

Local features

2. Smooth Overlap of Atomic Orbitals (SOAP)

- Atomic density field $\rho^Z(\vec{r}) = \sum_i w_{Z_i} \delta(\vec{r} - \vec{R}_i) \delta_{Z_i, Z}$, or

- Gaussian smeared atomic density field $\rho^Z(\vec{r}) = \sum_i e^{-\frac{1}{2\sigma^2} |\vec{r} - \vec{R}_i|^2} \delta_{Z_i, Z}$.

- $\rho^Z(\vec{r})$ in a local region surrounding the point of interest as origin expanded as

$$\rho^Z(\vec{r}) = \sum_{n=1}^{n_{\max}} \sum_{l=1}^{l_{\max}} \sum_{m=-l}^l c_{nlm}^Z g_{nl}(r) Y_{lm}(\theta, \phi). \quad C_{nlm}^Z \text{ s are coefficients.}$$

- $Y_{lm}^l(\theta, \phi) = (-1)^l \sqrt{\frac{(2l+1)(l-m)!}{4\pi(l+m)!}} P_{lm}(\cos \theta) e^{im\phi}$, angular basis functions.

- If g_{nl} 's are ortho-normal functions

- $c_{nlm} = \int_V g_n(r) Y_{lm}^*(\theta, \phi) \rho(\vec{r}) dv.$

Two types of radial basis functions:

- Spherical primitive Gaussian orbitals as

$$g_{nl}(r) = \sum_{n'=1}^{n_{\max}} \beta_{nn'l} \Phi_{n'l}(r) \quad \text{with} \quad \Phi_{nl}(r) = r^l e^{-\alpha_n r^2}$$

- α_n is the decay parameter, chosen such that each $\Phi_{nl}(r)$ decays to 10^{-3} at a cutoff radius r_{cut} .

- $g_{nl}(r)$'s are not orthonormal

- They can be made orthonormal by Löwdin orthonormalization procedure:

$$\bullet \text{ Overlap integral } S_{nn'} \equiv \int_0^{\infty} r^2 (\Phi_{nl}(r)) (\Phi_{n'l}(r)) dr = \int_0^{\infty} r^2 (r^l e^{-\alpha_n r^2}) (r^l e^{-\alpha_{n'} r^2}) dr$$

- Choose $\beta = \mathbf{S}^{-1/2}$

$$\bullet \int r^2 g_{ml}(r) g_{nl}(r) dr = \int r^2 \left(\sum_{m'} \beta_{mm'l} \Phi_{m'l}(r) \right) \left(\sum_{n'} \beta_{nn'l} \Phi_{n'l}(r) \right) dr$$

$$\bullet \beta_{mm'l} \beta_{nn'l} \int r^2 \Phi_{m'l}(r) \Phi_{n'l}(r) dr = \beta_{mm'l} \beta_{nn'l} S_{m'n'} = \beta_{mm'l} S_{m'n'} \beta_{n'n'l} \text{ matrices are real symmetric}$$

$$\bullet = S_{mm'}^{-1/2} S_{m'n'} S_{n'n}^{-1/2} = \delta_{mn}.$$

- Dscribe (we will work with later) includes up to $l \leq 9$.

- Polynomial functions

$$\bullet g_n(r) = \sum_{n'=1}^{n_{\max}} \beta_{nn'} \Phi_{n'}(r) \quad \text{with} \quad \Phi_{nl}(r) = (r - r_{\text{cut}})^{n+2}$$

