

A Guide to the REG analysis

(and REG.py documentation)

Leonardo José Duarte.

REG.py
Version 0.1
March 2020

This is a guide on how to perform the REG analysis using the REG.py library. If you use it, please cite: Thacker, J.C.R. & Popelier, P.L.A. Theor Chem Acc (2017) 136: 86.

To report bugs, ask questions and send suggestions, please send an email to
leo.j.duarte@hotmail.com.br or l156213@unicamp.br

This manual was written under the supervision of Prof Dr. Paul Popelier.

REG.py is available to download at: github.com/ljduarte

Contents

1	The IQA Energy Decomposition Scheme	3
1.1	Intra- and Inter-atomic Contributions	3
1.2	Starting from the Reduced Density-Matrices	4
1.2.1	The atoms within the electronic density	5
1.3	Grouping Atoms	5
2	Relative Energy Gradient (REG) Approach	8
2.1	Potential Energy Surface: Splitting into Segments	8
2.2	The Additive Property of the REG coefficient	9
3	AIMALL Tips	11
3.1	AIMAll and DFT	11
3.2	AIMAll Output	13
4	REG.py	15
4.1	List of Functions - reg.py	15
4.2	List of Functions - aimall_utils.py	17
4.3	List of Functions - gaussian_utils.py	18
4.4	List of Functions - morphy_utils.py	19
4.5	List of Functions - reg_vis.py	19
4.6	Data Visualisation Examples	20
4.7	Workflow	22
4.8	Script Example	23
5	Appendix A	25
5.1	Defining the Control Coordinate - An S_N2 reaction	25
5.1.1	STEP 1 - Optimise reagents geometries	25
5.1.2	STEP 2 - Scan through the Potential Energy Surface (PES)	25
5.1.3	STEP 3 - Releasing the constraints I	26
5.1.4	STEP 4 - Releasing the constraints II	27
5.1.5	STEP 5 - Finding the IRC	29
6	Appendix B	30
6.1	Demonstrating $C = 0$	30

1 The IQA Energy Decomposition Scheme

From the basis of the Quantum Theory of Atom in Molecules (QTAIM), the Interacting Quantum Atoms (IQA) emerges as an energy decomposition scheme that split the molecular, or supramolecular, system energy, E_{IQA}^{sys} , into a sum of individual atomic energies composed of *intra*-atomic and *inter*-atomic energy components.

Such terms provide a very detailed description of a system’s electronic density, allowing the proper quantification of chemical interactions. The IQA approach has already been applied to a variety of chemical problems, such as S_N2 reactions, hydrogen-bonded systems, the fluorine gauche effect, the rotation barriers on biphenyl, etc. In the next sections, we are going to explore the fundamentals of IQA, presenting the principal equations and how to manipulate them for use in different situations.

1.1 Intra- and Inter-atomic Contributions

The total system energy is recovered by summing each individual atom energy, E_{IQA}^A , according to:

$$E_{IQA}^{sys} = \sum_{A=1}^N E_{IQA}^A \quad (1)$$

where N is the total number of atoms in the system and A is the atomic label. Each atomic term can be expanded as a sum of intra- and interatomic contributions. Equation 1 becomes:

$$E_{IQA}^{sys} = \sum_{A=1}^N E_{Intra}^A + \sum_{A=1}^N \sum_{B>A}^N V_{Inter}^{AB} \quad (2)$$

with V_{Inter}^{AB} being the potential energy between atoms A and B . Alternatively, the double summation on the right-hand side of Equation 2 can be rewritten as:

$$\sum_{A=1}^N \sum_{B>A}^N V_{Inter}^{AB} \equiv \frac{1}{2} \sum_{A=1}^N \sum_{B \neq A}^N V_{Inter}^{AB} \quad (3)$$

it is important to keep such definition in mind, since both forms are used by programs like AIMAll and MORPHY.

The intratomic term can be split in its kinetic, T , electron-electron, V_{ee} , and electron-nucleus, V_{en} , potential energies:

$$E_{Intra}^A = T^A + V_{ee}^{AA} + V_{en}^{AA} \quad (4)$$

E_{intra}^A is a measure of the intrinsic stability of an atom in the system and behaves as the familiar steric repulsion. The total intratomic energy is obtained directly from the three-dimensional integration of the electronic density over the topological atom volume (i.e. atomic basin). The interatomic contribution can also be split into several terms:

$$V_{Inter}^{AB} = V_{nn}^{AB} + V_{en}^{AB} + V_{ne}^{AB} + V_{ee}^{AB} \quad (5)$$

where the subscripts e and n stand for *electron* and *nucleus* respectively. Note that the order of the superscripts and subscripts does matter. V_{en}^{AB} is equivalent to “potential energy between the electrons from atom A and the nucleus of atom B ”, whereas V_{ne}^{AB} means “potential energy between the nucleus of atom A and electrons from atom B ”. As a general reminder: $V_{en}^{AB} \neq V_{ne}^{AB} \equiv V_{en}^{BA} \neq V_{ne}^{BA}$.

The electron-electron potential energy can also be subdivided in two contributions: *classic Coulombic energy* and the *exchange-correlation interaction*, according to:

$$V_{ee}^{AB} = V_{coul}^{AB} + V_{xc}^{AB} \quad (6)$$

V_{xc}^{AB} is associated to the degree of covalency between the atoms A and B . All classic terms can be grouped into the classical term: $V_{cl}^{AB} = V_{nn}^{AB} + V_{en}^{AB} + V_{ne}^{AB} + V_{coul}^{AB}$. Using this definition, Equation 5 can be written as:

$$V_{Inter}^{AB} = V_{cl}^{AB} + V_{xc}^{AB} \quad (7)$$

All interatomic terms are obtained direct from a six-dimension integration of the electronic density over both, A and B , topological atoms volume.

1.2 Starting from the Reduced Density-Matrices

Once the desired multi-electronic wavefunction, Ψ , is found, the first and second-order density-matrices can be obtained following:

$$\rho_1(\mathbf{r}_1, \mathbf{r}_{1'}) = N_e \int \Psi(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N_e}) \Psi^*(\mathbf{x}'_1, \mathbf{x}_2, \dots, \mathbf{x}_{N_e}) d\mathbf{x}_2 d\mathbf{x}_3 \dots d\mathbf{x}_{N_e} \quad (8)$$

and;

$$\rho_2(\mathbf{r}_1, \mathbf{r}_2) = \frac{N_e(N_e - 1)}{2} \int \Psi(\mathbf{x}_1, \dots, \mathbf{x}_{N_e}) \Psi^*(\mathbf{x}_1, \dots, \mathbf{x}_{N_e}) d\mathbf{x}_3, \dots, d\mathbf{x}_{N_e} \quad (9)$$

where \mathbf{x}_n is the vector of the spatial, denoted by \mathbf{r}_n and spin coordinates for electron n , N_e is the total number of electrons.

According to the Born-Oppenheimer approximation, the multi-electronic Hamiltonian, in atomic units, is given by:

$$\hat{H}_{N_e} = \hat{T} + \hat{V}_{en} + \hat{V}_{ee} = - \sum_{n=1}^{N_e} \frac{1}{2} \nabla_n^2 - \sum_{n=1}^{N_e} \sum_{A=1}^N \frac{Z_A}{r_{nA}} + \sum_{n=1}^{N_e} \sum_{l>n}^{N_e} \frac{1}{r_{nl}} \quad (10)$$

where \hat{T} is the one-electron kinetic energy operator. The total energy is given by:

$$E_{total} = E_{electronic} + V_{nn} \quad (11)$$

the potential energy between two nuclei is simply defined as:

$$V_{nn} = \sum_{A=1}^N \sum_{B>A}^N V_{nn}^{AB} = \sum_{A=1}^N \sum_{B>A}^N \frac{Z_A Z_B}{r_{AB}} \quad (12)$$

the electronic energy can be obtained solving the Schrödinger equation

$$\hat{H}_{N_e} \Psi(\mathbf{x}_1, \dots, \mathbf{x}_{N_e}) = E_{Electronic} \Psi(\mathbf{x}_1, \dots, \mathbf{x}_{N_e}) \quad (13)$$

which results in

$$E_{electronic} = \int_{\infty} \hat{T} \rho_1(\mathbf{r}_1, \mathbf{r}_{1'}) d\mathbf{x}_1 + \int_{\infty} \hat{V}_{en} \rho_1(\mathbf{r}_1, \mathbf{r}_{1'}) d\mathbf{x}_1 + \int_{\infty} \int_{\infty} \hat{V}_{ee} \rho_2(\mathbf{r}_1, \mathbf{r}_2) d\mathbf{x}_1 d\mathbf{x}_2 \quad (14)$$

In order to compute individual terms for each atom the AIM topological partitioning method is invoked.

