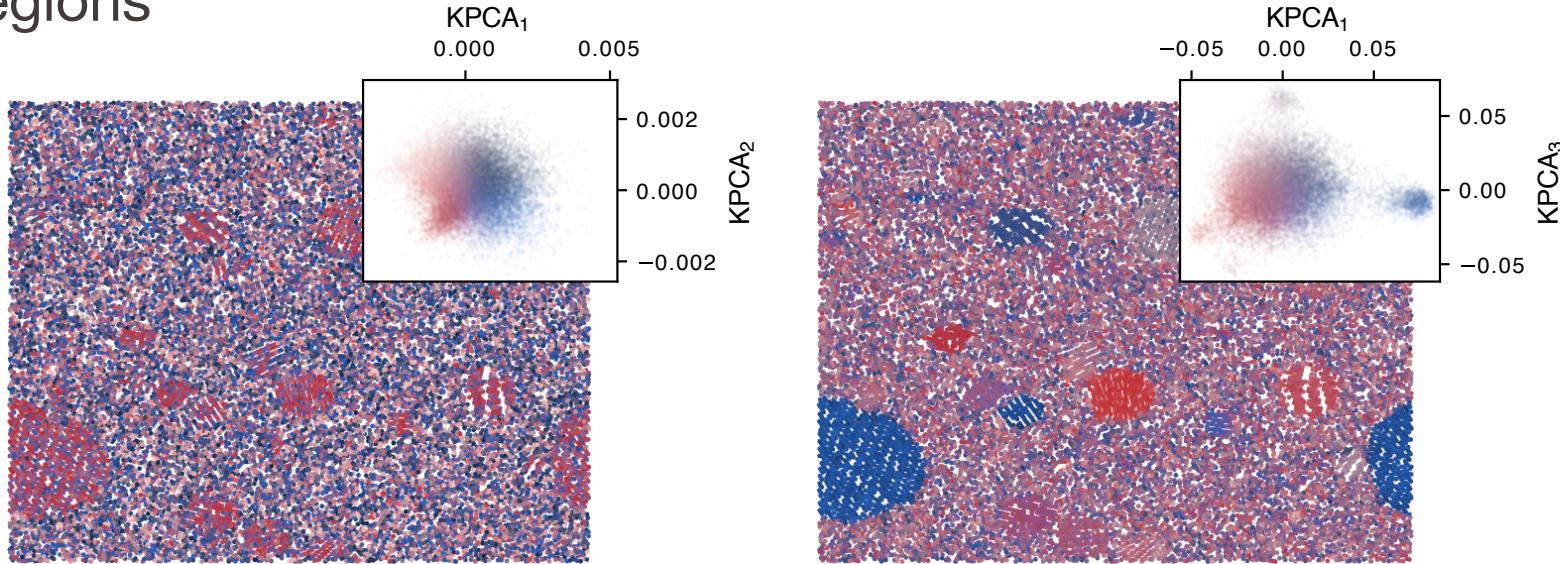
Application: Supercooled iron

Color atoms based on whether they are in liquid or solid regions



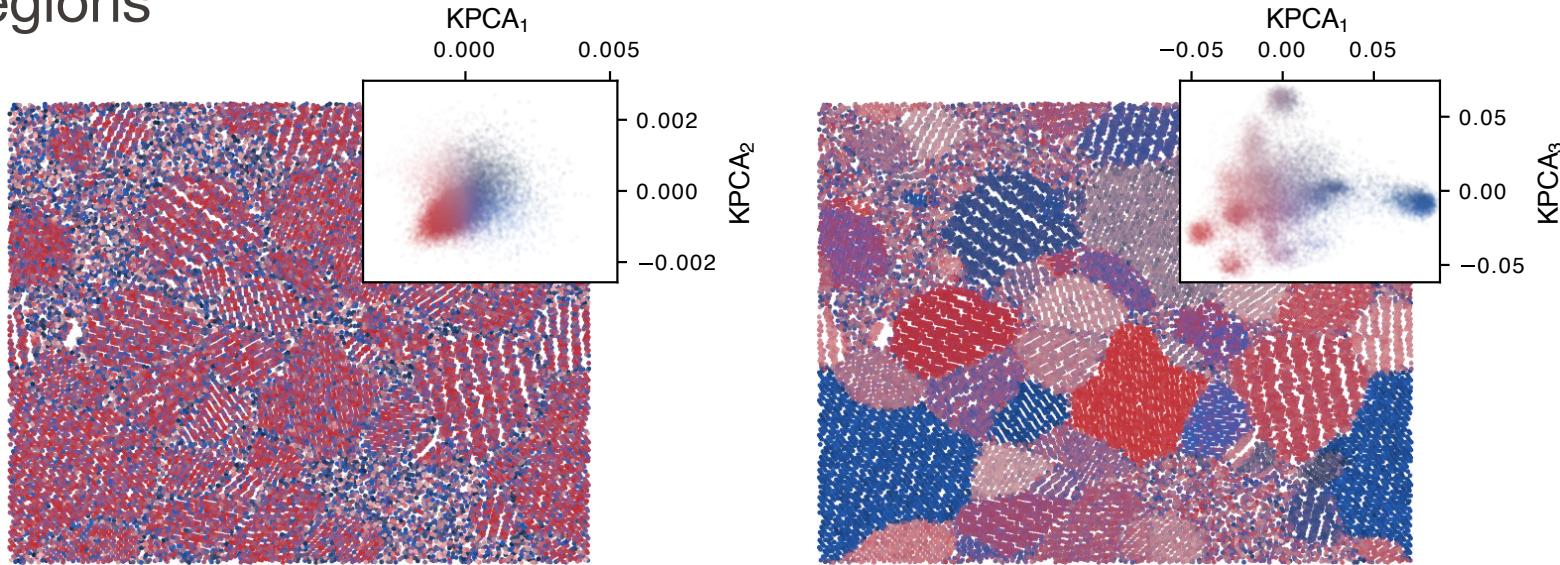
Application: Supercooled iron

Color atoms based on whether they are in liquid or solid regions



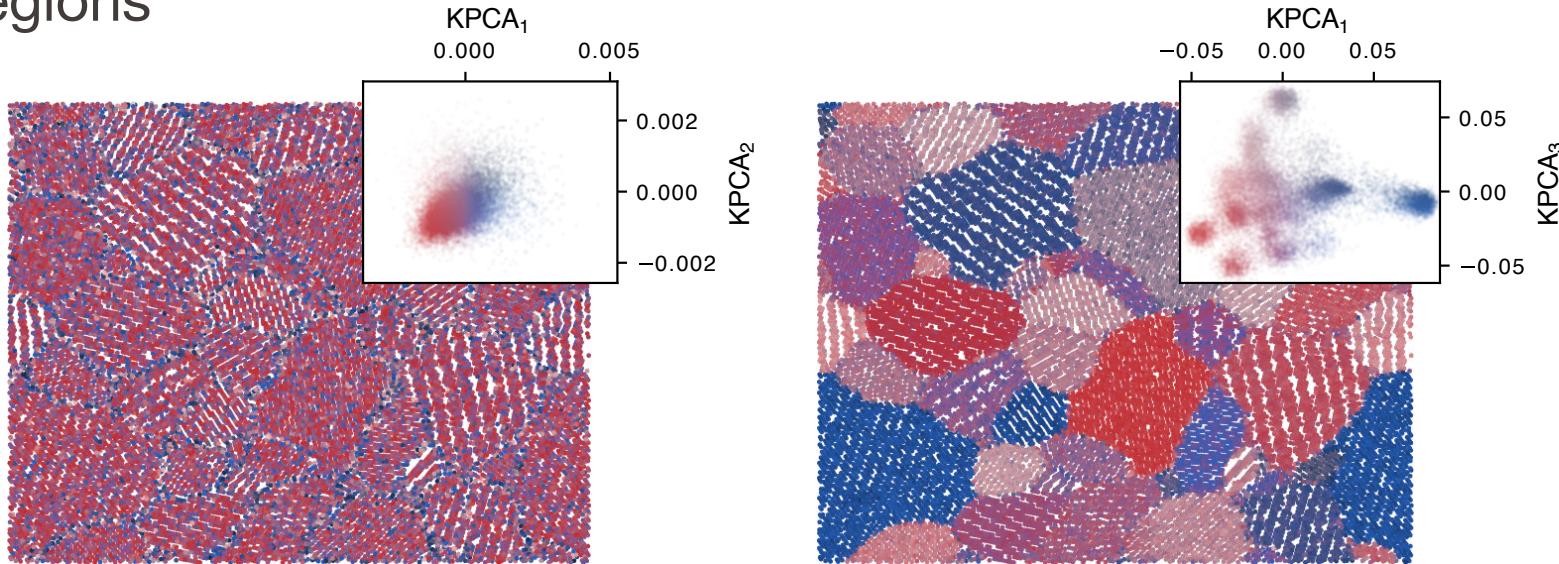
Application: Supercooled iron

Color atoms based on whether they are in liquid or solid regions



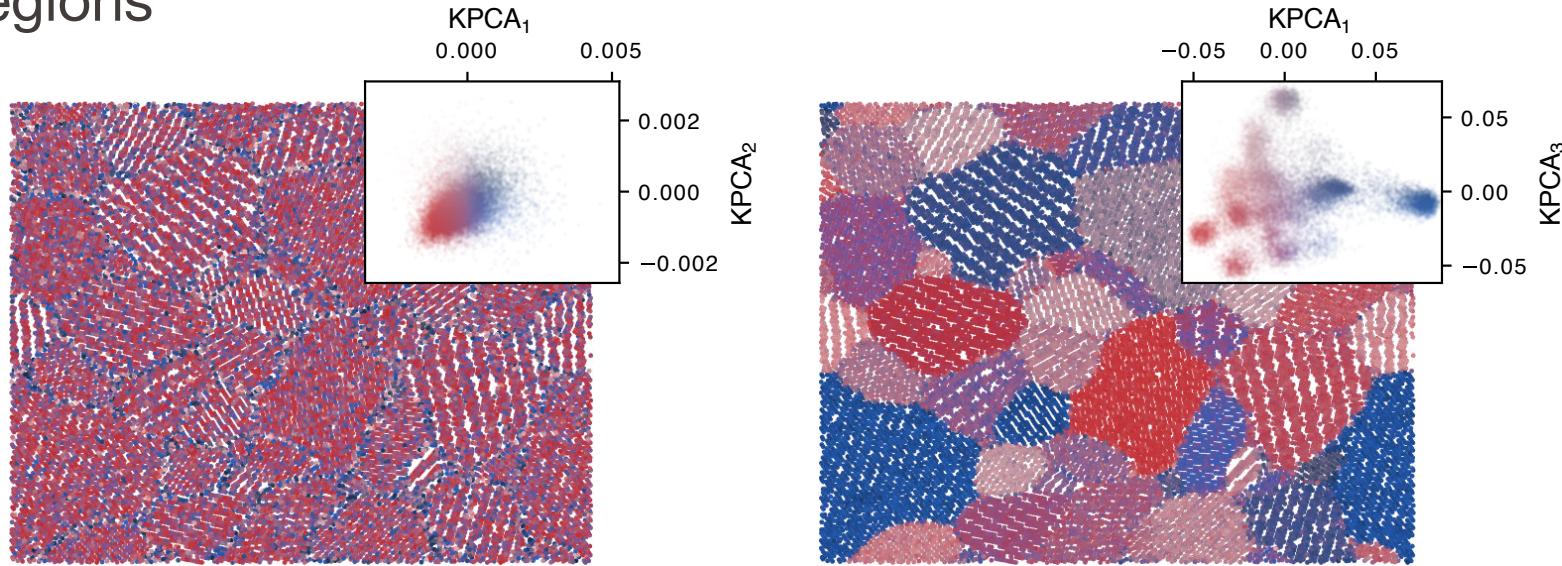
Application: Supercooled iron

Color atoms based on whether they are in liquid or solid regions



Application: Supercooled iron

Color atoms based on whether they are in liquid or solid regions

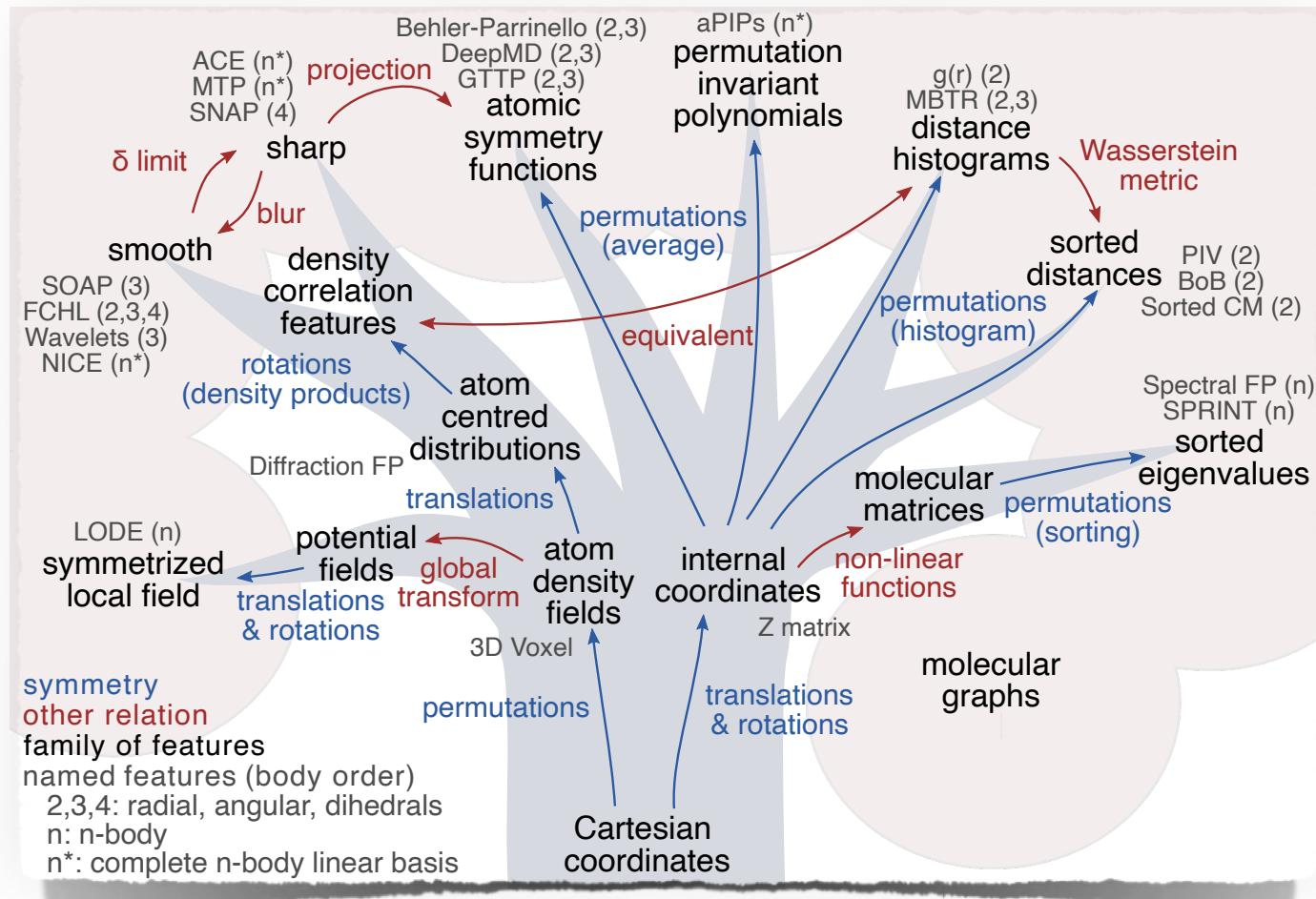


Rotational invariant

"Rotational non-invariant"

Understanding what goes into a representation is key to achieve meaningful results from automated data analytics

There are MANY related representations



Local descriptors - advantages

1. Invariant:

invariances to rotation and translation built in