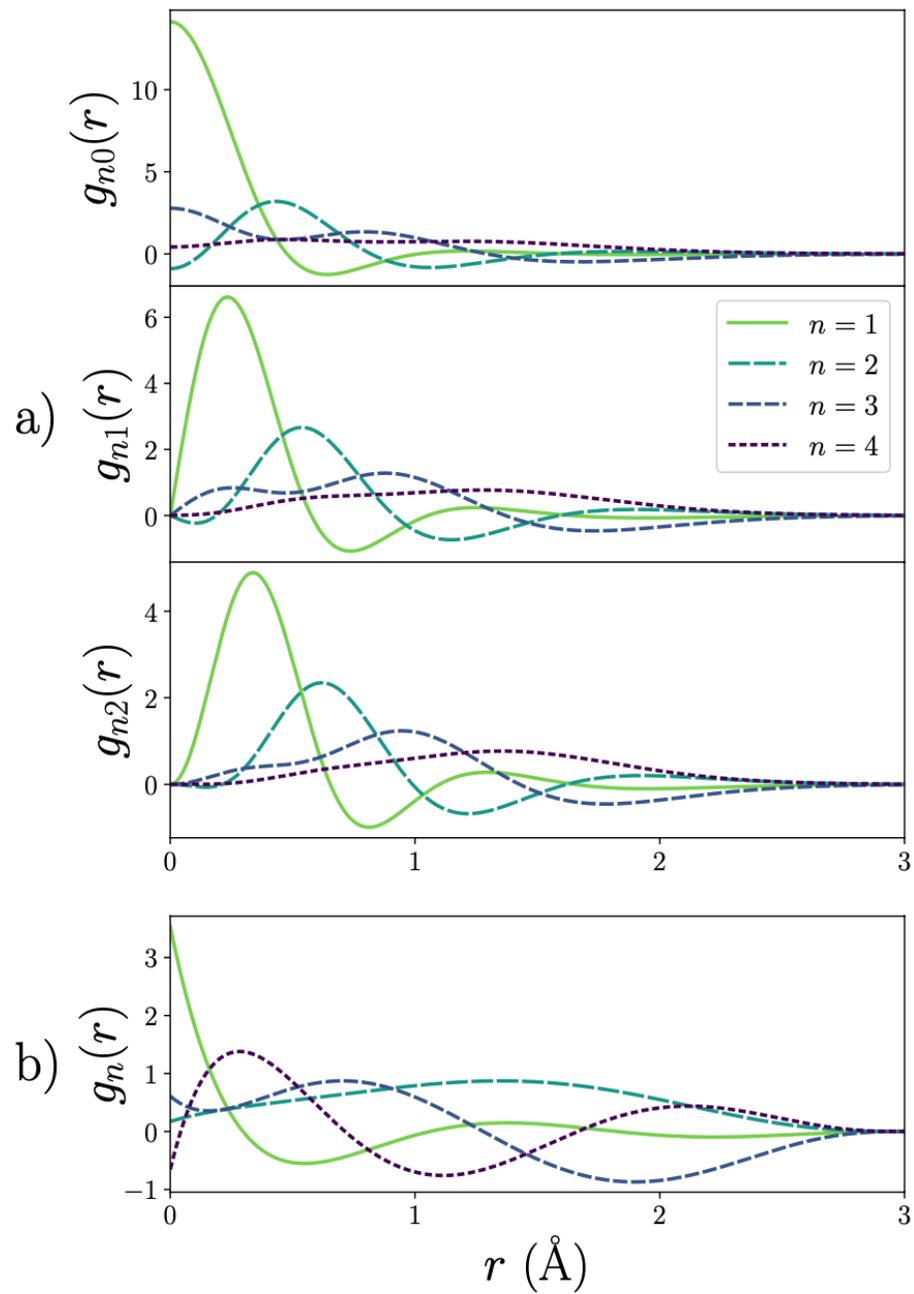
- α_n is the decay parameter, chosen such that each $\Phi_{nl}(r)$ decays to 10^{-3} at a cutoff radius r_{cut} .

- These can be made orthonormal by Lowdin orthonormalization procedure: choose

$$\bullet \mathbf{S}_{nn'} = \int_0^{\infty} r^2 (r - r_{\text{cut}})^{n+2} (r - r_{\text{cut}})^{n'+2} dr \quad \text{for each } l$$

$$\bullet \beta = \mathbf{S}^{-1/2}.$$

- Dscribe, cubic or higher order polynomials, $l \leq 20$.