1.2.1 The atoms within the electronic density

Using the vector field traced by the gradient, $\nabla\rho(\mathbf{r})$, of the electronic density, one obtains the so-called interatomic surfaces (denoted S), which are boundaries of an atomic basin (denoted Ω) and are defined by:

$$\nabla\rho(\mathbf{r}) \cdot \mathbf{n}(\mathbf{r}) = 0 \quad (15)$$

where $\mathbf{n}(\mathbf{r})$ is the vector normal to the interatomic surface, Ω . With all the atomic basins in the system properly identified, the atomic properties are extracted by integrating over the atomic space. The kinetic energy contribution for atom A is given by :

$$T^A = \int \hat{T} \rho_1^A(\mathbf{r}) d\mathbf{r} \quad (16)$$

The mono-electronic inter-atomic terms are obtained as follows:

$$V_{en}^{AB} = \int_{\Omega_A} \hat{V}_{en}^B \rho_1(\mathbf{r}_1, \mathbf{r}_1') d\mathbf{x}_1 = - \int_{\Omega_A} \frac{\rho_1(\mathbf{r}_1) Z^B}{r_{1B}} d\mathbf{r}_1 \quad (17)$$

The two-electron terms are obtained from the second-order density-matrix, which describes how electrons interact with each other:

$$V_{ee}^{AB} = \int_{\Omega_A} \int_{\Omega_B} \rho_2(\mathbf{r}_1, \mathbf{r}_2) r_{12}^{-1} d\mathbf{x}_1 d\mathbf{x}_2 \quad (18)$$

if $A = B$, then:

$$V_{ee}^{AA} = \frac{1}{2} \int_{\Omega_A} \int_{\Omega_A} \rho_2(\mathbf{r}_1, \mathbf{r}_2) r_{12}^{-1} d\mathbf{x}_1 d\mathbf{x}_2 \quad (19)$$

Equation 18 contains the Coulombic, exchange and correlation terms, which can be individually calculated by invoking the three types contributions to ρ_2 :

$$\begin{aligned} V_{ee}^{AB} &= \int_{\Omega_A} \int_{\Omega_B} \rho_2(\mathbf{r}_1, \mathbf{r}_2) r_{12}^{-1} d\mathbf{x}_1 d\mathbf{x}_2 = \\ &\int_{\Omega_A} \int_{\Omega_B} \rho(\mathbf{r}_1) \rho(\mathbf{r}_2) r_{12}^{-1} d\mathbf{r}_1 d\mathbf{r}_2 - \int_{\Omega_A} \int_{\Omega_B} \rho_1(\mathbf{r}_1, \mathbf{r}_2) \rho_1(\mathbf{r}_2, \mathbf{r}_1) r_{12}^{-1} d\mathbf{r}_1 d\mathbf{r}_2 + \\ &\int_{\Omega_A} \int_{\Omega_B} \rho_2^{corr} r_{12}^{-1} d\mathbf{r}_1 d\mathbf{r}_2 = V_{coul}^{AB} + V_{ex}^{AB} + V_{corr}^{AB} \end{aligned} \quad (20)$$

Exchange and correlation terms are typically lumped together in the exchange-correlation term. $V_{ex}^{AB} + V_{corr}^{AB} = V_{xc}^{AB}$. V_{ex}^{AB} is related to the covalency degree between atoms A and B , as well as the bond order. The classical electrostatic term, V_{coul}^{AB} , on the other hand, is related to the bond polarity and ionicity degree. The last term, V_{corr}^{AB} is responsible to increase the magnitude of the nucleus-electron potential energy and decreases the electron-electron repulsion. All these terms constitute the so-called *fine structure* of ρ_2 .

1.3 Grouping Atoms

Consider a system formed by N atoms divided into three groups J , K and L . The total energy of the system is recovered as follow:

$$E_{IQA}^{sys} = \sum_{X=1}^3 E_{IQA}^X = E_{IQA}^J + E_{IQA}^K + E_{IQA}^L \quad (21)$$

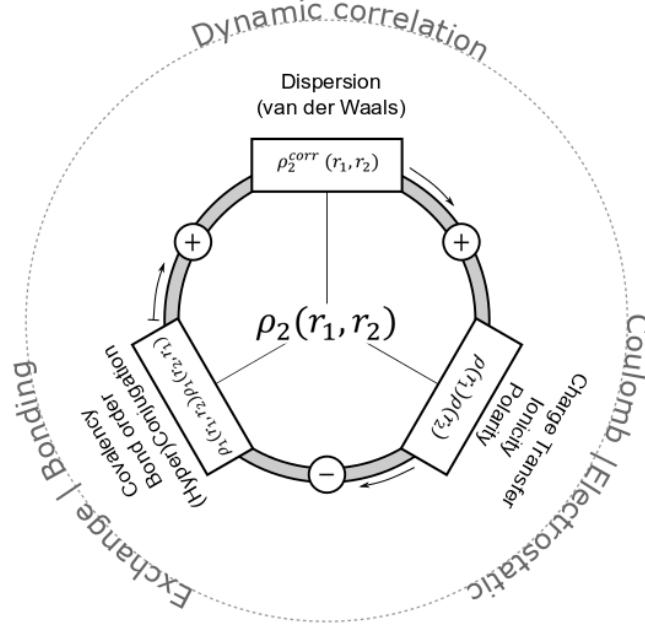


Figure 1: Representation of the fine structure of the second-order density-matrix ρ_2 . Each term is presented with its physical meaning. The “+” and “-” signs refer to the mathematical relationship between terms as shown in 20.

where $X = 1, 2$ or 3 correspond to J, K or L group. Each group contribution can be expanded in **intratomic-intragroup** ($Intra^2$), **interatomic-intragroup** ($Inter/Intra$) and **interatomic-intergroup** ($Inter^2$) terms.

The $Intra^2$ term corresponds to the sum of all intratomic energies of atoms that belong to group X and is defined as:

$$E_{Intra^2}^X = \sum_{A=1}^N E_{Intra}^A \quad \text{for } A \in X \quad (22)$$

$Intra^2$ is associated with the intrinsic stability of the group and the steric effects.

The $Inter^2$ term is equivalent to the sum of the interatomic energy between two atoms of distinct groups,

$$E_{Inter^2}^X = \sum_{A=1}^N \sum_{B>A}^N V_{Inter}^{AB} \quad \text{for } A \in X \text{ and } B \notin X \quad (23)$$

alternatively, using the definition in Equation 3, one can write:

$$E_{Inter^2}^X = \frac{1}{2} \sum_{A=1}^N \sum_{B \neq A}^N V_{Inter}^{AB} \quad \text{for } A \in X \text{ and } B \notin X \quad (24)$$

Similarly to individual atomic terms, it is possible to split $Inter^2$ in classical and exchange-correlation interaction contributions. The chemical insights of the individual atomic terms are preserved and apply to $E_{Inter^2}^X$.

The $E_{Inter/Intra}^X$ is defined in the same way as $E_{Inter^2}^X$ with one exception: the atoms A and B MUST belong to the same group X .

$$E_{Inter/Intra}^X = \sum_{A=1}^N \sum_{B>A}^N V_{Inter}^{AB} \quad \text{for } \{A, B\} \subset X \quad (25)$$

the alternative definition is also valid:

$$E_{Inter/Intra}^X = \frac{1}{2} \sum_{A=1}^N \sum_{B \neq A}^N V_{Inter}^{AB} \quad for \{A, B\} \subset X \quad (26)$$

The principal advantage of using the grouping approach to partition the energetic terms consists of preserving the chemical information in a small number of terms.

2 Relative Energy Gradient (REG) Approach

In the last section, the IQA partition scheme was presented. It should be clear that, for a small system formed by ten atoms, a modest IQA partition will result in a set of $10 \times E_{intra}^A$ atomic terms, $45 \times V_{cl}^{AB}$ interatomic classical terms and $45 \times V_{xc}^{AB}$ interatomic exchange-correlation terms. With the increase of the partitioning level and/or the increase of the size of the system, the number of terms will explode and a manual analysis of such data becomes impractical, not to say impossible. The problem becomes worst when dealing with an energy surface, where the IQA terms need to be analysed over a coordinate change.

To overcome this adversity, Thacker and Popelier proposed the Relative Energy Gradient method, which compares the derivatives of each contribution with the overall energy variation over a control coordinate and ranks them in order to find the set of contributions that preserve the majority of chemical insights into the total energy profile. As stated by Thacker and Popelier, the REG method has two objectives:

1. Determine the subset of energy contributions that preserves the chemical insights and best describes the overall profile of the total energy over any coordinate changes.
2. Extract chemical insight from the partitioned total energy.