2. Quantitative structure information

bond lengths, angles and higher orders encoded

3. Continuity

continuous for small structural changes

4. Locality

can resolve local environments (unlike global descriptors)

5. Accuracy

among the most accurate descriptors for regression

6. Versatility

works for molecules and periodic structures

1. **Hyperparameters:**

local descriptors often have many hyperparameters that need to be tuned, which is cumbersome and timely

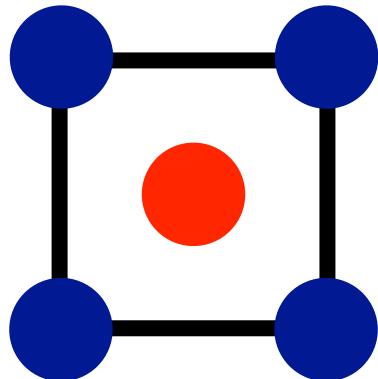
2. **Size:**

descriptors are large ($\sim 10^4$ features)

3. **Lack of interpretability:**

entries have no direct physical or chemical meaning

Descriptor based neural network potentials

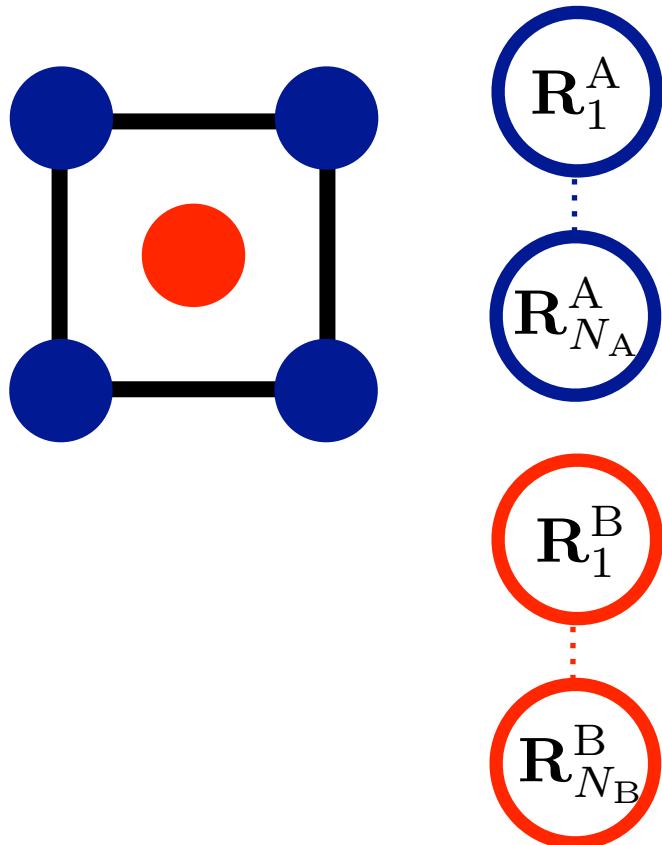


Locality as basis for inter-atomic potentials:

$$E = \sum_i^{N_{atoms}} E_i$$

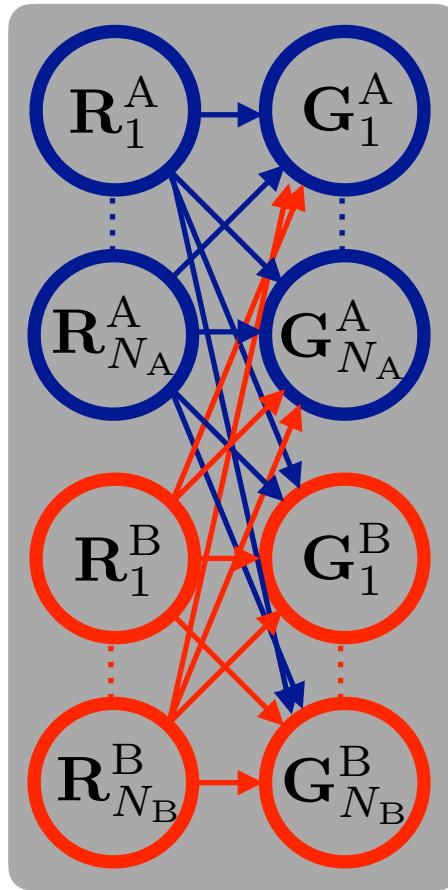
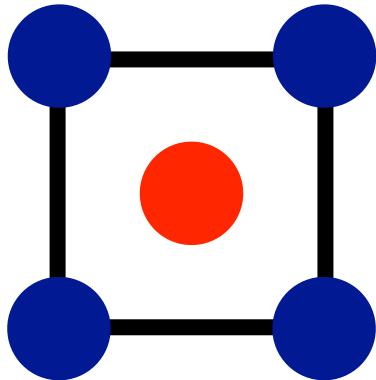
$$E = \sum_j^{N_{elements}} \sum_i^{N_{atoms,j}} E_i^j$$

