

# Machine learning for drug design, medicine and quantum chemistry : an introduction

Convolutional neural networks and solid harmonic wavelet transform

Robert Benda

CERMICS

November 13, 2019

# Outline

- 1 Machine learning for drug design and medicine : general introduction
  - Examples of applications
  - Molecular descriptors
- 2 Solid harmonic wavelet scattering for predictions of molecule properties (based on Refs. [1, 2])
  - General idea of the method
  - Wavelet scattering coefficients
  - Comparison to usual force fields
  - Why does it work so well : comparison to the multipolar expansion in polarizable force fields methods



Michael Eickenberg, Georgios Exarchakis, Matthew Hirn, Stéphane Mallat, and Louis Thiry.

Solid harmonic wavelet scattering for predictions of molecule properties.

*The Journal of Chemical Physics*, 148(24):241732, jun 2018.



Michael Eickenberg, Georgios Exarchakis, Matthew Hirn, and Stéphane Mallat.

Solid harmonic wavelet scattering: Predicting quantum molecular energy from invariant descriptors of 3d electronic densities.

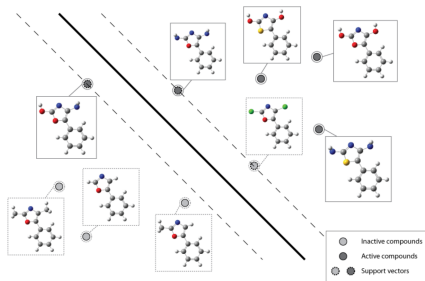
In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 6540–6549. Curran Associates, Inc., 2017.

# Outline

- 1 Machine learning for drug design and medicine : general introduction
  - Examples of applications
  - Molecular descriptors
- 2 Solid harmonic wavelet scattering for predictions of molecule properties (based on Refs. [1, 2])

# Examples of applications (1)

- Compound classification (e.g. with Support Vector Machines) according to their *functions* or *activities*



J.C. Gertrudes, *Machine Learning Techniques and Drug Design*, Current Medicinal Chemistry, **2012**, 19, 4289-4297.

## Examples of applications (2)

- Compound classification (e.g. with SVMs) according to their *functions* or *activities*

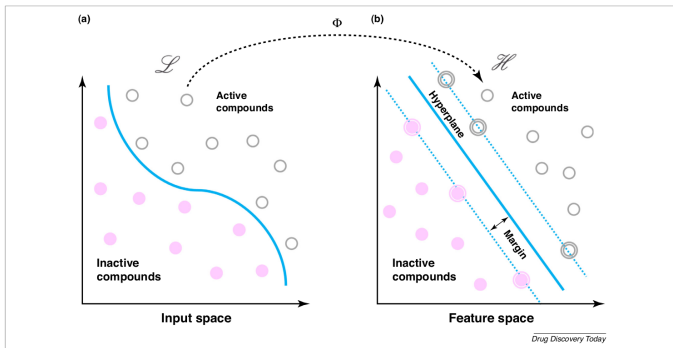


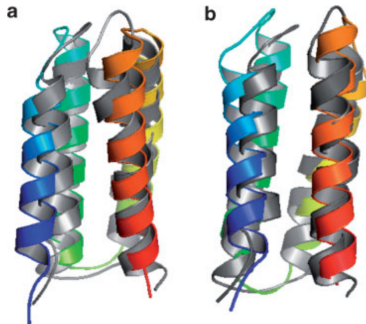
FIGURE 1

Projection into high-dimensional feature space. Using a mapping function  $\Phi$ , active (empty gray points) and inactive (filled pink points) compounds that are not linearly separable in low-dimensional input space  $\mathcal{L}$  (a) are projected into high-dimensional feature space  $\mathcal{H}$  (b) and separated by the maximum-margin hyperplane. Points intercepted by the dotted line are called 'support vectors' (circled points).

A. Lavecchia, *Machine-learning approaches in drug discovery: methods and applications*, Drug Discovery Today, Vol. 20, Nb 3 (2015).