Radial basis functions for SOAP: (a) spherical gaussian, (b) polynomial

- Effect of rotation, rotation operator \hat{R} acting on $\rho(\vec{r})$:

$$\hat{R}\rho(\vec{r}) = \hat{R} \sum_{nlm} c_{nlm} g_{nl}(r) Y_{lm}(\theta, \phi) = \sum_{nlm} c_{nlm} g_{nl}(r) (\hat{R} Y_{lm}(\theta, \phi))$$

- Spherical harmonics under rotation $\hat{R} Y_{lm}(\theta, \phi) = D_{mm'}^l(\hat{R}) Y_{lm'}(\theta, \phi)$, $D^l(\hat{R})$'s are $(2l+1)(2l+1)$ matrices, called Wigner matrix.

$$\hat{R}\rho(\vec{r}) = \sum_{n=1}^{n_{\max}} \sum_{l=0}^{l_{\max}} \sum_{m=-l}^{+l} \sum_{m'=-l}^{+l} c_{nlm} g_{nl}(r) (D_{mm'}^l(\hat{R}) Y_{lm'}(\theta, \phi))$$

$$\hat{R}\rho(\vec{r}) = \sum_{n=1}^{n_{\max}} \sum_{l=0}^{l_{\max}} \sum_{m'=-l}^{+l} \left(\sum_{m=-l}^{+l} c_{nlm} D_{mm'}^l(\hat{R}) \right) g_{nl}(r) Y_{lm'}(\theta, \phi)$$

$$= \sum_{n=1}^{n_{\max}} \sum_{l=0}^{l_{\max}} \sum_{m'=-l}^{+l} c_{nlm'} g_{nl}(r) Y_{lm'}(\theta, \phi); \text{ with } c_{nlm'} \equiv \sum_{m=-l}^{+l} c_{nlm} D_{mm'}^l(\hat{R})$$

- The rotation, thus, can be thought of in terms of the transformation of the coefficients c_{nlm} 's
- That is, the function is expressed in terms of an initial set of Y_{lm} 's (c_{nlm}) and a rotated set of Y_{lm} 's ($c_{nlm'}$).

- Wigner matrices are unitary
 - $\mathbf{D}^{l\dagger}\mathbf{D}^l = \mathbf{I}$, unit matrix.

Therefore, considering c_{nlm} for fixed (n, l) as a column vector $\mathbf{c}_{nl} = \begin{pmatrix} c_{n,l,-l} \\ c_{n,l,-l+1} \\ \vdots \\ c_{n,l,l} \end{pmatrix}$,

- $c_{nlm'} \equiv \sum_{m=-l}^{+l} c_{nlm} D_{mm'}^l(\hat{R})$ can be written as $\mathbf{c}'_{nl} = \mathbf{D}^l \mathbf{c}_{nl}$.
- $\mathbf{c}'_{nl\dagger} \mathbf{c}'_{nl} = (\mathbf{c}_{nl\dagger} \mathbf{D}^{l\dagger})(\mathbf{D}^l \mathbf{c}_{nl}) = \mathbf{c}_{nl\dagger} \mathbf{c}_{nl}$.
- Thus $\mathbf{c}'_{nl\dagger} \mathbf{c}'_{nl}$ is a quantity that remains invariant under rotations.
- This property is used to construct a rotation-invariant feature vector

$$p_{nn'l}^{Z_1, Z_2} = \sum_{m=-l}^{+l} (c_{nlm}^{Z_1})^* (c_{n'lm}^{Z_2}).$$

- Concatenating $p_{nn'l}^{Z_1, Z_2}$'s for all possible (Z_1, Z_2) pairs, and all possible (nl) combinations produce the SOAP feature vector.

- SOAP construction

- Atomic density field
- Expand in terms of radial and angular basis functions
- Usually spherical primitive Gaussian and polynomial radial basis
- Make them orthonormal
- Once they are made orthonormal, the expansion coefficients

$$c_{nlm} = \int_V g_n(r) Y^*(\theta, \phi) \rho(\vec{r}) dv$$

- Once c 's are known, construct $p_{nn'l}^{Z_1, Z_2} = \sum_{m=-l}^{+l} (c_{nlm}^{Z_1})^* (c_{n'lm}^{Z_2})$
- Concatenation of p 's is the SOAP