Behler-Parinello neural network potential



xyz coordinates of all atoms

Behler-Parinello neural network potential

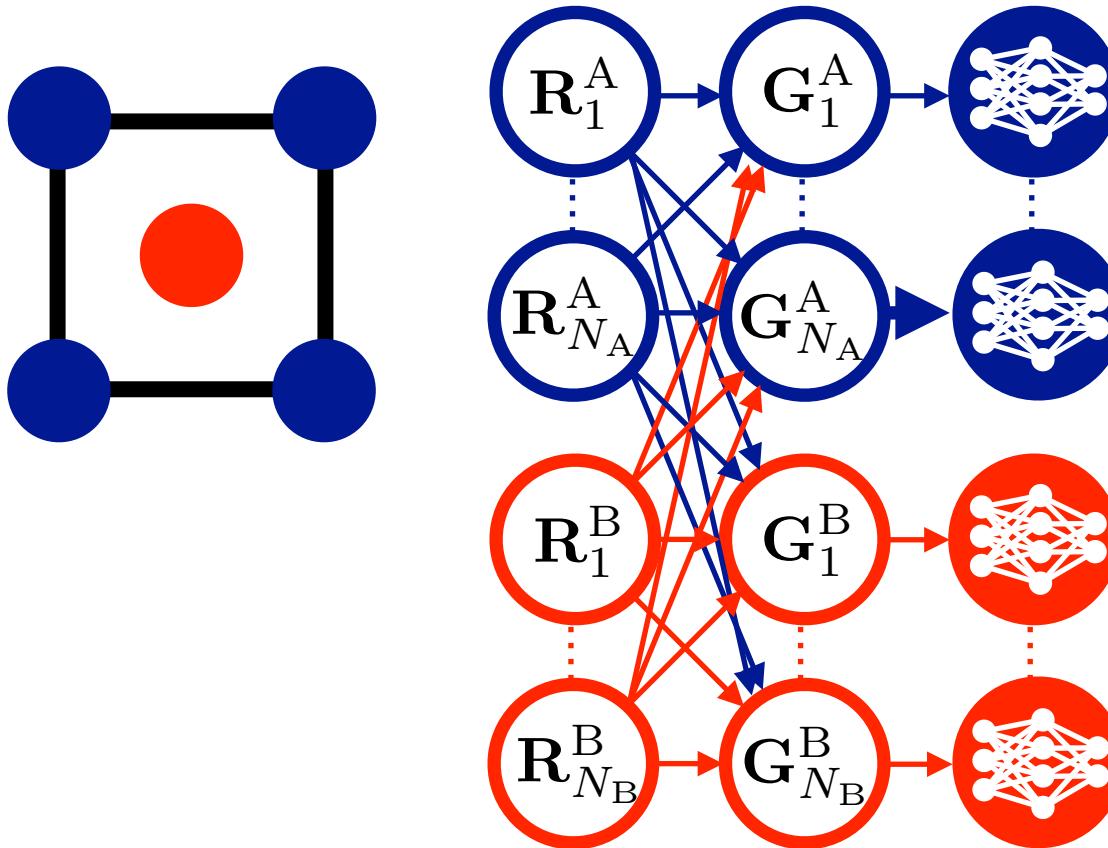


Mapping over via
descriptors functions

Descriptor:

- enforces invariances
- encodes local environments

Behler-Parinello neural network potential

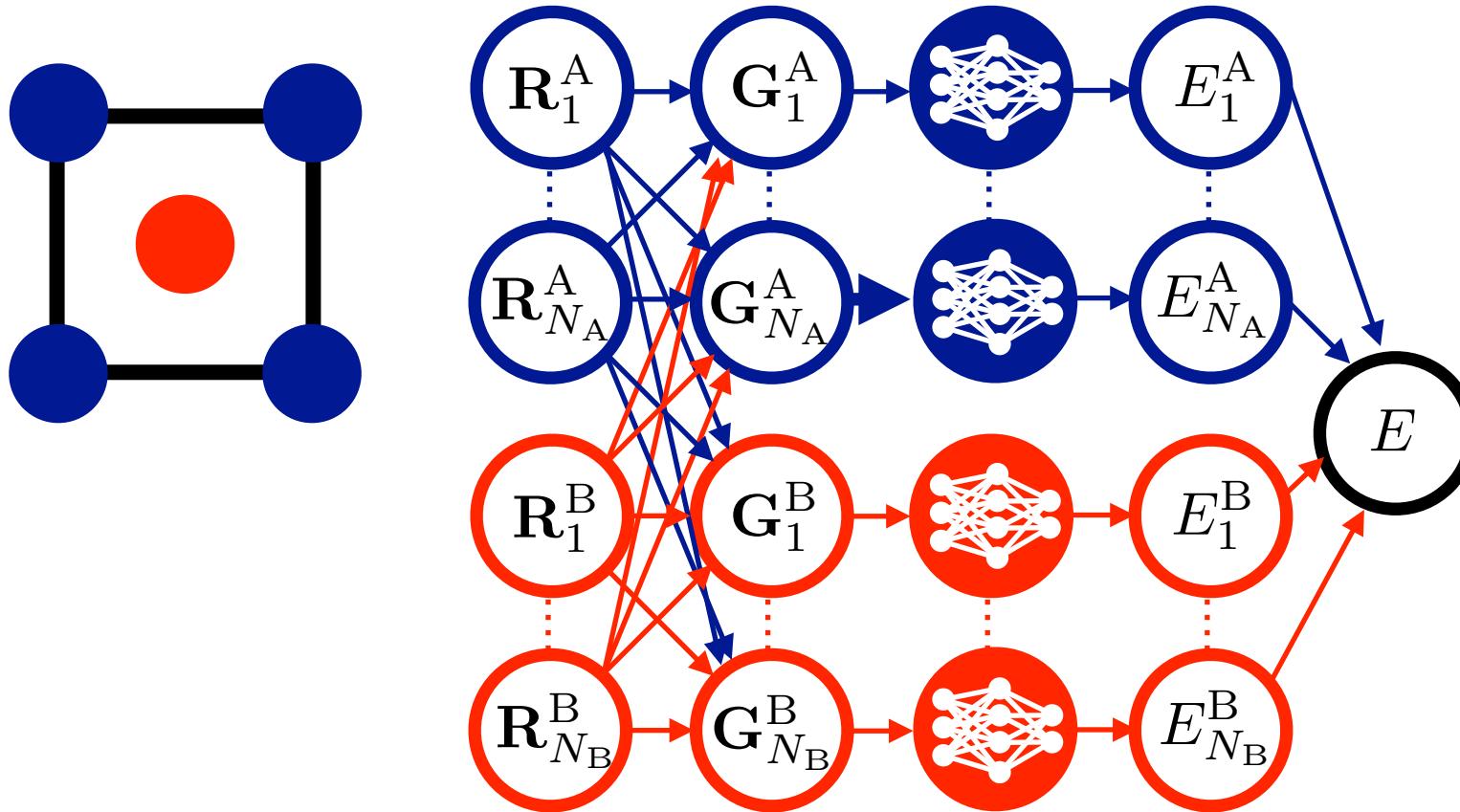


Add 1 neural network (NN) per atomic environment

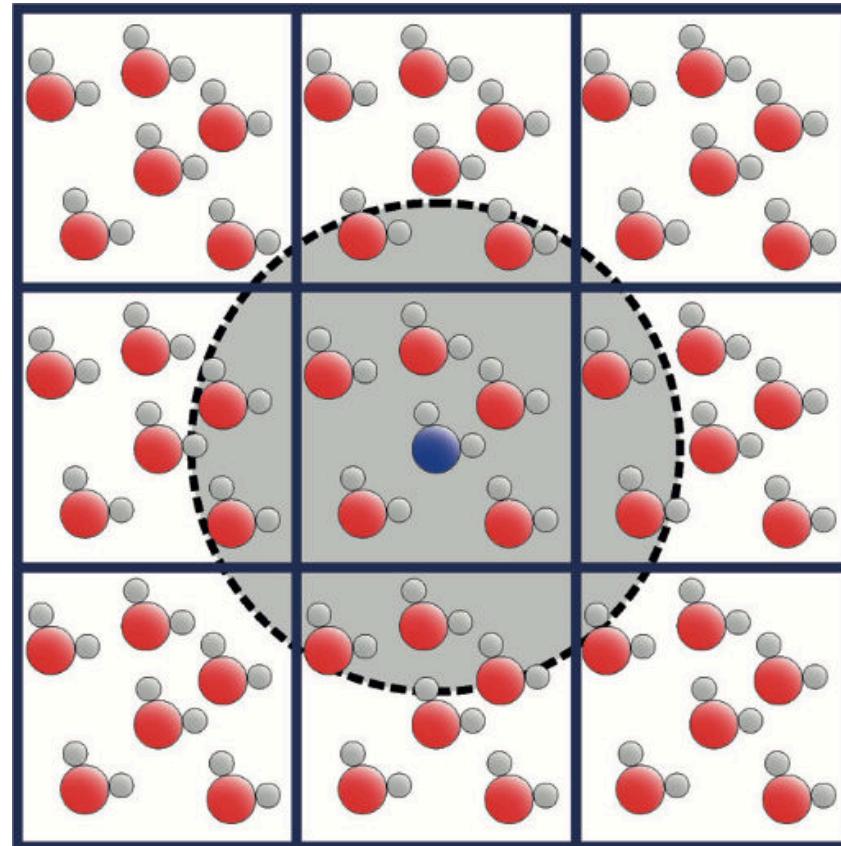
Only 1 NN per species

Difference comes in through local environments

Behler-Parinello neural network potential



Environment dependence



The **cut-off function** controls the number of atoms in the local environment.