## Examples of applications (3)

- Protein structure prediction  
(e.g. folding from sequence of amino-acids)



C.A. Floudas, *Computational Methods in Protein Structure Prediction*, Biotechnology and Bioengineering, Vol. 97, No. 2, June 1, 2007.

## Examples of applications (4)

- QSAR (Quantitative Structure-Activity Relationship) analysis
- "Virtual screening" (search in the (huge) compound space to meet a given target property : *e.g. drugs*)
- Relationship between genome variation (*e.g. mutation*) and disease risk (genotype-phenotype relation, predictive medicine, etc.)

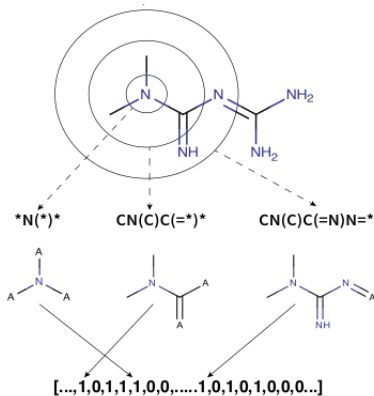


# Molecular descriptors

- Numerical representation necessary (e.g. for input in a neural network) : "common language" to represent molecules.
- Chemical graph theory / chemical fingerprint (e.g. by fragments).
- Example : Coulomb matrix

$$C_{ij} = \begin{cases} 0.5Z_i^{2.4} & \text{if } i=j \\ \frac{Z_i Z_j}{|\vec{r}_i - \vec{r}_j|} & \text{otherwise} \end{cases}$$

Yu-Chen Lo et al., *Machine learning in chemoinformatics and drug discovery*, Drug Discovery Today, Vol. 23, Nb 8 (2018).



# Outline

- 1 Machine learning for drug design and medicine : general introduction
- 2 Solid harmonic wavelet scattering for predictions of molecule properties (based on Refs. [1, 2])
  - General idea of the method
  - Wavelet scattering coefficients
  - Comparison to usual force fields
  - Why does it work so well : comparison to the multipolar expansion in polarizable force fields methods

# General idea : regression of scattering coefficients (1)

- $\underbrace{\mathbf{x}}_{\text{molecule}} = \left\{ \left( \underbrace{r_k}_{\text{position}}, \underbrace{z_k}_{\text{nuclear-charge}} \right) \in \mathbb{R}^3 \times \mathbb{Z} \right\}_{k \in \llbracket 1, N \rrbracket} \mapsto f(\mathbf{x}) \in \mathbb{R}$
- $f(\mathbf{x})$  : physical property of interest of the molecule (e.g. total energy)
- $\mathbf{x}$  (all atomic positions and element types) : *possible* molecular descriptor (not unique : cf. drug design context)
- Approximation of  $f(\mathbf{x})$  by multilinear regression :

$$\tilde{f}_{a_1, \dots, a_q}(\mathbf{x}) = b + \underbrace{\sum_{i=1}^M \underbrace{a_i}_{\in \mathbb{R}} \mathcal{C}_{\mathbf{p}_i}[\mathbf{x}]}_{\text{linear}} \left( + \underbrace{\sum_{q=2}^r \left[ \sum_{i_1 < \dots < i_q} \underbrace{a_{i_1, \dots, i_q}}_{\in \mathbb{R}} \mathcal{C}_{\mathbf{p}_{i_1}}[\mathbf{x}] \dots \mathcal{C}_{\mathbf{p}_{i_q}}[\mathbf{x}] \right]}_{\text{multilinear}} \right) \quad (1)$$

- $\mathcal{C}_{\mathbf{p}_i}[\mathbf{x}]$  : (solid harmonic wavelet) scattering coefficient *i.e.* relevant *molecular descriptor*.