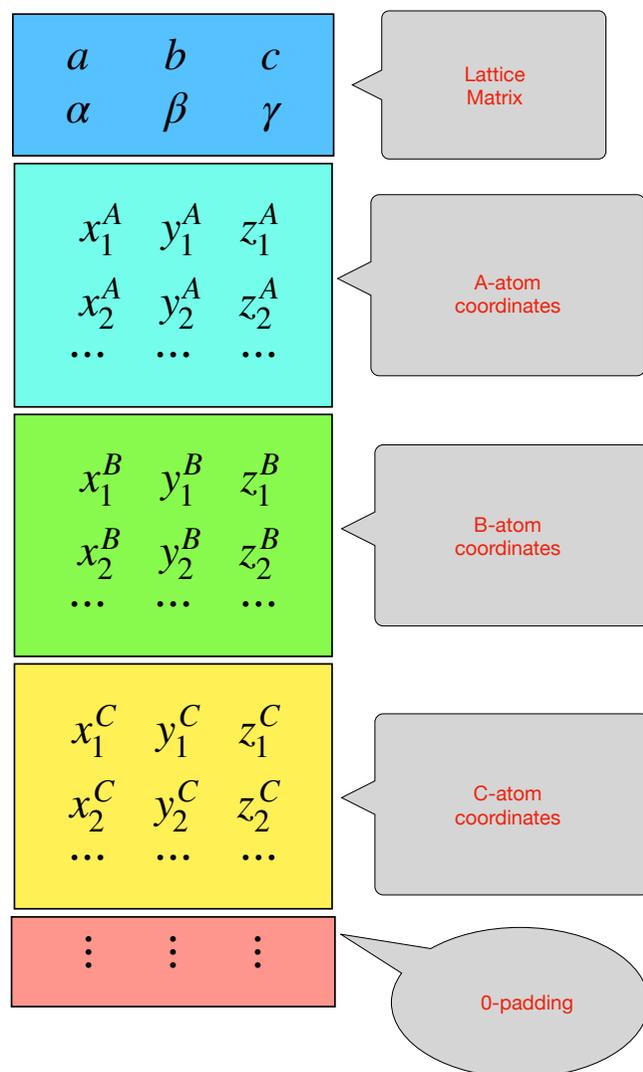
Global features

AnBmCl

- **Point cloud representation:** pictorial representations
 - 3D voxel image, values of some quantity on a 3D grid