2.1 Potential Energy Surface: Splitting into Segments

In a Potential Energy Surface (PES), each maximum critical point corresponds to an *energy barrier* separating two minimum critical points. However, the word "barrier" in this context is associated with an increase of energy, and energy segments go down when analysed from left to right. The sign of m_{REG} is independent of the direction of analysis (left to right or right to left). It is better to refer to an energy *segment* in a neutral way, hence the word "segment" instead of "barrier". A segment ideally, begins and ends at a critical (minimum, maximum and/or saddle points) or at a lateral limit of the function.

Consider the Lennard-Jones potential plot in Figure 2. The critical point corresponds to the equilibrium distance between two particles. From the critical point, the potential curve is divided in two segments, one at the left side of the critical point (segment 1), and one at the right (segment 2).

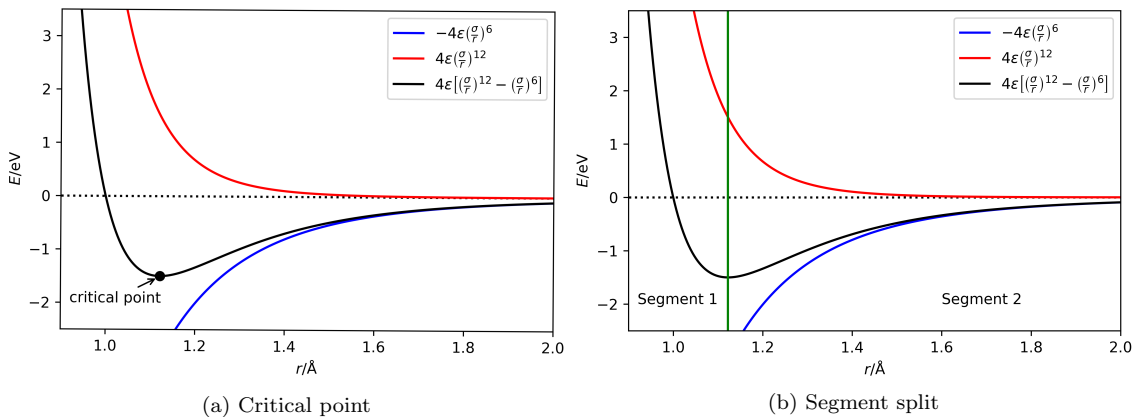


Figure 2: Lennard-Jones potential where $\sigma = 1.00$ and $\epsilon = 1.50$. The blue, red and black lines correspond to the attractive contribution, repulsive contribution and the total energy respectively. In Figure 2a is the critical point corresponds to the equilibrium point of the system. In Figure 2b, the green line split the potential curve in two segments

In this example, segment 1 is dominated by the Pauli repulsion term, $(\frac{\sigma}{r})^{12}$, while segment 2 is dominated by the dispersion interaction, $(\frac{\sigma}{r})^6$. A way to visualise the aforementioned behaviour is to compare, for each segment, the ratio between each component

derivative and the total energy derivative with respect to changes in the control coordinate, calculating the REG coefficient, m_{REG} , as follows:

$$\frac{dE_i}{ds} = m_{REG,i} \frac{dE_{total}}{ds} \quad (27)$$

The greater the value of $m_{REG,i}$ the more important the contribution E_i to the overall behaviour of the potential energy surface.

However, Equation 27 is valid only for a perfect correlation. For all other cases, m_{REG} is approximated as the angular coefficient of the linear regression between E_i and E_{total} , as given:

$$E_i(s) = m_{REG,i} E_{total}(s) + c_i \quad (28)$$

where:

$$m_{REG,i} = \frac{n \sum_{s=1}^n (E_{total}(s) E_i(s)) - \sum_{s=1}^n E_{total}(s) \sum_{s=1}^n E_i(s)}{n \sum_{s=1}^n E_{total}^2(s) - (\sum_{s=1}^n E_{total}(s))^2} \quad (29)$$

where n is the number of points in the PES and s is the control coordinate. c_i is the intercept, which is not relevant for the REG analysis. The approximation in 29 is better the closer to 1 the currently absolute value of the Pearson coefficient, r , is.

The m_{REG} values in Figure 3 confirms the statement that the segment 1 is dominated by the Pauli repulsion term, while the segment 2 is driven by the dispersion attractive term.

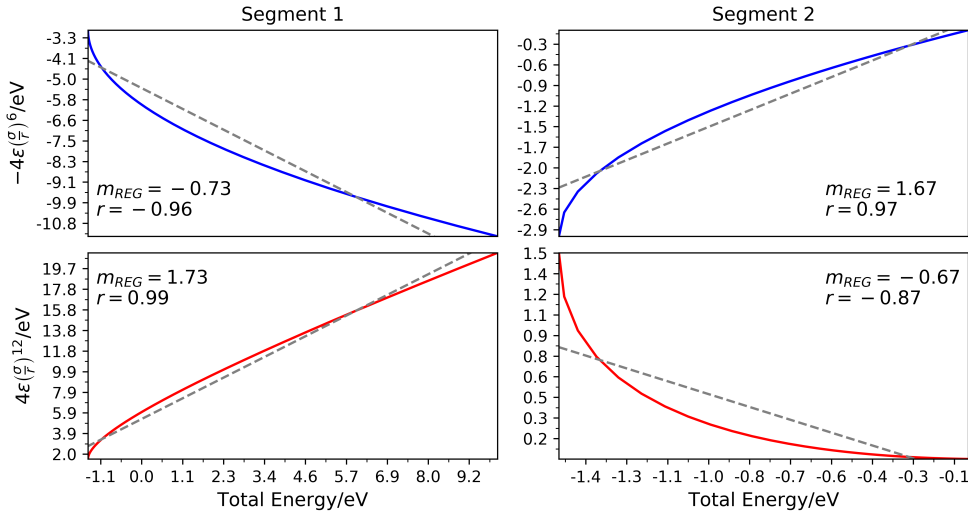


Figure 3: m_{REG} and r values for the attractive and repulsive contributions of the Lennard-Jones potential for both segments. The dashed gray lines correspond to the least square linear regression between the component and the total energy.

2.2 The Additive Property of the REG coefficient

The additive property of the REG coefficient states that if E_{total} is equal to a sum of m energy terms, $E_{total} = \sum_{l=0}^m E_l$, then $\sum_{l=0}^m m_{REG,l} = 1$. This property is easily derived from Equation 27 since the derivative of a function that is equal to a sum of other functions equals the sum of their derivatives.

$$\frac{dE_{total}}{ds} = \frac{d \sum_{l=0}^m E_l}{ds} = \sum_{l=0}^m \frac{dE_l}{ds} \quad (30)$$

then from Equation 27, it follows that

$$\sum_{l=0}^m \frac{dE_l}{ds} = \sum_{l=0}^m m_{REG,l} \frac{dE_{total}}{ds} \quad (31)$$

which leads to:

$$\sum_{l=0}^m m_{REG,l} = \frac{\sum_{l=0}^m \frac{dE_l}{ds}}{\frac{dE_{total}}{ds}} = \frac{\frac{dE_{total}}{ds}}{\frac{dE_{total}}{ds}} = 1 \quad (32)$$

The same is valid for the approximation in Equation 28, since:

$$\sum_{l=0}^m E_l(s) = \sum_{l=0}^m m_{REG,l} E_{total}(s) + \sum_{l=0}^m c_l \quad (33)$$

$$E_{total}(s) = E_{total}(s) \sum_{l=0}^m m_{REG,l} + C \quad (34)$$

since, $C = 0$, resulting in (see section: 6):

$$\sum_{l=0}^m m_{REG,l} = 1 \quad (35)$$

This property is useful when grouping IQA terms. As a consequence of the additive property, the REG coefficient of a group of terms equals the sum of the REG coefficient of each term, or: *The m_{REG} of a sum of each contribution is equal to the sum of each individual m_{REG} value.*

3 AIMALL Tips

AIMALL is a program used to integrate the wavefunction in order to extract the QTAIM properties of the system. Given a system of N atoms, AIMAll will create N ".int" files that contains all monoatomic properties and, if requested, $\frac{N^2-N}{2}$ ".int" files containing the interatomic properties. If you are running AIMAll via a command line, you may use the command **-encomp=4** to request all the IQA terms necessary to the REG analysis. An example is given:

```
/mnt/iusers01/pp01/n48687ld/AIMAll/aimqb.ish -nogui -encomp=4 -nproc=8  
-delmog=false wavefunction.wfn
```

The **delmog** option is set to false, in order to prevent .mog files from being deleted. MOG files are used to reduce the integration error if necessary, by performing an integration with a better grid over the atoms that were poorly integrated.