Machine learning interatomic potentials (MLIP) can be trained for smaller structures and then applied to larger.

MLIPs scale usually linear with the number of particles

EPFL How to get the forces and the stress

59

Force on an atom: $F_i = -\frac{\partial E}{\partial \mathbf{R}_i}$ ↪ without efforts by auto grad ✨

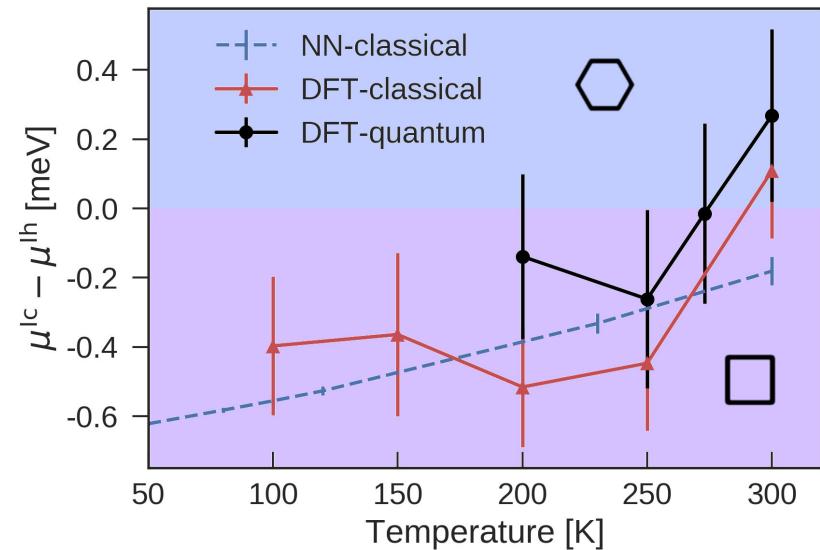
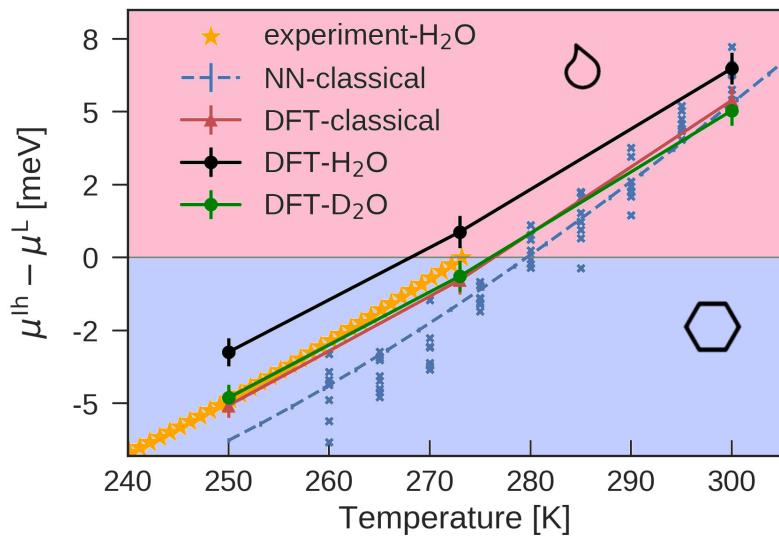
Stress: $V\sigma = \frac{\partial E}{\partial \epsilon} \Big|_{\epsilon=0} = \sum_i^{N_{\text{atoms}}} \mathbf{R}_i \otimes \frac{\partial E}{\partial \mathbf{R}_i}$
 ϵ is the strain

Forces and stress are usually part of the loss function in training.

EPFL Ab initio (thermo)dynamics made easy

60

Simulating matter at finite temperature, including quantum nuclear effects and dynamics is now much more affordable!

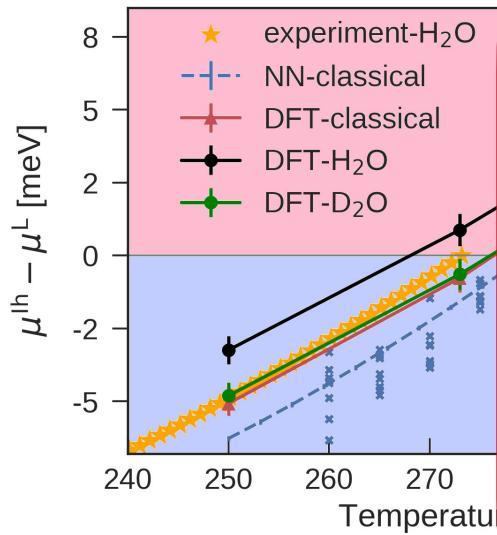


EPFL

Ab initio (thermo)dynamics made easy

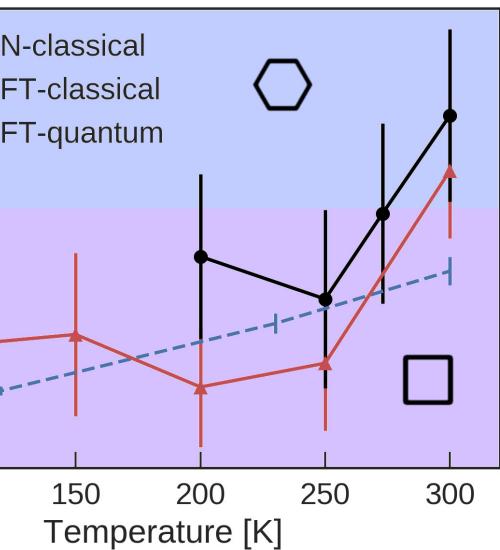
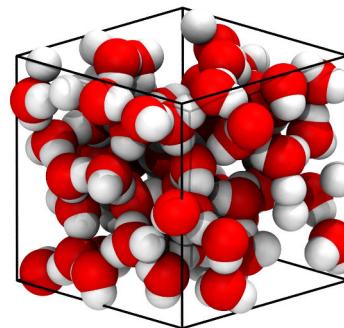
60

Simulating matter at finite temperature, including quantum nuclear effects and dynamics is now much more affordable!



Dataset

1,593 diverse reference structures of
64 molecules of liquid water at revPBE0 level.



1. **invariant:**

invariance of the energy to rotation and translation built in

2. **flexible functional form**

species NNs can adjust to training data; apart from the symmetry functions for the local environments, no functional form for the potential is constraining

3. **locality**

trained on smaller structures, but applied to larger, more complex ones

4. **accuracy**

significantly more accurate than conventional force fields

5. **versatility**

works for molecules and periodic structures

1. **training data:**

require significantly more training data than conventional force fields

2. **transferability:**

often generalise poorly to new or different environments

3. **scaling with species:**

scale exponentially with the number of chemical species.

End-to-end neural network potentials

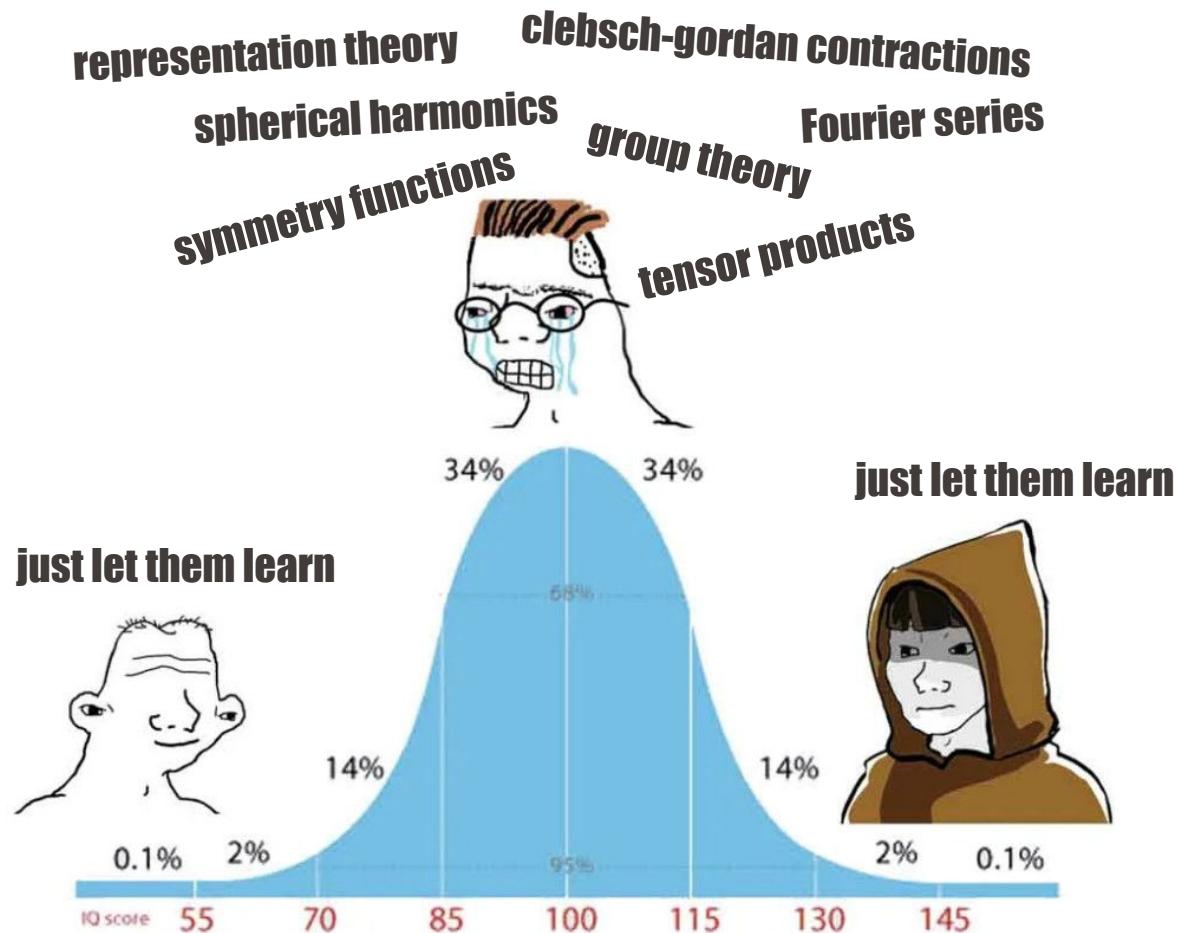
Can we do better?