## General idea : regression of scattering coefficients (2)

- $(a_1^*, \dots, a_p^*) = \arg \min_{(a_1, \dots, a_p)} \underbrace{\sum_{\nu} \left| \tilde{f}_{a_1, \dots, a_p}(\mathbf{x}_{\nu}) - f(\mathbf{x}_{\nu}) \right|^2}_{\text{Training-set}} : (\text{locally})$   
optimal regression coefficients.
- Training set  $\left\{ \left( \underbrace{\mathbf{x}_{\nu}}_{\text{geometry}}, \underbrace{f(\mathbf{x}_{\nu})}_{\text{energy}} \right) \right\}_{\nu}$  : from ab-initio (DFT) calculations.
- Validation set : similar type of data (but different from the training set)

## General idea (3)

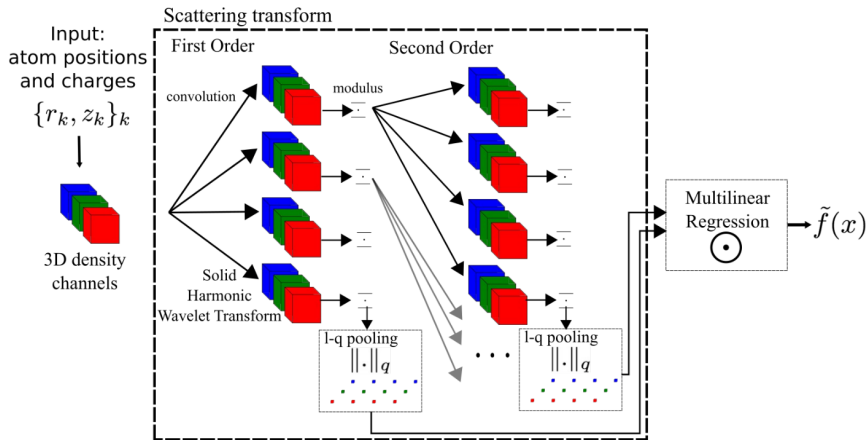


Figure: Convolutional network (two wavelet scattering transform steps) [1, 2]

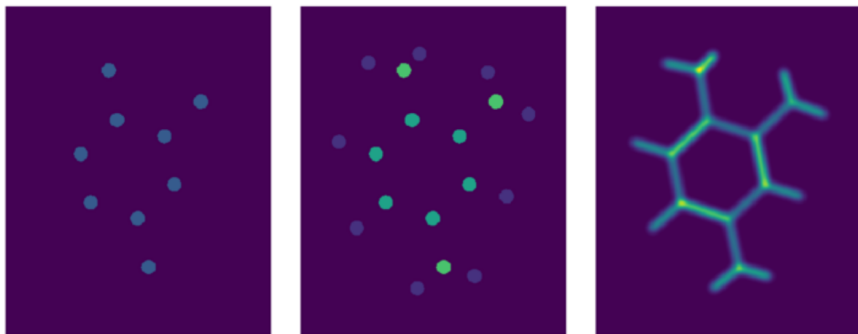
# Computation of wavelet scattering coefficients : outline (1)

- $\underbrace{x}_{\text{geometry}} \longrightarrow \underbrace{\rho_x(\cdot)}_{\text{density "channels"}} \longrightarrow \underbrace{S\rho_x}_{\text{wavelet-transform}} \quad (\text{set of coefficients})$
- Three density channels (if available in the dataset) : core / valence / bond density
- "Naïve" ("surrogate") core/valence density :

$$\rho_x^{\text{core/valence}}(\mathbf{u}) = \sum_k \gamma_k g(\mathbf{u} - \mathbf{r}_k) \quad (2)$$

- $g(\mathbf{v}) = Ke^{-\frac{v^2}{2\sigma^2}}$ ,  $\gamma_k$  number of core/valence electrons at atom  $k$  (so that  $\rho_x(\cdot)$  integrates to  $\sum_k \gamma_k$ )
- **Invariance on permutation of atom indexes.**

## Computation of wavelet scattering coefficients : outline (2)



**Figure:** Three density "channels" (or "guesses") possibly used as input to the convolutional neural network (core, valence and bond channel) [1, 2] – the "full" density channel is also added