To improve the integration grid, the following command line can be used:

```
/mnt/iusers01/pp01/n48687ld/AIMAll/aimqb.ish -nogui -encomp=4 -nproc=8  
-bim=proaim -boaq=gs60 -iasmesh=superfine -delmog=false filename.wfn
```

-bim determines which integration method will be used, **-boaq** will define the basin outer angular quadrature and **-iasmesh** will define the integration mesh spacing.

To check the integration error in a specific atom, you may search for the L value in that atom's .int file. Usually a $L \leq 10^{-4}$ is sufficiently good. However, if you need to reintegrate a specific atom, you can use the following command line:

```
/mnt/iusers01/pp01/n48687ld/AIMAll/aimqb.ish -nogui -encomp=4 -nproc=8  
-bim=proaim -boaq=gs60 -iasmesh=superfine -delmog=false -atoms=1,2,4  
filename.wfn
```

In this example, AIMAll will run the integration over atom 1, 2 and 4. For more information about AIMAll command lines, please check the AIMAll documentation page.

3.1 AIMAll and DFT

The following DFT methods are implemented in AIMAll:

- LSDA
- B3LYP
- M06-2X
- PBE
- PBE0

However, when using DFT methods, the wavefunction file needs to be modified. In the first line of a .wfn or .wfx file, one can read:

```
GAUSSIAN          9 MOL ORBITALS    106 PRIMITIVES      8 NUCLEI
```

in order to use the DFT method, you need to write the method *exactly three spaces* after "NUCLEI". An example for B3LYP is given below:

Modifying all wfn (or wfx) files manually is not recommended, since a small mistake can waste a huge amount of time. It is recommended using the following **sed** command in a terminal to modify the .wfn file:

```
$sed -i 's/NUCLEI/NUCLEI METHOD/g' filename.wfn
```

changing METHOD to LSDA, B3LYP, M06-2X, PBE or PBE0. If the wavefunction file is not modified properly, the integration error will be huge (around 200 kJmol^{-1}).

At end of any .int file, there are three important notes that the user should be aware of. They are reproduced below:

2EDM Note:

For post-HF, natural orbital "wavefunctions" the Muller approximation of the two-electron density matrix (2EDM) in terms of natural orbitals of the one-electron density matrix (1EDM) is used to calculate 2EDM-dependent properties (i.e., Vee contributions and electronic localization and delocalization properties). For HF single-determinant wavefunctions the expression used for the 2EDM is exact. 2EDM-dependent properties are calculated with spin-orbital "self-interaction" terms included.

DFT Note:

VeeX represents the exchange-correlation functional for the wavefunction's underlying model. For Hartree-Fock wavefunctions, VeeX is the two-electron Hartree-Fock exchange functional. For wavefunctions of supported DFT models (currently LSDA, B3LYP and M062X), VeeX is the actual exchange-correlation functional of the corresponding DFT model and atomic contributions VeeX(A) to VeeX are explicitly calculated and unambiguous. However, since VeeX is at least partly just a one-electron functional for DFT models, the partitioning of VeeX(A) into interatomic (VeeX(A,B) and VeeX(A,A')) and intraatomic (VeeX(A,A)) contributions is ambiguous. Currently, interatomic contributions VeeX(A,B) and VeeX(A,A') are calculated using the Hartree-Fock exchange functional while the intratomic contribution VeeX(A,A) is calculated as VeeX(A) - VeeX(A,A'). For wavefunctions of non-supported DFT models, VeeX is (incorrectly) assumed to be the Hartree-Fock exchange functional and the atomic energies E_IQA(A) will not be correct and will not sum to the correct molecular energy.

Vnn(A) - W(A) Note:

For Virial-Based atomic energies, Vnn(A) - W(A) is the contribution of atom A to the nuclear repulsion energy, Vnn(A), minus the contribution of atom A to the nuclear virial of the energy-gradient-based forces on the nuclei, W(A). Vnn(A) - W(A) is defined by the 3-term sum given in equation (6.85.5) of R.F.W. Bader's 1990 AIM book. For approximate wavefunctions Vnn(A,Mol) - W(A) is origin-dependent to the extent that the atomic Ehrenfest force theorem is not satisfied. However, an atom-specific origin can always be defined so that the atomic virial theorem is satisfied. Each Vnn(A,Mol) - W(A) value given here is determined using such an origin and is equivalent to the following definition:

$$\text{Vnn(A) - W(A)} = -2\text{T(A)} - \text{Vee(A)} - \text{Ven(A)}.$$

For stationary point geometries, such as equilibrium and transition-state geometries,

$W(A)$ is zero and the total atomic energy equals the atomic electronic energy, $E(A) = E_e(A)$.

3.2 AIMAll Output

The AIMAll output presents all the IQA terms presented in the section 1 of this manual, however the notation used by the program is different from the notation utilized in our papers. Table 1 shows the correspondence between the terms:

AIMAll Notation	Our Notation
$T(A)$	T^A
$V_{neen}(A,A)/2 = V_{ne}(A,A)$	V_{en}^{AA}
$V_{ne}(A,Mol)/2$	$\frac{1}{2} \sum_{B=0}^N V_{ne}^{AB}$
$V_{en}(A,Mol)/2$	$\frac{1}{2} \sum_{B=0}^N V_{en}^{AB}$
$V_{neen}(A,Mol)/2$	$\frac{1}{2} \sum_{B=0}^N V_{ne}^{AB} + \frac{1}{2} \sum_{B=0}^N V_{en}^{AB}$
$V_{ne}(A,A')/2$	$\frac{1}{2} \sum_{B \neq A}^N V_{ne}^{AB}$
$V_{en}(A,A')/2$	$\frac{1}{2} \sum_{B \neq A}^N V_{en}^{AB}$
$V_{neen}(A,A')/2$	$\frac{1}{2} \sum_{B \neq A}^N V_{ne}^{AB} + \frac{1}{2} \sum_{B \neq A}^N V_{en}^{a,b}$
$V_{ee}(A,A) + V_{ee}(A,A')/2$	$V_{en}^{AA} + \frac{1}{2} \sum_{B \neq A}^N V_{ee}^{AB}$
$V_{eeC}(A,A) + V_{eeC}(A,A')/2$	$V_{coul}^{AA} + \frac{1}{2} \sum_{B \neq A}^N V_{coul}^{AB}$
$V_{eeX}(A,A) + V_{eeX}(A,A')/2$	$V_{xc}^{AA} + \frac{1}{2} \sum_{B \neq A}^N V_{xc}^{AB}$
$V_{ee}(A,A)$	V_{ee}^{AA}
$V_{eeC}(A,A)$	V_{coul}^{AA}
$V_{eeX}(A,A)$	V_{xc}^{AA}
$V_{ee}(A,A')/2$	$\frac{1}{2} \sum_{B \neq A}^N V_{ee}^{AB}$
$V_{eeC}(A,A')/2$	$\frac{1}{2} \sum_{B \neq A}^N V_{coul}^{AB}$
$V_{eeX}(A,A')/2$	$\frac{1}{2} \sum_{B \neq A}^N V_{xc}^{AB}$
$V_IQA(A)$	$\sum_{B=0}^N V_{cl}^{AB} + \sum_{B=0}^N V_{xc}^{AB}$
$VC_IQA(A)$	$\sum_{B=0}^N V_{cl}^{AB}$
$VX_IQA(A)$	$\sum_{B=0}^N V_{xc}^{AB}$
$V_IQA(A,A)$	$V_{coul}^{AA} + V_{xc}^{AA}$
$VC_IQA(A,A)$	V_{cl}^{AA}
$VX_IQA(A,A)$	V_{xc}^{AA}
$V_IQA(A,A')/2$	$\frac{1}{2} \sum_{B \neq A}^N V_{cl}^{AB} + \frac{1}{2} \sum_{B \neq A}^N V_{xc}^{AB}$
$VC_IQA(A,A')/2$	$\frac{1}{2} \sum_{B \neq A}^N V_{cl}^{AB}$
$VX_IQA(A,A')/2$	$\frac{1}{2} \sum_{B \neq A}^N V_{xc}^{AB}$
$E_IQA(A)$	$T^A + V_{ee}^{AA} + V_{en}^{AA} + \frac{1}{2} \sum_{B \neq A}^N V_{coul}^{AB} + \frac{1}{2} \sum_{B \neq A}^N V_{xc}^{AB}$
$E_IQA_Intra(A)$	$E_{Intra}^A = T^A + V_{ee}^{AA} + V_{en}^{AA}$
$E_IQA_Inter(A)$	$\frac{1}{2} \sum_{B \neq A}^N V_{coul}^{AB} + \frac{1}{2} \sum_{B \neq A}^N V_{xc}^{AB}$
$V_{ne}(A,B)/2$	$\frac{1}{2} V_{ne}^{AB}$
$V_{en}(A,B)/2$	$\frac{1}{2} V_{en}^{AB}$
$V_{neen}(A,B)/2$	$\frac{1}{2} V_{ne}^{AB} + \frac{1}{2} V_{en}^{AB}$
$V_{ee}(A,B)/2$	$\frac{1}{2} V_{ee}^{AB}$
$V_{nn}(A,B)/2$	$\frac{1}{2} V_{nn}^{AB}$
$V_{eeC}(A,B)/2$	$\frac{1}{2} V_{coul}^{AB}$
$V_{eeX}(A,B)/2$	$\frac{1}{2} V_{xc}^{AB}$
$V_IQA(A,B)/2$	$\frac{1}{2} V_{Inter}^{AB} = \frac{1}{2} V_{cl}^{AB} + \frac{1}{2} V_{xc}^{AB}$
$VC_IQA(A,B)/2$	$\frac{1}{2} V_{cl}^{AB}$
$VX_IQA(A,B)/2$	$\frac{1}{2} V_{xc}^{AB}$