1 H Hydrogen 1.008	3 Li Lithium 6.941	4 Be Beryllium 9.012	11 Na Sodium 22.990	12 Mg Magnesium 24.305	19 K Potassium 39.098	20 Ca Calcium 40.078	21 Sc Scandium 44.956	22 Ti Titanium 47.867	23 V Vanadium 50.942	24 Cr Chromium 51.996	25 Mn Manganese 54.938	26 Fe Iron 55.845	27 Co Cobalt 58.933	28 Ni Nickel 58.693	29 Cu Copper 63.546	30 Zn Zinc 65.38	5 B Boron 10.811	6 C Carbon 12.011	7 N Nitrogen 14.007	8 O Oxygen 15.999	9 F Fluorine 18.998	10 Ne Neon 20.180
37 Rb Rubidium 85.468	38 Sr Strontium 87.62	39 Y Yttrium 88.906	40 Zr Zirconium 91.224	41 Nb Niobium 92.906	42 Mo Molybdenum 95.95	43 Tc Technetium 98.907	44 Ru Ruthenium 101.07	45 Rh Rhodium 102.906	46 Pd Palladium 106.42	47 Ag Silver 107.868	48 Cd Cadmium 112.411	49 In Indium 114.818	50 Sn Tin 118.711	51 Sb Antimony 121.760	52 Te Tellurium 127.6	53 I Iodine 126.904	54 Xe Xenon 131.294					
55 Cs Cesium 132.905	56 Ba Barium 137.328	57-71 [RE]	72 Hf Hafnium 178.49	73 Ta Tantalum 180.948	74 W Tungsten 183.84	75 Re Rhenium 186.207	76 Os Osmium 190.23	77 Ir Iridium 192.217	78 Pt Platinum 195.085	79 Au Gold 196.967	80 Hg Mercury 200.592	81 Tl Thallium 204.383	82 Pb Lead 207.2	83 Bi Bismuth 208.980	84 Po Polonium [208.982]	85 At Astatine 209.987	86 Rn Radon 222.018					
87 Fr Francium 223.020	88 Ra Radium 226.025	89-103 [RE]	104 Rf Rutherfordium [261]	105 Db Dubnium [262]	106 Sg Seaborgium [266]	107 Bh Bohrium [264]	108 Hs Hassium [269]	109 Mt Meitnerium [278]	110 Ds Darmstadtium [281]	111 Rg Roentgenium [280]	112 Cn Copernicium [285]	113 Nh Nihonium [286]	114 Fl Flerovium [289]	115 Mc Moscovium [289]	116 Lv Livermorium [293]	117 Ts Tennessee [294]	118 Og Oganesson [294]					

57 La Lanthanum 138.905	58 Ce Cerium 140.116	59 Pr Praseodymium 140.908	60 Nd Neodymium 144.242	61 Pm Promethium 144.913	62 Sm Samarium 150.36	63 Eu Europium 151.964	64 Gd Gadolinium 157.25	65 Tb Terbium 158.925	66 Dy Dysprosium 162.500	67 Ho Holmium 164.930	68 Er Erbium 167.259	69 Tm Thulium 168.934	70 Yb Ytterbium 173.055	71 Lu Lutetium 174.967
89 Ac Actinium 227.028	90 Th Thorium 232.038	91 Pa Protactinium 231.036	92 U Uranium 238.029	93 Np Neptunium 237.048	94 Pu Plutonium 244.064	95 Am Americium 243.061	96 Cm Curium 247.070	97 Bk Berkelium 247.070	98 Cf Californium 251.080	99 Es Einsteinium [254]	100 Fm Fermium 257.095	101 Md Mendelevium 258.1	102 No Nobelium 259.101	103 Lr Lawrencium [262]



Can we do better?

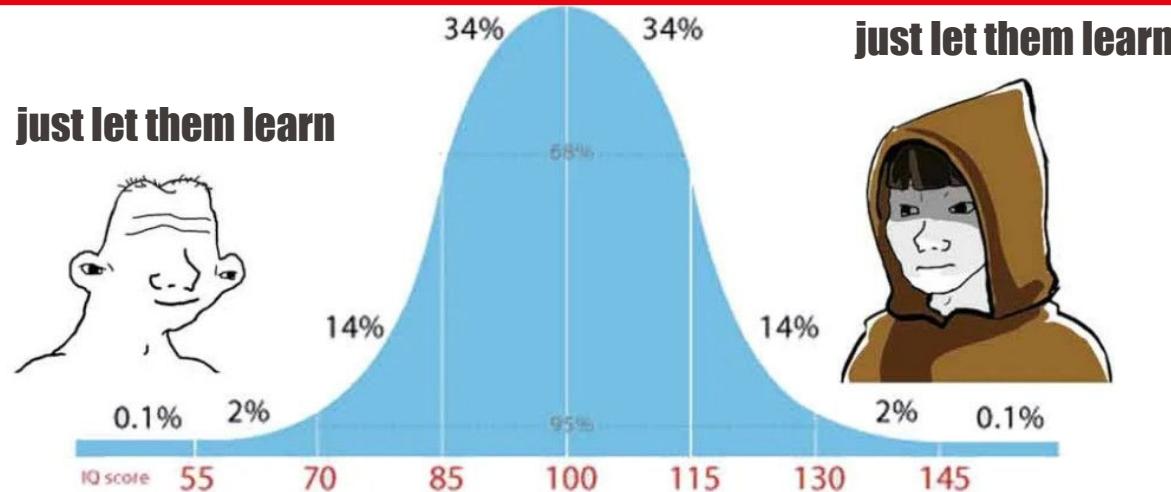


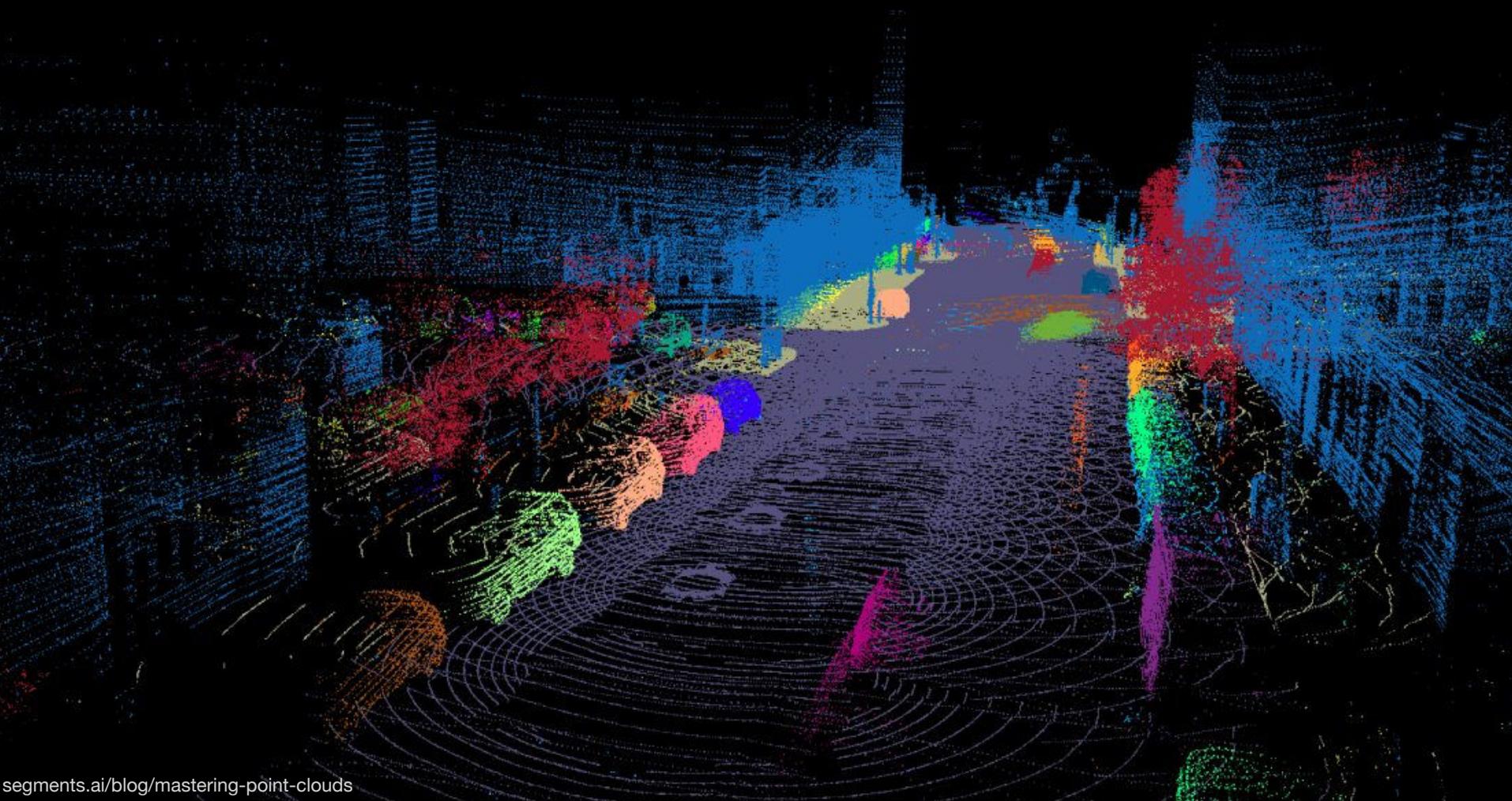
Can we do better?