# Computation of wavelet scattering coefficients : outline (3)

- Surrogate bond density :

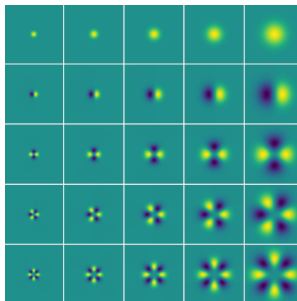
$$\rho_x^{bonds}(\mathbf{u}) = C \sum_{i \leftrightarrow j}^{bonds} \frac{\gamma_{ij}}{|\mathbf{r}_i - \mathbf{r}_j|} e^{-\frac{d_{ij}(\mathbf{u})^2}{2d_0^2}} \quad (3)$$

- $d_{ij}(\cdot)$  = distance to the bond (line)  $i - j$
- $\gamma_{ij}$  : number of electrons involved in the bond  $i - j$
- No prior, precise knowledge of the electronic density required ("rough" guess).



# Solid harmonic wavelets (1)

- Wavelets  $\psi(\cdot)$  : localized both in time and in frequency.
- Wavelets coefficients / decomposition (over a family of wavelets scaled and translated from a "mother" wavelet) : analogous to Fourier decomposition.
- More efficient to represent signal with discontinuities than Fourier.



**Figure:** Real parts of 2D solid harmonic wavelets  $\psi_{l,j}^{2D}$  ( $\psi_l^{2D}(r, \theta) = \frac{1}{2\pi} e^{-\frac{1}{2}r^2} r^l e^{il\phi}$ ) : **angular momentum**  $l = 0 \dots 4$  from top to bottom and **scale**  $j = 0 \dots 4$  from left to right [2]

## Solid harmonic wavelets (2)

- Solid harmonic wavelets in 3D :

$$\psi_l^m(\mathbf{u}) = \frac{1}{(\sqrt{2\pi})^3} e^{-\frac{|\mathbf{u}|^2}{2}} |\mathbf{u}|^l Y_l^m \left( \frac{\mathbf{u}}{|\mathbf{u}|} \right) \in \mathbb{C} \quad (4)$$

- The Fourier transform of a wavelet is a wavelet :

$$\hat{\psi}_l^m(\omega) = (-i)^l e^{-\frac{1}{2}|\omega|^2} |\omega|^l Y_l^m \left( \frac{\omega}{|\omega|} \right) \quad (5)$$

- Wavelets *scaled* at scale  $2^j$  :

$$\psi_{l,j}^m(\mathbf{u}) = \frac{1}{(2^j)^3} \psi_l^m \left( \frac{\mathbf{u}}{2^j} \right) \quad (6)$$

# Wavelet scattering transform : step 1 (1)

$$\bullet \quad \underbrace{\rho(\cdot)}_{\text{density-channel}} \xrightarrow{\underbrace{\quad}_{1^{\text{st}} \text{convolution}}} \rho * \psi_{l,j}^m \xrightarrow{\underbrace{\quad}_{\text{modulus-operator}}} U[j, l] \rho \longrightarrow S\rho[j, l, q]$$

$$U[j, l] \rho : \mathbf{u} \mapsto \sqrt{\sum_{m=-l}^l \left| (\rho * \psi_{l,j}^m)(\mathbf{u}) \right|^2} \quad (7)$$

$$\boxed{S\rho[j, l, q] = \int_{\mathbb{R}^3} |U[j, l] \rho(\mathbf{u})|^q d\mathbf{u} \in \mathbb{R}} \quad (8)$$

- **First order solid harmonic wavelet scattering coefficients.**
- **Rotational invariant.**
- Coefficients  $S\rho[j, l, q]$  for  $q = 2$  exponent : encode pairwise interactions (e.g. Coulomb) ? Why ? No correlations  $\rho(\mathbf{u})\rho(\mathbf{u}')$  ...
- Encode both short (small  $j$ ) and long-range (large  $j$ ) interactions contributions to the energy. Why ?