E_IQA_Inter(A,B)/2	$\frac{1}{2}V_{Inter}^{AB} = \frac{1}{2}V_{cl}^{AB} + \frac{1}{2}V_{xc}^{AB}$
Vne(B,A)/2	$\frac{1}{2}V_{ne}^{BA}$
Ven(B,A)/2	$\frac{1}{2}V_{en}^{BA}$
Vneen(B,A)/2	$\frac{1}{2}V_{ne}^{BA} + \frac{1}{2}V_{en}^{BA}$
Vee(B,A)/2	$\frac{1}{2}V_{ee}^{BA}$
Vnn(B,A)/2	$\frac{1}{2}V_{nn}^{BA}$
VeeC(B,A)/2	$\frac{1}{2}V_{coul}^{BA}$
VeeX(B,A)/2	$\frac{1}{2}V_{xc}^{BA}$
V_IQA(B,A)/2	$\frac{1}{2}V_{Inter}^{BA} = \frac{1}{2}V_{cl}^{BA} + \frac{1}{2}V_{xc}^{BA}$
VC_IQA(B,A)/2	$\frac{1}{2}V_{cl}^{BA}$
VX_IQA(B,A)/2	$\frac{1}{2}V_{xc}^{BA}$
E_IQA_Inter(B,A)/2	$\frac{1}{2}V_{Inter}^{BA} = \frac{1}{2}V_{cl}^{BA} + \frac{1}{2}V_{xc}^{BA}$
Vne(A,B)	V_{ne}^{AB}
Ven(A,B)	V_{en}^{AB}
Vneen(A,B)	$V_{ne}^{AB} + V_{en}^{AB}$
Vee(A,B)	V_{ee}^{AB}
Vnn(A,B)	V_{nn}^{AB}
VeeC(A,B)	V_{coul}^{AB}
VeeX(A,B)	V_{xc}^{AB}
V_IQA(A,B)	$V_{cl}^{AB} + V_{xc}^{AB}$
VC_IQA(A,B)	V_{cl}^{AB}
VX_IQA(A,B)	V_{xc}^{AB}
E_IQA_Inter(A,B)	$V_{Inter}^{AB} = V_{cl}^{AB} + V_{xc}^{AB}$
Vne(B,A)	V_{ne}^{BA}
Ven(B,A)	V_{en}^{BA}
Vneen(B,A)	$V_{ne}^{BA} + V_{en}^{BA}$
Vee(B,A)	V_{ee}^{BA}
Vnn(B,A)	V_{nn}^{BA}
VeeC(B,A)	V_{coul}^{BA}
VeeX(B,A)	V_{xc}^{BA}
V_IQA(B,A)	$V_{cl}^{BA} + V_{xc}^{BA}$
VC_IQA(B,A)	V_{cl}^{BA}
VX_IQA(B,A)	V_{xc}^{BA}
E_IQA_Inter(B,A)	$V_{Inter}^{BA} = V_{cl}^{BA} + V_{xc}^{BA}$

Table 1: IQA terms from AIMAll output

4 REG.py

REG.py is a python library that contains all functions needed to extract the IQA terms from the AIMAll outputs and perform a REG analysis. REG.py is an alternative to the ANANKE program and REG.py's advantage relies on its modular characteristics, i.e. the user can organize the functions according to his/her needs. There are five modules in the REG.py library:

1. **reg.py**: contains functions to perform the REG analysis.
2. **aimall_utils.py**: contains functions to get data from AIMAll outputs.
3. **gaussian_utils.py**: contains functions to get data from .wfn/.wfx files generated by GAUSSIAN09 or GAUSSIAN16
4. **morph_utils.py**: contains functions to get data from MORPHY/MORFI outputs (MORFI is used to split V_{xc} into $V_{correlation} + V_{exchange}$. See MORFI manual for more details.
5. **reg_vis.py** : contains functions that built some graphics to display the REG results.

The functions in each module are listed in the following sections.

4.1 List of Functions - reg.py

- **Function:** regression(A, B, mode=None)
Perform a linear regression between A and B ($B = \text{slope} \cdot A + \text{intercept}$)
Input: A, B and mode
A = X values
B = Y values
mode = None (default) : for regular linear regression
= "norm" : use normalised values of A and B
= "std" : use standardised values for A and B
Output: [slope, intercept, pearson]
slope = angular coefficient of the linear regression
intercept = linear coefficient of the linear regression
pearson = Pearson correlation coefficient
Error:
"Arrays must have the same size" : A and B have different sizes.
"Mode not recognized. Use None, 'norm' or, 'std'" : invalid value for 'mode'
- **Function:** find_critical(Y, X, min_points=5, use_inflex=False)
Takes the X and Y values of a function and returns the critical points that match the criteria.
Input: Y, X, min_points, use_inflex
Y = ordinate values of a function
X = abscissa values of a function
min_point = 5 (default): minimum amount of points between two critical points
use_inflex = False (default) : Not search for inflexion points (second derivative = 0)
= True : Search for inflexion points (second derivative = 0)
Output: critical_point_list
critical_point_list = Array containing the index of critical points (zero indexed)
Error:

"Arrays must have the same size" : X and Y have different sizes.
 "Invalid value. Use True or False" : invalid value for use_inflex
 "Too many points between two critical points" : min_points must be lower than the X array length

- **Function:** split_segm(A, critical_point_list)

Takes the A array and divide it into N arrays according to the number of critical points. Each array corresponds to a segment of the REG analysis.

Input: A, critical_point_list

A = Array containing the abscissa values of a function
 critical_points_list = list of the index of the critical points (zero indexed)

Output: segm

segm = Array of arrays containing the values of A for each segment.

Error:

"Invalid critical point index" : index of critical point out of range of array.

- **Function:** reg(wfn_energy, control_coord, terms, critical=True, np=5, inflex=True, critical_index= [], mode="norm")

Perform the REG analysis over all contributions inside "terms" array

Input: wfn_energy, control_coord, terms, critical, np, inflex, critical_index, mode

wfn_energy = Array containing the energy values for each point (PES - Energy)
 control_coord = Array containing the control coordinate values for each point (PES - Coordinate)

terms = Array of arrays, each one corresponding to one IQA contribution.

critical = True(default): Search for critical points

= False : user must provide the critical point index array

np = 5(default): : minimum amount of points between two critical points (used if critical == True)

use_inflex = False (default) : Not search for inflexion points (second derivative = 0) (used if critical == True)

= True : Search for inflexion points (second derivative = 0 (used if critical == True)

critical_index = [](default) : list of critical points provided by the user. (used if critical == False)

mode = None : for regular linear regression

= "norm"(default) : use normalised values of A and B

= "std" : use standardised values for A and B

Output: [REG_values, pearson_values][segment]

REG_values = array of REG coefficients for each contribution. Each array correspond to a segment

pearson_values = array of Pearson coefficients for each contribution. Each array correspond to a segment

Error:

"PES abscissa and ordenate must have same number of points" : wfn_energy and control_coord have different sizes

"Contributions arrays must have same size" : arrays inside terms have different sizes.

- **Function:** integration_error(wfn_energies, IQA_energies)

Calculates integration error for each PES point

Input: wfn_energies, IQA_energies

`wfn_energies` = list of total energies from wfn files
`IQA_energies` = list of total energy obtained from IQA terms
Output: [error, RMSE]
`error` = list of `wfn_energies` - `IQA_energies`
`RMSE` = Root mean square error
Error:
 "Arrays must have the same size" : `wfn_energies` and `IQA_energies` have different sizes.