representation theory clebsch-gordan contractions
spherical harmonics group theory
Fourier series

Current state-of-the-art MLIP:

Equivariant message passing graph neural networks





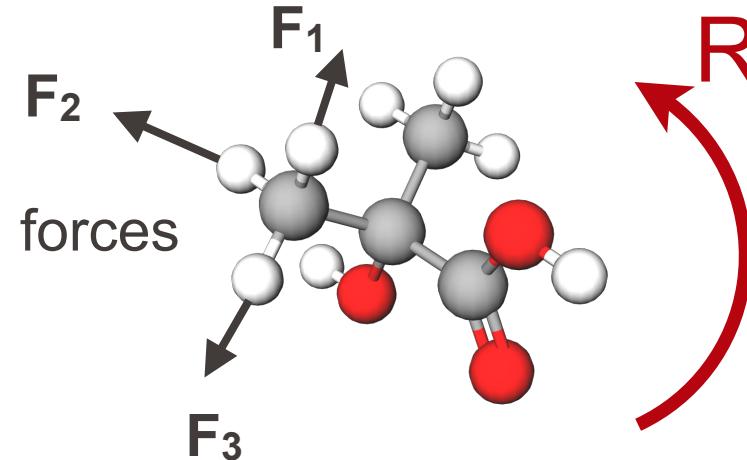
Invariance and equivariance

Equivariant message passing graph neural networks

Invariance with respect to rotation R :

$$E(Rx) = E(x) \text{ or in general } f(Rx) = f(x) \text{ for all } R \in \mathbb{R}^{(3 \times 3)}$$

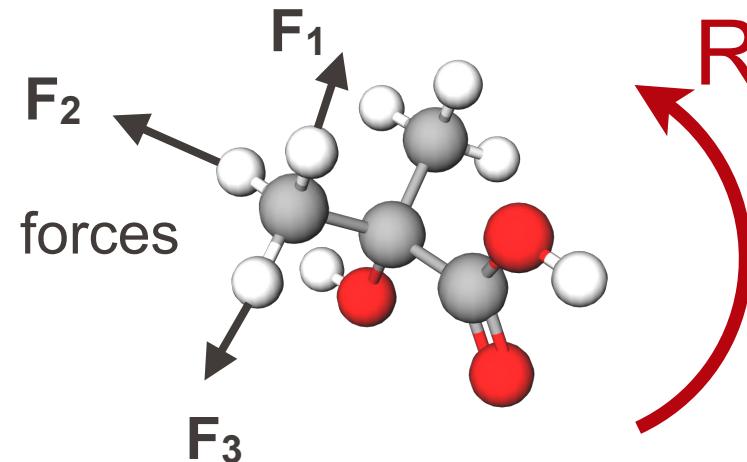
Rotational invariance works for energy, but not forces!



Vectors (and tensors) change under rotations

Equivariance with respect to rotation R :

$$f(Rx) = Rf(x) \text{ for all } R \in \mathbb{R}^{(3 \times 3)}$$



Vectors (and tensors) change under rotations

Equivariance with respect to rotation R :

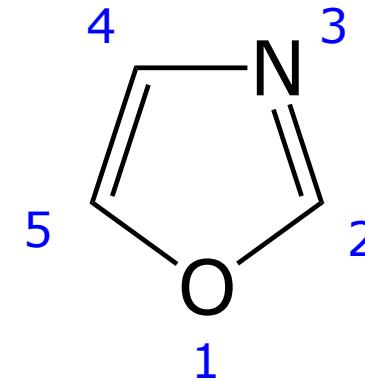
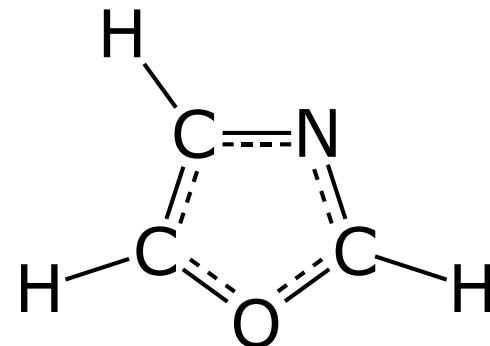
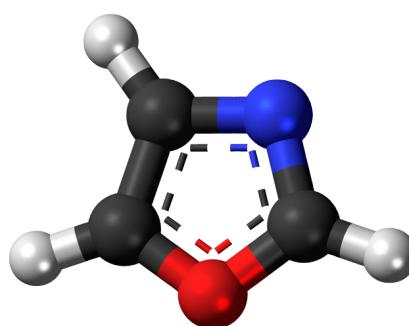
$$f(Rx) = Rf(x) \text{ for all } R \in \mathbb{R}^{(3 \times 3)}$$

If we built equivariance explicitly into the architecture, it would be more data efficient, because it does not have to learn equivariance from data.

F_3

Graph-based representation

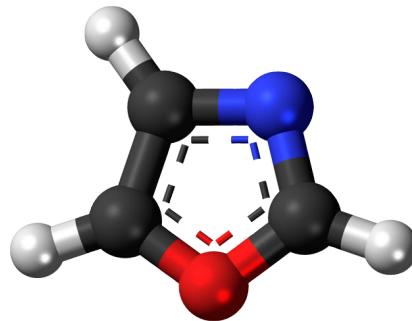
Equivariant message passing **graph neural networks**



Representing molecules as graphs
(similar as molecular representations) - works also for solids

Graph-based representation

Equivariant message passing **graph neural networks**



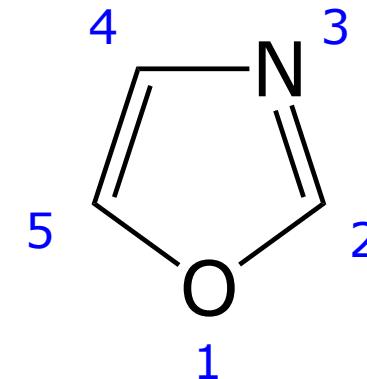
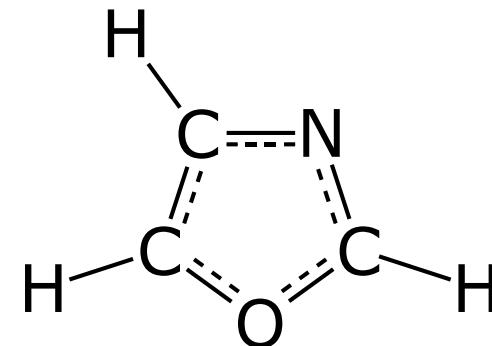
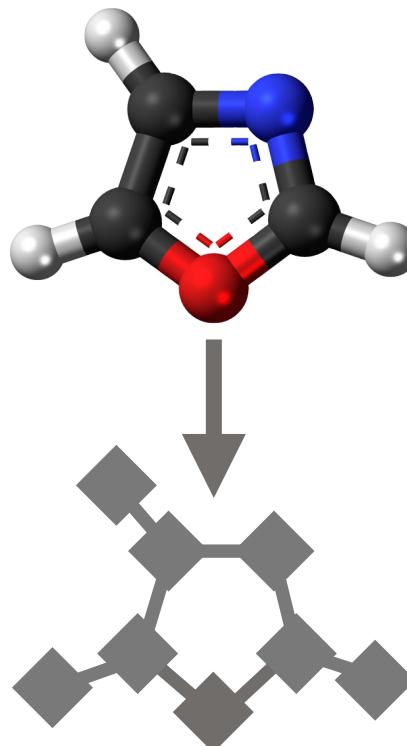
adjacency matrix
encodes graph
structure

Adjacency matrix H :

	1	2	3	4	5	6	7	8
1	0	1	0	0	1	0	0	0
2	1	0	1	0	0	0	1	0
3	0	1	0	1	0	0	0	0
4	0	0	1	0	1	0	0	1
5	1	0	0	1	0	1	0	0
6	0	0	0	0	1	0	0	0
7	0	1	0	0	0	0	0	0
8	0	0	0	1	0	0	0	0

Graph-based representation

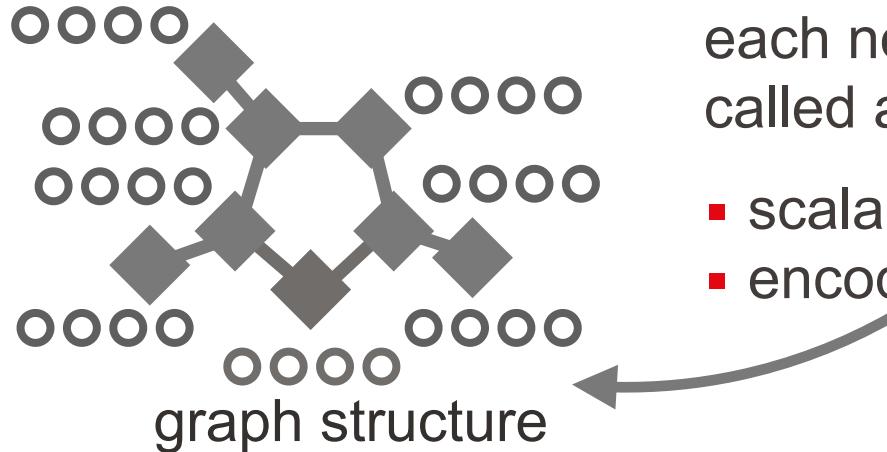
Equivariant message passing **graph neural networks**