## Wavelet scattering transform : step 1 (2)

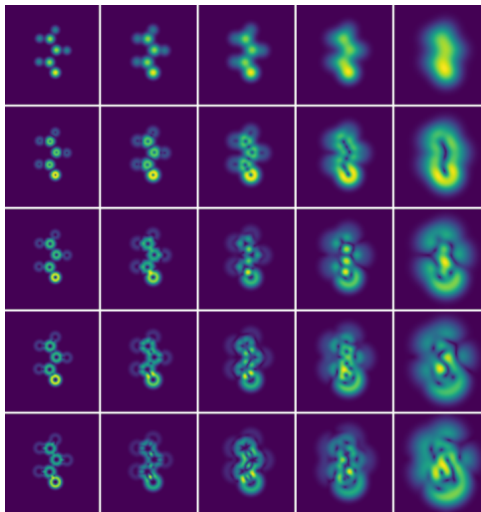


Figure: Solid harmonic wavelet scattering coefficient (moduli  $U[j, l]\rho$ ) [2] : reminiscent of interference patterns / molecular orbitals ?

## Wavelet scattering transform : step 1 (3)

- Modulus coefficients (functions)  $U[j, l]\rho(\cdot)$  are similar to multipole moments :

$$U[j, l]\rho(\mathbf{u}) = \sqrt{\sum_{m=-l}^l \left| \int_{\mathbb{R}^3} \rho(\mathbf{v}) \psi_{lj}^m(\mathbf{u} - \mathbf{v}) d\mathbf{v} \right|^2} \quad (9)$$

$$U[j, l]\rho(\mathbf{u}) = \frac{1}{2^{3j}} \frac{1}{2^{lj}} \sqrt{\sum_{m=-l}^l \left| \underbrace{\int_{\mathbb{R}^3} |\mathbf{v} - \mathbf{u}|^l Y_l^m \left( \frac{\mathbf{u} - \mathbf{v}}{|\mathbf{u} - \mathbf{v}|} \right) \rho(\mathbf{v}) e^{-\frac{1}{2} \left( \frac{|\mathbf{u} - \mathbf{v}|}{2^j} \right)^2} d\mathbf{v}}_{\rightarrow_{j \rightarrow +\infty} (Q_l^m)_{\mathbf{u}}^c} \right|^2} \quad (10)$$

- Multipole moments (of a charge distribution) of order  $l$  with respect to expansion center  $\mathbf{u}$  (using **real** spherical harmonics...) :

$$(Q_l^m)_{\mathbf{u}} = \int_{\mathbb{R}^3} |\mathbf{v} - \mathbf{u}|^l Y_l^m \left( \frac{\mathbf{v} - \mathbf{u}}{|\mathbf{v} - \mathbf{u}|} \right) \rho(\mathbf{v}) d\mathbf{v} \in \mathbb{R} \quad (11)$$

## Wavelet scattering transform : step 1 (4)

- In which sense modulus coefficients (functions)  $U[j, l]\rho(\cdot)$  are "analogous to localized multipole moments" [1] ?
- If  $\forall \mathbf{w} \in \{\mathbf{v} | \rho(\mathbf{v}) \neq 0\}$ ,  $|\mathbf{u} - \mathbf{w}| \ll 2^j$  (i.e. at a large enough scale):

$$U[j, l]\rho(\mathbf{u}) \approx \frac{1}{2^{3j}} \frac{1}{2^{lj}} \sqrt{\sum_{m=-l}^l |(Q_l^m)_{\mathbf{u}}^{\mathbb{C}}|^2} \quad (12)$$

where  $(Q_l^m)_{\mathbf{u}}^{\mathbb{C}}$  are defined analogously to (real) multipole moments but with complex spherical harmonics.

- Otherwise, for all  $j$ ,  $U[j, l]\rho(\mathbf{u})$  relates to the multipole moments (with respect to center  $\mathbf{u}$ ) of the charge distribution  $\rho(\cdot)$  localized by a gaussian.