- Function:** `group_intra(intra_energies, intra_energies_header, groups)`
 Group IQA intra-atomic terms into the user-defined groups
Input: `intra_energies`, `intra_energies_header`, `groups`
`intra_energies` = list of IQA intra-atomic terms
`intra_energies_header` = header of IQA intra atomic terms.
`groups` = list of list of atomic labels, each internal list correspond to a different group
Output: [new_terms, new_header]
`new_terms` = list of the grouped IQA terms values
`new_header` = list of the grouped terms header
- Function:** `group_inter(inter_energies, inter_energies_header, groups)`
 Group IQA interatomic terms into the user-defined groups
Input: `inter_energies`, `inter_energies_header`, `groups`
`inter_energies` = list of IQA interatomic terms
`inter_energies_header` = header of IQA interatomic terms.
`groups` = list of list of atomic labels, each internal list correspond to a different group
Output: [new_terms, new_header]
`new_terms` = list of the grouped IQA terms values
`new_header` = list of the grouped terms header

4.2 List of Functions - `aimall_utils.py`

- Function:** `get_atom_list(wfn_file)`
 Get atomic labels from wfn file
Input: `wfn_file`
`wfn_file` = Any wfn file of the desired PES
Output: [atom_list]
`atom_list` = list of each atom label for all atoms in molecule
Error:
 "Atomic labels not found" : Atom list does not exist in `wfn_file`
- Function:** `get_aimall_wfn_energies(A)`
 Get all wfn energies from the wavefunction files (`wfn`)
Input: [A]
 [A] = list of all wfn files of the PES
Output: [wfn_energy]
 [wfn_energy] = list of energies for each PES point
Error:
 "Energy values not found in file: file_name" : No energy values in file_name

- Function:** `intra_property_from_int_file(folders, prop, atom_list)`
 Get IQA intratomic properties from int files
Input: `folders, prop, atom_list`
`folders` = path to *.atomicfiles folders
`prop` = list of IQA for each atoms e.g.: `['T(A)', 'Vee(A,A)', 'Vne(A,A)']`
`atom_list` = list of atomic labels e.g.: `[n1, c2, h3, ...]`
Output: `[intra_properties, contributions_list]`
`intra_properties` = array of array containing the IQA values for each atom for each geometry
`contributions_list` = list of contributions organised in the same order as in `intra_properties`
- Function:** `inter_property_from_int_file(folders, prop, atom_list)`
 Get IQA interatomic properties from int files
Input: `folders, prop, atom_list`
`folders` = path to *.atomicfiles folders
`prop` = list of IQA for each atoms e.g.: `['T(A)', 'Vee(A,A)', 'Vne(A,A)']`
`atom_list` = list of atomic labels e.g.: `[n1, c2, h3, ...]`
Output: `[intra_properties, contributions_list]`
`intra_properties` = array of array containing the IQA values for each atom for each geometry
`contributions_list` = list of contributions organised in the same order as in `intra_properties`

4.3 List of Functions - `gaussian_utils.py`

- Function:** `get_atom_list_wfn_g09(wfn_file)`
 Get atomic labels from g09 wfn file
Input: `wfn_file`
`wfn_file` = Any wfn file of the desired PES
Output: `atom_list`
`atom_list` = list of each atom label for all atoms in molecule
Error:
`"Atomic labels not found"` : Atom list does not exist in `wfn_file`
- Function:** `get_atom_list_wfx_g09(wfx_file)`
 Get atomic labels from g09 wfx file
Input: `wfx_file`
`wfx_file` = Any double wavefunction file of the desired PES
Output: `atom_list`
`atom_list` = list of each atom label for all atoms in molecule
Error:
`"Atomic labels not found"` : Atom list does not exist in `wfn_file`
- Function:** `get_atom_list_wfn_g16(wfn_file)`
 Get atomic labels from g16 wfn file
Input: `wfn_file`
`wfn_file` = Any wfn file of the desired PES
Output: `atom_list`
`atom_list` = list of each atom label for all atoms in molecule

Error:

"Atomic labels not found" : Atom list does not exist in wfn_file

4.4 List of Functions - morphy_utils.py

- **Function:** read_morphy_outputs(file_list)
Get correlation data from Morphy output
Input: [file_list]
[file_list] = array that contains the path for the outputs of Morphy
Output [interactions_temp, final_values]
interactions_temp = list of all interatomic correlation energies headers
final_values = list of all interatomic correlation energies values

4.5 List of Functions - reg_vis.py

- **Function:** plot_violin(df_list, save = False, file_name='reg_violin.png')
Creates a violin plot for each segment.
Input: df_list, save, file_name
df_list = list of dataframes containing REG and R for each segment.
save = False (default) : not save figure to file.
= True : save figure to file named reg_violin.png
file_name = ./reg_violin.png (default) : output file path if save = True
Output: figure
figure : violin plot
- **Function:** generate_data_vis(df, df_list, n_term, save = False, file_name='data_vis.png', title = "REG data visualization")
Creates a graphical summary of the REG analysis
Input: df, df_list, n_term, save, file_name
df = segment dataframe [containing REG and R^2 values]
df_list = list of dataframes for all segments [containing REG and R^2 values]
n_term = number of relevant terms to be highlighted
save = False (default) : not save figure to file.
= True : save figure to file named data_vis.png
file_name = ./data_vis.png :output file path if save = True
title = "REG data visualization" (default) : output file path if save = True
Output: figure
figure : data visualization summary figure
- **Function:** plot_segment(coordinate, wfn_energy, critical_points, label = False, color=True, annotate=True, title='REG segments', y_label='Energy', x_label='Coordinate', save = False, file_name='segments.png')
Plot the PS showing the segment division
Input: coordinate, wfn_energy, critical_points, label, color, annotate, title, y_label, x_label, save, file_name

coordinate	=	control coordinate array
wfn_energy	=	array containing the energies from the wfn/wfx files
critical_points	=	array containing the list of critical points
label	=	False (default) : not label each point
	=	True : label each point
color	=	True (default) : color each segment
	=	False : not color each segment
annotate	=	True(default) : label each segment
	=	False : not label each segment
save	=	False (default) : not save figure to file.
	=	True : save figure to file named ./segments.png
title	=	"REG segments" (default) : output file path if save = True
file_name	=	./segments.png :output file path if save = True
y_label	=	y-axis label
x_label	=	x-axis label

Output: figure

figure : data visualization summary figure

4.6 Data Visualisation Examples

1. Violin Plot:

The violin plots are used to compare the R^2 distribution for each segment. If the division of segments is done correctly, one can expect that the distribution of R^2 values is similar in each segment. However, a bad definition of segments will causes discrepant distributions, which can easily be seen in a violin plot. Figure 4 presents examples of violin plots, showing a bad and a good segmentation of the PES.

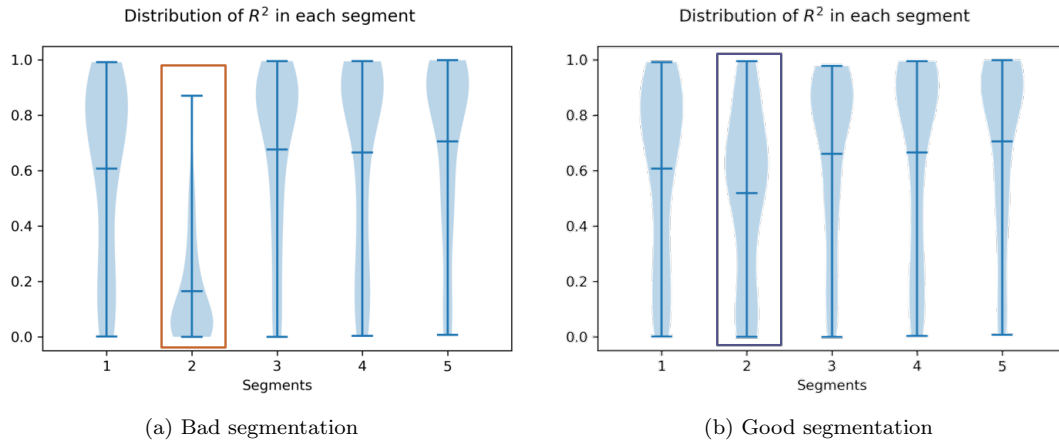


Figure 4: Example of a bad segmentation (a) and a good segmentation(b). Note that for segment 2, in (a), the vast majority of the IQA presents low values of R^2 .

2. Summary of Results:

Figure 5 contains a summary of the REG analysis for one segment. In (a): the violin plot comparing the R^2 distribution of segment 2 (green) with the R^2 distribution of all other segments. (b): histogram showing the distribution of R^2 in segment 2. (c): REG coefficient lines, showing the m_{REG} plotted against R^2 . (d): The most relevant IQA terms according to the REG coefficient value.