Each atom becomes a **node** in the graph

Each “bond” becomes an **edge** in the graph

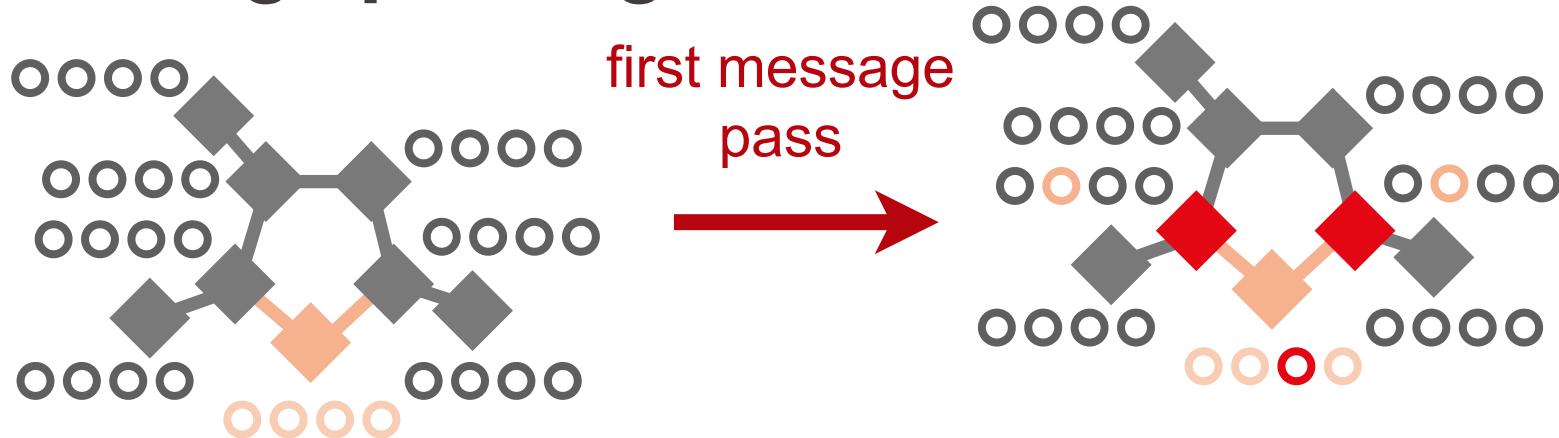
Graph-based representation



each node carries information,
called attributes:

- scalars, vectors, tensors,...
- encoded as neurons

Message passing



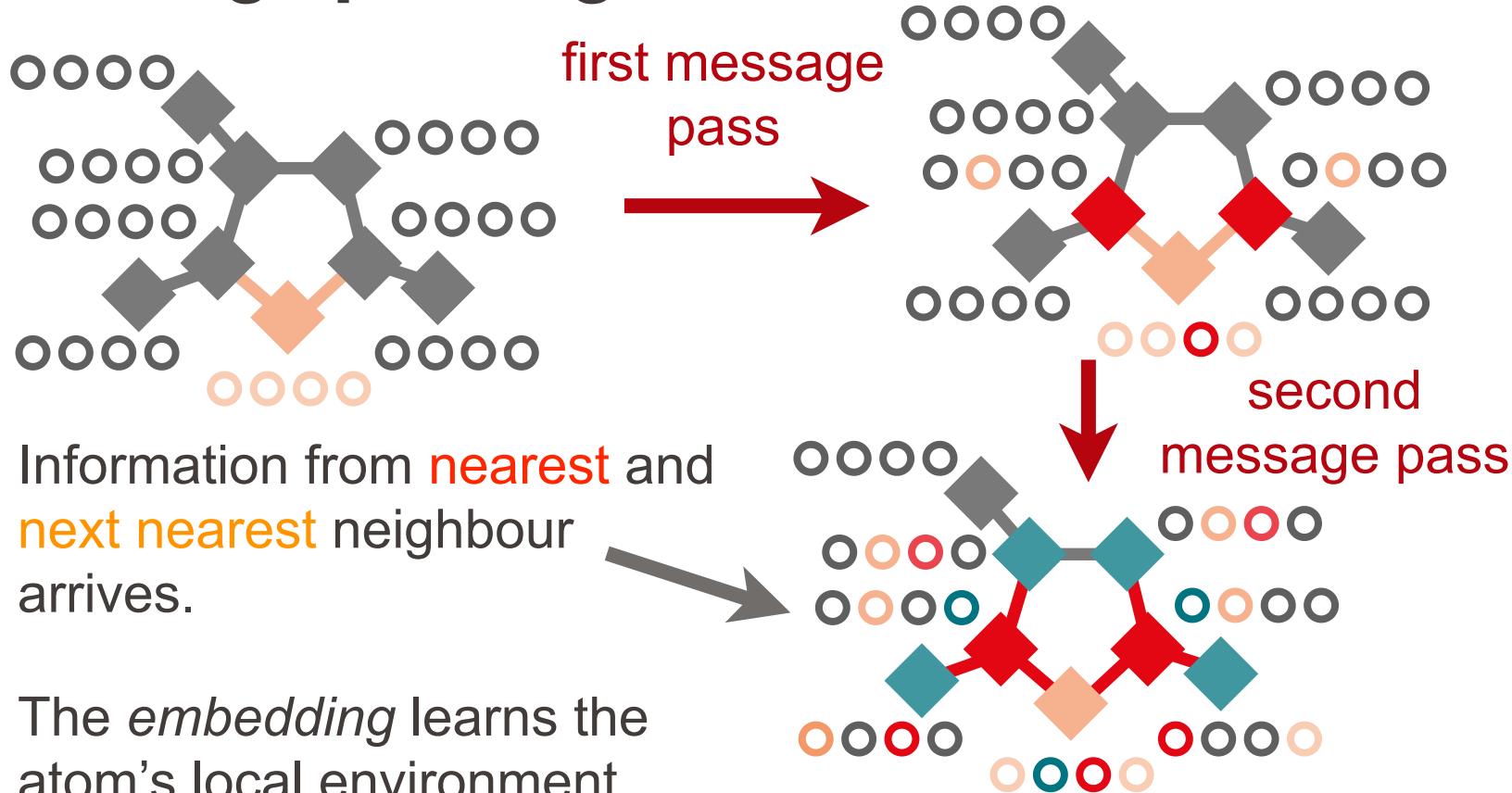
first message
pass

Each node passes its attributes to its neighbours
(according to adjacency matrix) and

learnable update and **interaction function** over the edges

$$\mathbf{h}_i^{(l+1)} = f_{\text{update}} \left(\mathbf{h}_i^{(l)}, \sum_{j \in \mathcal{N}_i} f_{\text{int}} \left[\mathbf{h}_j^{(l)}, \mathbf{e}_{(ij)}^{(l)} \right] \right)$$

Message passing



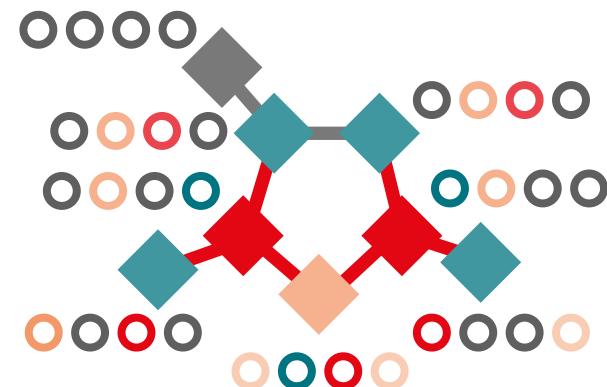
Message passing

Message passing builds up the embeddings,
i.e. local environments, of each atom.

These embeddings are our atomic descriptors.

In descriptor based NN, embeddings were mostly input.
Now they are mostly learned.

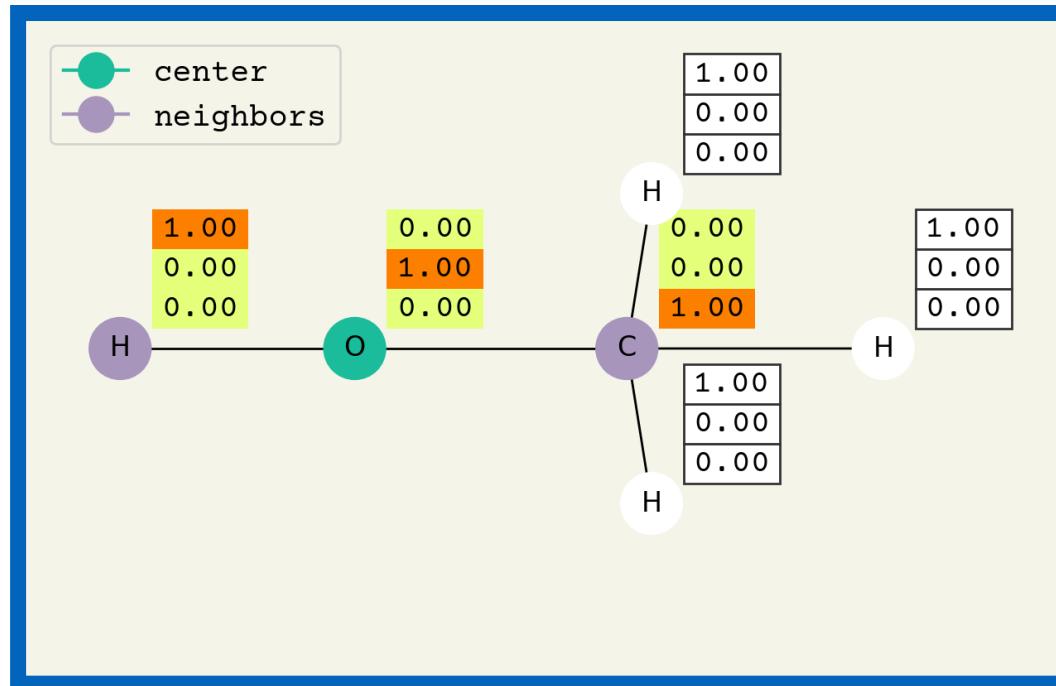
Same approach as before
to get the energy/forces etc.