## Wavelet scattering transform : step 2

- $U[j, l]\rho \xrightarrow{(*)} (U[j, l]\rho) * \psi_{l,j'}^m \xrightarrow{\text{modulus}} U[j', l](U[j, l]\rho) \stackrel{\mathcal{D}}{=} U[j, j', l]\rho$
- Step  $(*)$  : second convolution, at larger scales ( $j' > j$ )

$$U[j, j', l]\rho(\mathbf{u}) = U[j', l](U[j, l]\rho)(\mathbf{u}) = \left( \sum_{m=-l}^l |(U[j, l]\rho * \psi_{l,j'}^m)(\mathbf{u})|^2 \right)^{\frac{1}{2}} \quad (13)$$

- Final step :

$$U[j, j', l]\rho \longrightarrow S\rho[j, j', l, q] = \int_{\mathbb{R}^3} |U[j, j', l]\rho(\mathbf{u})|^q d\mathbf{u} \quad (14)$$

- **Second order solid harmonic wavelet scattering coefficients.**
- "Multiscale coupling" coefficient (substructures scales  $2^j$  and  $2^{j'}$ )  
*interpreted as van der Waals (dispersion) interaction terms*  $\propto C \frac{\alpha_1 \alpha_2}{R^6} \dots$

# Final regression

- Case of linear regression :

$$\tilde{f}_{\vec{p}}(\mathbf{x}) = b + \underbrace{\sum_j \sum_l \sum_q \overset{\text{Scale Moment}}{w_{l,q}^j} \underbrace{S\rho[j, l, q]}_{\text{local-terms+Coulomb?}}}_{\text{local-terms+Coulomb?}} + \underbrace{\sum_j \sum_{j' > j} \sum_l \sum_q \overset{\text{Scale Scale Moment}}{w_{l,q}^{j,j'}} \underbrace{S\rho[j, j', l, q]}_{\text{induction+dispersion?}}}_{\text{induction+dispersion?}} \quad (15)$$

- "Analogous to perturbation expansion" (cf. SAPT theory)
- $\vec{p} = \left\{ w_{l,q}^j, w_{l,q}^{j,j'} \right\}_{l=0..L, j=0..J, j < j', q=0..Q}$  : set of parameters to optimize
- In fact : sum also on the three different density channels  $\rho_x(.)$  ?
- Relatively small number of scattering coefficients.
- Bilinear regression : additional products of scattering coefficients  $S\rho[j_1, l_1, q_1] S\rho[j_2, j'_2, l_2, q_2]$ , more coefficients to fit.
- Linear vs. bilinear vs. trilinear regression ? Overfitting ? What is the right "functional form" to assume ?



# Comparison to usual "force fields"

- "Physical" functional form near equilibrium :

$$\begin{aligned}
 U_{\mathbf{p}}(\vec{r}_1, \dots, \vec{r}_N) = & \sum_{\text{bonds}} k_b (b - b_0)^2 + \sum_{\text{angles}} k_{\theta} (\theta - \theta_0)^2 + \sum_{\text{dihedrals}} k_{\phi} (1 + \cos(n\phi - \delta)) + \sum_{\text{impropers}} k_{\chi} (\chi - \chi_0)^2 \\
 & + \underbrace{\sum_{i < j | d(i,j) \geq M} c \frac{q_i q_j}{4\pi\epsilon r_{ij}}}_{\text{electrostatic}} + \underbrace{\sum_{i < j | d(i,j) \geq M} \epsilon_{ij} \left[ \left( \frac{\sigma_{ij}}{r_{ij}} \right)^{12} - \left( \frac{\sigma_{ij}}{r_{ij}} \right)^6 \right]}_{\text{dispersion(vdW)}}
 \end{aligned} \tag{16}$$

- Set of parameters to optimize :

$$\mathbf{p} = \left( \{k_b^{ij}\} \{b_0^{ij}\}, \{k_{\theta}^{ijk}\} \{\theta_0^{ijk}\}, \{k_{\phi}^{ijkl}\} \{n^{ijkl}\}, \{k^{abcd}\}, \{\chi_0^{abcd}\}, \{q_i\}, \{\epsilon_{ij}, \sigma_{ij}\}_{i < j | d(i,j) \geq 3} \right) \tag{17}$$

- Cheap computationnaly.
- Used to perform long simulations of molecules (molecular dynamics).