computational resources

2-dimensional Point Cloud Representation



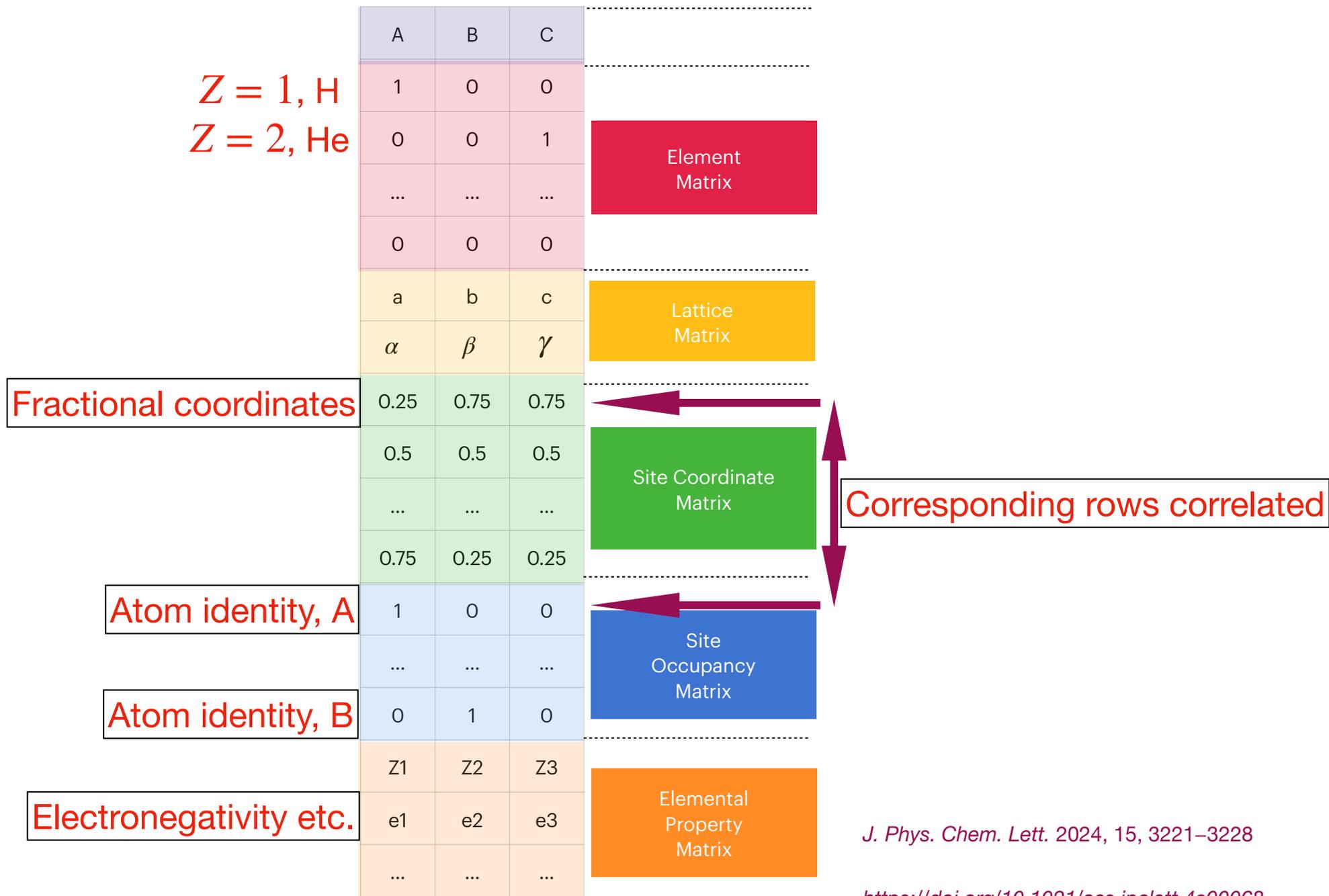
2D point cloud representation
of a ternary material

Accompanied by 'composition condition'

$(N_{\max} + 2) \times 3$ image

N_{\max} is max number of atoms in the unit cell
in the entire training set

IRCR Representation



J. Phys. Chem. Lett. 2024, 15, 3221–3228

<https://doi.org/10.1021/acs.jpcllett.4c00068>

A

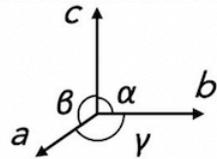
FTCP Representation

real-space features | reciprocal-space features

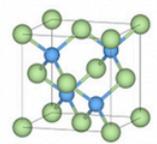
Element Matrix



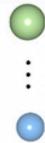
Lattice Matrix



Site Coordinate Matrix



Site Occupancy Matrix



Elemental Property Matrix

electronegativity
atomic radius
⋮

0	0	<u>0</u>	<u>0</u>	...	
⋮	⋮	⋮	⋮		
0	0	<u>0</u>			
a	b	c			
α	β	γ			
0	0	0			
...			
0.25	0.25	0.25			
1	0	<u>0</u>			
...			⋮
0	1	<u>0</u>		...	<u>0</u>
<u>0</u>	<u>0</u>	<u>0</u>	d_{hkl}	...	
0	0	<u>0</u>		...	
⋮	⋮	⋮	⋮	⋮	⋮
0	0	<u>0</u>			

³¹Ga Galium ³³As Arsenic (100) ... (201)

Underlines indicate zero paddings, e.g., 0, ...

Distance of (hkl) from (000)

0	1	<u>0</u>	0	0	<u>0</u>
⋮	⋮	⋮	⋮	⋮	⋮
0	0	<u>0</u>	0.25	0.25	0.25

$$F_{hkl} = \sum_{i=1}^N z_i \cdot \frac{j}{2} \ln(e^{-j2\pi(hx_i + ky_i + lz_i)})$$

FTCP Matrix

Ren et al., Matter 5, 314–335
<https://doi.org/10.1016/j.matt.2021.11.032>