EPFL A very simple graph neural network

76

Node features: one-hot based on the element

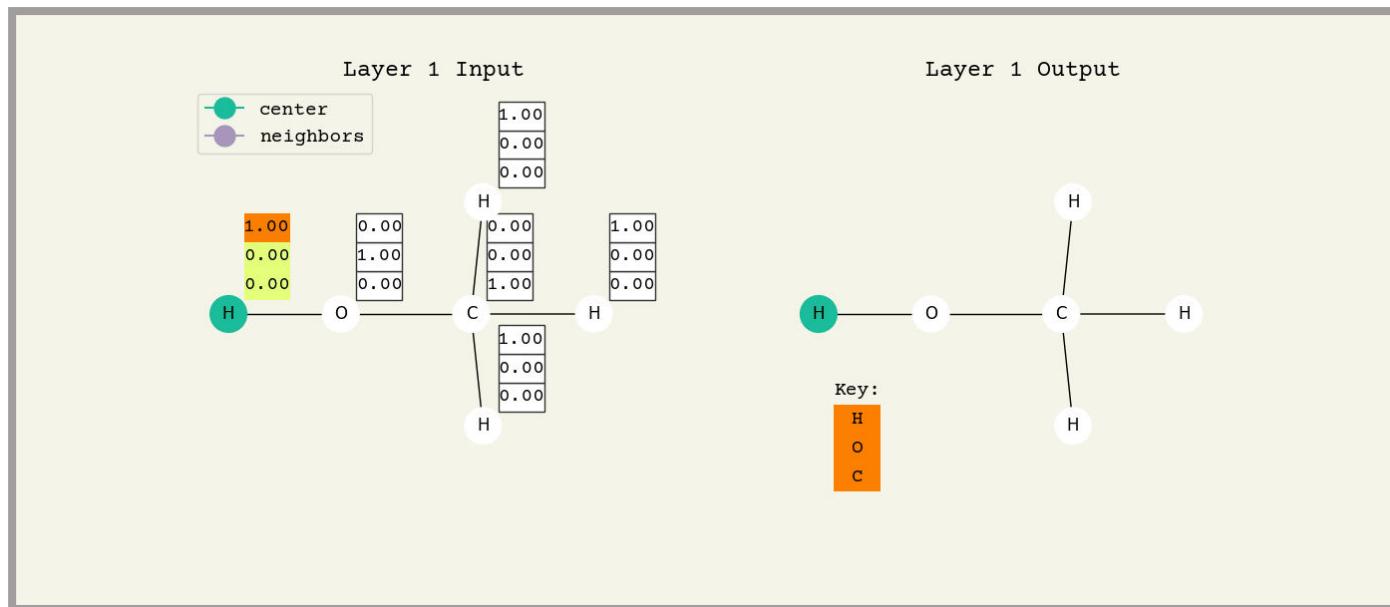


EPFL A very simple graph neural network

77

Node features: one-hot based on the element

Inference: Update via the averaging over neighbors



How successful are GNNs models



Matbench Discovery

Model	CPS ↑	Acc ↑	F1 ↑	DAF ↑	Prec ↑	MAE ↓	R ² ↑	κ_{SRME} ↓	RMSD ↓	Training Set	Params	Targets	Date Added	r _{cut}	Org
eSEN-30M-OAM	0.888	0.977	0.925	6.069	0.928	0.018	0.866	0.170	0.061	6.6M (113M) OMat24+MPtrj+sAlex	30.2M	EFS _G	2025-03-17	6 Å	
Nequip-OAM-L	0.870	0.967	0.893	5.823	0.890	0.022	0.865	0.166	0.065	6.6M (113M) OMat24+sAlex+MPtrj	9.6M	EFS _G	2025-09-08	6 Å	
ORB v3	0.861	0.971	0.905	5.912	0.904	0.024	0.821	0.210	0.075	6.47M (133M) MPtrj+Alex+OMat24	25.5M	EFS _G	2025-04-05	6 Å	
SevenNet-MF-ompa	0.845	0.969	0.901	5.825	0.890	0.021	0.867	0.317	0.064	6.6M (113M) OMit24+sAlex+MPtrj	25.7M	EFS _G	2025-03-13	6 Å	
Allegro-OAM-L	0.840	0.966	0.895	5.674	0.867	0.022	0.868	0.319	0.065	6.6M (113M) OMit24+sAlex+MPtrj	9.7M	EFS _G	2025-09-08	7 Å	
GRACE-2L-OAM	0.837	0.963	0.880	5.774	0.883	0.023	0.862	0.294	0.067	6.6M (113M) OMit24+sAlex+MPtrj	12.6M	EFS _G	2025-02-06	6 Å	ICAMS
DPA-3.1-3M-FT	0.802	0.963	0.884	5.667	0.866	0.023	0.869	0.469	0.069	163M OpenLAMMPS	3.27M	EFS _G	2025-06-05	6 Å	
eSEN-30M-MP	0.797	0.946	0.831	5.260	0.804	0.033	0.822	0.340	0.075	146k (1.58M) MPtrj	30.1M	EFS _G	2025-03-17	6 Å	
MACE-MPA-0	0.795	0.954	0.852	5.582	0.853	0.028	0.842	0.412	0.073	3.37M (12M) MPtrj+sAlex	9.06M	EFS _G	2024-12-09	6 Å	

Usually models trained on very large datasets

The field moves very fast

The Open Catalyst 2025 (OC25) Dataset and Models for Solid-Liquid Interfaces



Sushree Jagriti Sahoo¹, Mikael Maraschin², Daniel S. Levine¹, Zachary Ulissi¹, C. Lawrence Zitnick¹, Joel B Varley⁴, Joseph A. Gauthier^{2,†}, Nitish Govindarajan^{3,4†}, Muhammed Shuaibi^{1,†}

¹FAIR at Meta, ²Department of Chemical Engineering, Texas Tech University, Lubbock, TX 79409, USA, ³School of Chemistry, Chemical Engineering and Biotechnology, Nanyang Technological University, 21 Nanyang Link, Singapore 637371, Singapore, ⁴Materials Science Division, Lawrence Livermore National Laboratory, Livermore, CA 94550, USA

[†]Co-corresponding Author



Designing quantum chemistry algorithms with just-in-time compilation

Xiaojie Wu^{1,†}, Qiming Sun¹, Yuanheng Wang¹

¹ByteDance Seed

[†]Corresponding authors



SimpleFold: Folding Proteins is Simpler than You Think

Yuyang Wang, Jiarui Lu*, Navdeep Jaitly, Josh Susskind, Miguel Angel Bautista

Apple

Protein folding models have achieved groundbreaking results typically via a combination of integrating domain knowledge into the architectural blocks and training pipelines. Nonetheless, given the success of generative models across different but related problems, it is natural to question whether these architectural designs are a necessary condition to build performant models. In this paper, we introduce *SimpleFold*, the first flow-matching based protein folding model that solely uses general purpose transformer blocks. Protein folding models typically employ computationally expensive modules involving triangular updates, explicit pair representations or multiple training objectives curated for this specific domain.

MatterSim: A Deep Learning Atomistic Model Across

Elements, Temperatures and Pressures

Han Yang^{①*†}, Chenxi Hu^{②†}, Yichi Zhou^{1†}, Xixian Liu^{③†}, Yu Shi^{④†},
 Jielan Li^{⑤*†}, Guanzhi Li^{⑥†}, Zekun Chen^{⑦†}, Shuizhou Chen^{⑧†},
 Claudio Zeni^⑨, Matthew Horton^⑩, Robert Pinsler^⑪, Andrew Fowler^⑫,
 Daniel Zügner^⑬, Tian Xie^⑭, Jake Smith^⑮, Lixin Sun^⑯, Qian Wang^⑰,
 Lingyu Kong^⑱, Chang Liu^⑲, Hongxia Hao^{⑳*}, Ziheng Lu^{⑳*}



*Corresponding author(s). E-mail(s): hanyang@microsoft.com; jielanli@microsoft.com; hongxiahao@microsoft.com; zihenglu@microsoft.com;

[†]These authors contributed equally to this work.

1. **invariant & equivariant:**
energy invariance and force equivariance built in
2. **flexible functional form**
representations are learned internally
3. **locality**
allows to be trained on smaller structures, but applied to larger, more complex ones
4. **accuracy**
significantly more accurate than conventional force fields
5. **versatility**
works for molecules and periodic structures
6. **efficiency**
quite data efficient due to complex architectures

1. **training data:**

require significantly more training data than conventional force fields, but less than descriptor based NNPs

2. **run time:**

execution time is higher than descriptor based NNPs

3. **long range:**

long range effects (like van der Waals or electrostatics) were not included in early generations, but are being added now.

4. **charge:**

explicit charging was not included in early generations, but is being added now.

1. Descriptors (attribute or structure based) are *currently* a requirement to perform atomistic machine learning
2. Descriptors need certain properties to be expressive: atom-centered, invariant, equivariant, ...
3. In combination with neural networks, descriptors allow the exact modeling of the PES
4. Modern state-of-the-art descriptors are learnable