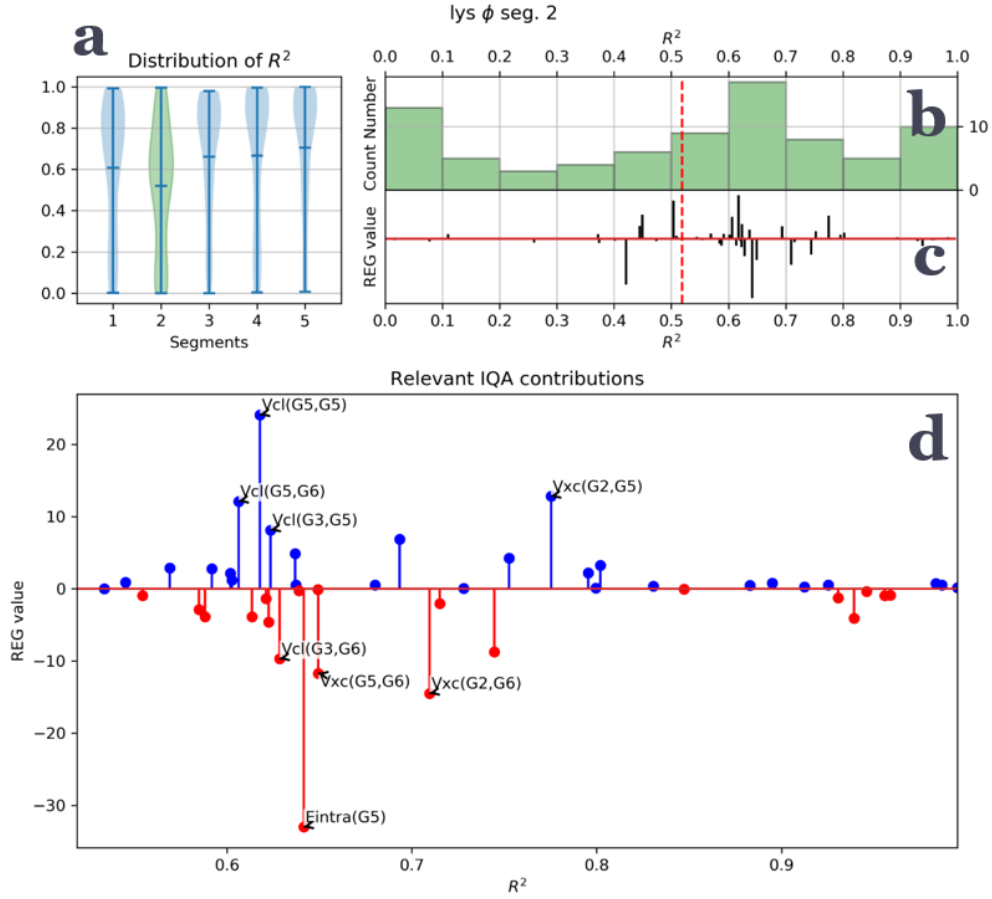


Figure 5: Data visualization summary generated by *generate_data_vis* function.

4.7 Workflow

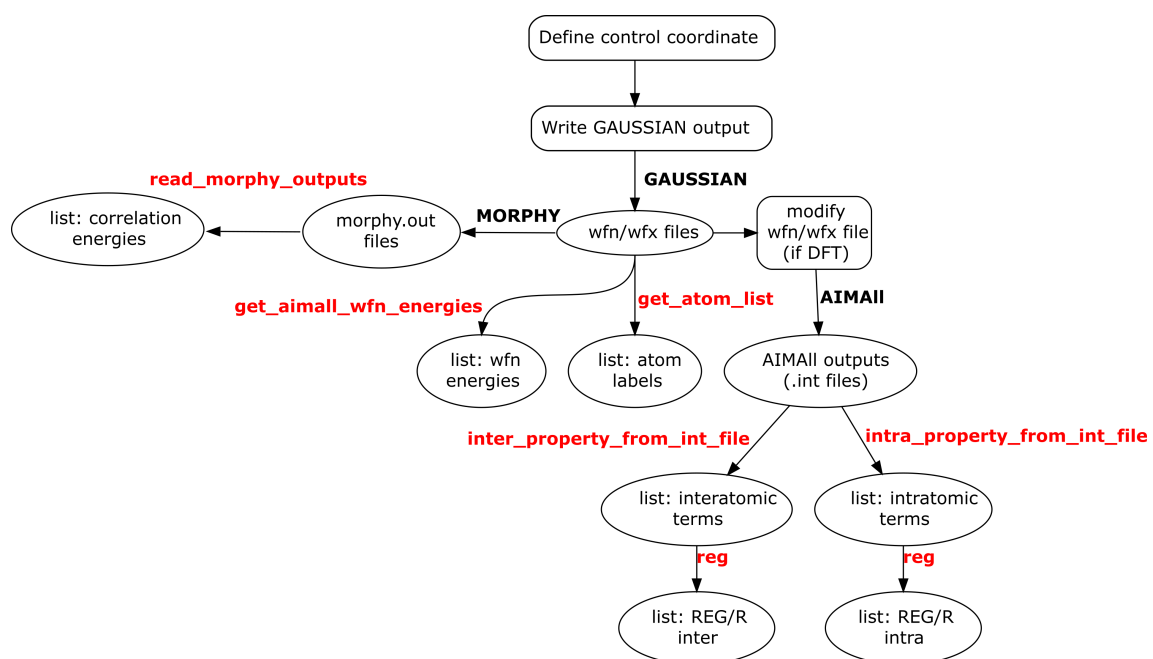


Figure 6: Basic workflow chart of how to perform a REG analysis. In red, functions available in REG.py

4.8 Script Example

In this section, a script is presented as a guide to help the user to create its own. Download the full example at [github](#).

```
# IMPORT MODULES
import reg
import aimall_utils
import numpy as np
import pandas as pd
import reg_vis as rv
##### OPTIONS #####
SYS = 'bbr3'
EQ_DIST = 0
POINTS = 4
INFLEX = False
WRITE = True
AUTO = True
tp = [7, 15]
SAVE_FIG = True
ANNOTATE = True
DETAILED_ANALYSIS = True
LABELS = True
n_terms = 4
##### SETUP FILES #####
#DEFINE PATHS AND FILES:
path = "../final/b3lyp_nh3/" + SYS + '/'
A = [i for i in range(0, 40)]
#A.pop(15)
#DEFINE THE CONTROL COORDINATES:
cc = [1+i*0.1 for i in A]
wfn_files = [path + str(format(i, '02d')) + '/' + str(format(i, '02d')) + '.wfn' for i in A]
folders = [path + str(format(i, '02d')) + '/' + str(format(i, '02d')) + '_atomicfiles' for i in A]

#####
# Note: In this case we decided to group some atoms in specific groups.
#####
groups = [['b1'], #G1
          ['br2', 'br3', 'br4'], #G2
          ['n5'], #G3
          ['h6', 'h7', 'h8']] #G4
#GET ATOM LIST FROM ANY .WFN FILE:
atoms = aimall_utils.get_atom_list(wfn_files[0])
#DEFINE THE DESIRED TERMS:
intra_prop = ['E_IQA_Intra(A)'] #intra atomic properties
inter_prop = ['VC_IQA(A,B)', 'VX_IQA(A,B)'] #inter atomic properties
#GET TOTAL ENERGY FROM THE .WFN FILES:
total_energy_wfn = aimall_utils.get_aimall_wfn_energies(wfn_files)
#GET INTRA-ATOMIC TERMS:
iqa_intra_temp = aimall_utils.intra_property_from_int_file(folders,
                                                           intra_prop, atoms)[0]
iqa_intra_header_temp = np.array(aimall_utils.intra_property_from_int_file(
    folders, intra_prop, atoms)[1]) #used for reference
#GET INTER-ATOMIC TERMS:
iqa_inter_temp = aimall_utils.inter_property_from_int_file(folders,
                                                           inter_prop, atoms)[0]
iqa_inter_header_temp = np.array(aimall_utils.inter_property_from_int_file(
    folders, inter_prop, atoms)[1]) #used for reference
#CONVERT LISTS TO ARRAYS.
total_energy_wfn = np.array(total_energy_wfn)
iqa_inter_temp = np.array(iqa_inter_temp)
```



```

iqa_intra_temp = np.array(iqa_intra_temp)
#CALCULATED THE IQA TOTAL ENERGY
total_energy_iqa = sum(iqa_inter_temp) + sum(iqa_intra_temp) # used to
                                                           calculate the integration error.

#GROUP INTRA ATOMIC TERMS:
iqa_intra = np.array(reg.group_intra(iqa_intra_temp, iqa_intra_header_temp,
                                     groups)[0])
iqa_intra_header = np.array(reg.group_intra(iqa_intra_temp,
                                             iqa_intra_header_temp, groups)[1])

#GROUP INTER ATOMIC TERMS:
iqa_inter = np.array(reg.group_inter(iqa_inter_temp, iqa_inter_header_temp,
                                     groups)[0])
iqa_inter_header = np.array(reg.group_inter(iqa_inter_temp,
                                             iqa_inter_header_temp, groups)[1])

##### REG ANALYSIS #####
#INTRA ATOMIC CONTRIBUTION
reg_intra = reg.reg(total_energy_wfn, cc, iqa_intra, np=POINTS, critical=
                    AUTO, inflex=INFLEX, critical_index=
                    tp)

#INTER ATOMIC CONTRIBUTION
reg_inter = reg.reg(total_energy_wfn, cc, iqa_inter, np=POINTS, critical=
                    AUTO, inflex=INFLEX, critical_index=
                    tp)

#CALCULATE INTEGRATION ERROR:
rmse_integration = reg.integration_error(total_energy_wfn, total_energy_iqa
)
print('Integration error [kJ/mol](RMSE)')
print(rmse_integration[1])