# Why does multilinear regression of wavelet scattering coefficients work so well (1) ?

- Electrostatic (**long-range** interaction) energy terms in (e.g.) **polarizable** force fields (multipolar expansion of the interaction energy between two multipoles  $[q_i, \vec{d}_i, \bar{\bar{Q}}^i]$  and  $[q_j, \vec{d}_j, \bar{\bar{Q}}^j]$  at atoms  $i$  and  $j$ ) :

$$\begin{aligned}
 & \underbrace{\frac{q_i q_j}{|\vec{r}_i - \vec{r}_j|}}_{\text{charge-charge}} - \underbrace{\frac{(q_j \vec{d}_i - q_i \vec{d}_j) \cdot (\vec{r}_i - \vec{r}_j)}{2 |\vec{r}_i - \vec{r}_j|^3}}_{\text{charge-dipole}} + \underbrace{\frac{\vec{d}_i \vec{d}_j}{|\vec{r}_i - \vec{r}_j|^3} - 3 \frac{[\vec{d}_i \cdot (\vec{r}_i - \vec{r}_j)] [\vec{d}_j \cdot (\vec{r}_i - \vec{r}_j)]}{|\vec{r}_i - \vec{r}_j|^5}}_{\text{dipole-dipole}} \\
 & - \underbrace{\frac{q_i \text{Tr}(\bar{\bar{Q}}^j) + q_j \text{Tr}(\bar{\bar{Q}}^i)}{|\vec{r}_i - \vec{r}_j|^3} + 3 \frac{q_j (\vec{r}_i - \vec{r}_j)^T \bar{\bar{Q}}^i (\vec{r}_i - \vec{r}_j) + q_i (\vec{r}_i - \vec{r}_j)^T \bar{\bar{Q}}^j (\vec{r}_i - \vec{r}_j)}{|\vec{r}_i - \vec{r}_j|^5}}_{\text{charge-quadrupole}}
 \end{aligned} \tag{18}$$

- $q_i$  (atomic charge),  $\vec{d}_i$  (dipole),  $\bar{\bar{Q}}^i$  (quadrupole) are (local) **molecular descriptors**.

# Why does multilinear regression of wavelet scattering coefficients work so well (2) ?

- Equation 18 is an approximation (at angular moment of order 2) of :

$$\frac{1}{2} \left( \int \rho_i(\vec{r}) V_j(\vec{r}) d\vec{r} + \int \rho_j(\vec{r}) V_i(\vec{r}) d\vec{r} \right) \quad (19)$$

- Bilinear** regression should work better than **linear** regression : products charge-dipole, dipole-dipole, charge-quadrupole in the "physical" expression of the electrostatic energy *i.e.* **products of two molecular descriptors** only ?
- Trilinear** should not (physically) give better results than **bilinear** regression – at least concerning description of long-range electrostatic interactions ?

# Discussion

- Transferability of the approximated function  $\tilde{f}(\mathbf{x})$  to larger molecules ? Increase number of scales ( $j$ ) ? ( $\implies$  more scattering terms)
- Transferability to other chemical contexts (e.g. chemical reaction  $\implies$  no training on the bond channel) ?
- Training data set : symmetric configurations of the molecule ? Or impose "by hand" the symmetries (e.g. by symmetry invariant coefficients) as here.
- Interpretation of the first/second order scattering coefficients as physical interaction terms (cf. short-range / long-range interactions) ?
- To what extent (and how) do scattering coefficients "account for different types of interactions at different scales" ?

Questions ?