```

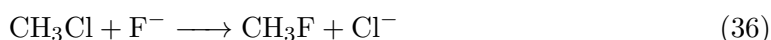
5 Appendix A

5.1 Defining the Control Coordinate - An S_N2 reaction

The control coordinate for a REG analysis is defined by a curve connecting two points on the potential energy surface of a system. Such a line can correspond to a chemical reaction, a vibrational movement, a conformation change, a phase change, a dissociation process, etc. The definition of the control coordinate will be different for each application.

As an example, we define a control coordinate using the intrinsic reaction coordinate (IRC) of an S_N2 reaction using GAUSSIAN09. To do that, we first need to find the transition state of the reaction.

Finding transition states (TS) can be a tricky process since different systems require different approaches. In this section we describe in detail all steps used to find the TS of gas-phase S_N2 reactions without using the Synchronous Transit-Guided Quasi-Newton Methods (QST2, QST3). The examples given are based on the following reaction:



Input files are reproduced when necessary.

5.1.1 STEP 1 - Optimise reagents geometries

In chemical reactions, the transition state is a point in the reaction path that can only be accessed when the reagents collide at a specific angle with a determined minimum amount of energy. In this way, optimizing the reactants' geometry is a good starting point. After the optimisation step, an IR frequency calculation is performed to attest for a minimum energy structure. If all frequencies are real numbers we reached a minimum point (global or local) on the potential energy surface. Each molecule is optimised separately. The optimisation input is reproduced below:

```
%nprocshared=12
%mem=16GB
#opt mp2/aug-cc-pvtz density=current nosymm
```

CH3Cl Optimization

O 1			
C	0.000 003 00	-0.000 002 00	-1.122 804 00
H	-0.000 005 00	1.028 877 00	-1.464 294 00
H	-0.891 036 00	-0.514 435 00	-1.464 294 00
H	0.891 038 00	-0.514 441 00	-1.464 289 00
Cl	0.000 000 00	0.000 000 00	0.657 605 00

In this example, the keyword **density=current** is used to request the density matrix of the current method (MP2).

5.1.2 STEP 2 - Scan through the Potential Energy Surface (PES)

Once we found the reagents optimized geometries, we can begin to explore the PES. To do that, we need to use the **Scan** method of GAUSSIAN09. To access this method we add **#opt=ModRedundant** to the input file and specify what coordinate we want to scan over.

In this step we bring together both reactants geometries, building an initial guess for the TS geometry. Then, we move the nucleophile away from the electrophile (or *vice versa*). At each point of the scan coordinate we get an energy value. A maximum point

within this coordinate will give an approximation of the TS geometry. We'll use this approximation in **Step 3**. The GAUSSIAN scan input is reproduced below:

```
%nprocshared=12
%mem=16GB
#opt=(Modredundant) mp2/aug-cc-pvtz density=current

CH3Cl + F- scan

-1 1
C      0.000 074 00      -0.000 043 00      -0.645 149 00
H      -0.000 061 00       1.038 991 00      -0.987 064 00
H      -0.899 823 00      -0.519 443 00      -0.987 064 00
H       0.899 832 00      -0.519 518 00      -0.987 069 00
F      -0.000 022 00       0.000 012 00       0.808 569 00
Cl     -0.000 007 84       0.000 004 56      -2.328 519 00

B 1 5 S 20 0.05
A 1 5 6 F
```

The last two lines of the input file correspond to the **Modredundant** parameters. The first parameter states that the coordinate along the **Bond** between atoms **1** and **5** will be **Scanned** in **20** steps with a displacement of **0.05 Å** each. The second parameter states that the **Angle** defined by atoms **1**, **5** and **6** are **Freeze**, that is, it may not be altered during the scanning process. If the nucleophile has more than one atom, more constraints will be necessary to perform the scanning.

Also, it is important to be aware of the correct basis set. Figure 7 compares the results of a scan calculation using the same geometry and constraints, changing the basis-set only. One can see that using **aug-cc-pVTZ** leads to a maximum point on the PES while the same does not occurs when using **6-31G** basis set.

5.1.3 STEP 3 - Releasing the constraints I

The necessary information to make the first guess for finding the TS is obtained from the procedure described in **STEP 2**. Figure 8A shows the energy against the scan coordinate, note that there is a maximum point, which will be the initial guess geometry for the TS. Figure 8B shows the gradient of the former curve and at the maximum point the gradient is equal to zero.

The geometry at the maximum point needs to be optimised without releasing all the constraints. To do that, we keep the angular constraint between nuclei 1(*C*), 5(*F*) and 6(*Cl*). The GAUSSIAN input is reproduced:

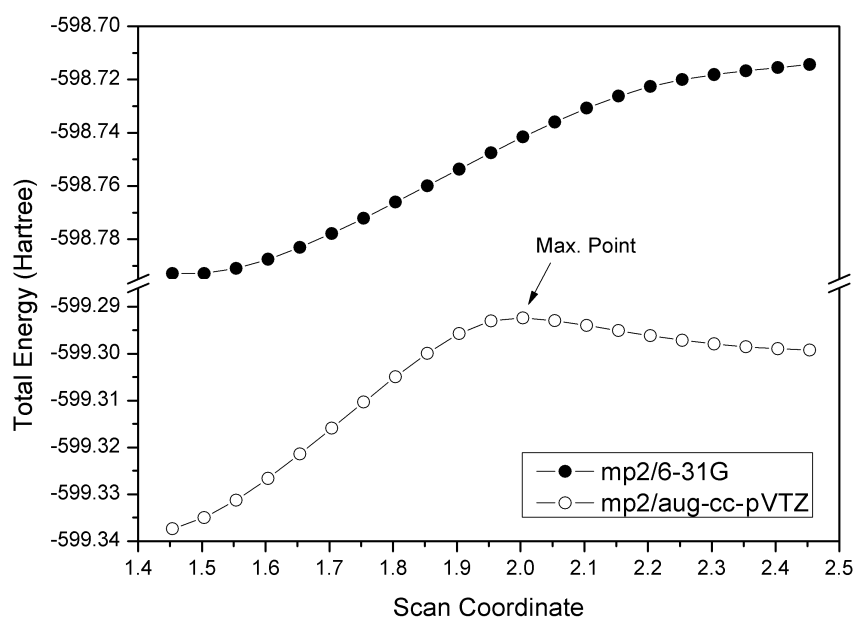


Figure 7: Total Energy versus the scan coordinate for the approximation of the fluoride ion and CH_3Cl . The black dots were obtained using the 6-31G basis set. The white dots were obtained using the aug-cc-pVTZ basis set, which leads to a maximum point.

```
%nprocshared=12
%mem=16GB
#opt=(TS,Modredundant,calcall,noeigentest) mp2 aug-cc-pvtz density=current

Constrained transition state optimisation

-1 1
C      0.495 553 00      0.000 114 00      -0.000 104 00
H      0.610 015 00     -0.492 687 00      0.942 588 00
H      0.610 001 00     -0.569 918 00     -0.898 210 00
H      0.609 934 00      1.062 923 00     -0.044 693 00
F      2.499 271 00     -0.000 056 00      0.000 054 00
Cl     -1.605 689 00     -0.000 029 00      0.000 027 00

A 1 5 6 F
```

The **TS** keyword means that the Bernaly algorithm will be used to perform the optimisation, searching for a maximum point instead of searching for a local minimum. **Calcall** means that the Hessian matrix will be recalculated in each optimization step. **Noeigentest** prevents that the occurrence of more than one negative eigenvalue causes an error.

In this step, if no constraints are imposed, the optimisation algorithm can lead to another reaction TS. In this example, an input without constraints will find the TS of methyl chloride deprotonation (Figure 9).

5.1.4 STEP 4 - Releasing the constraints II

Once the best guess for the TS was found, a final optimisation without constraints is made. This process should result in the correct TS geometry. The GAUSSIAN09 input

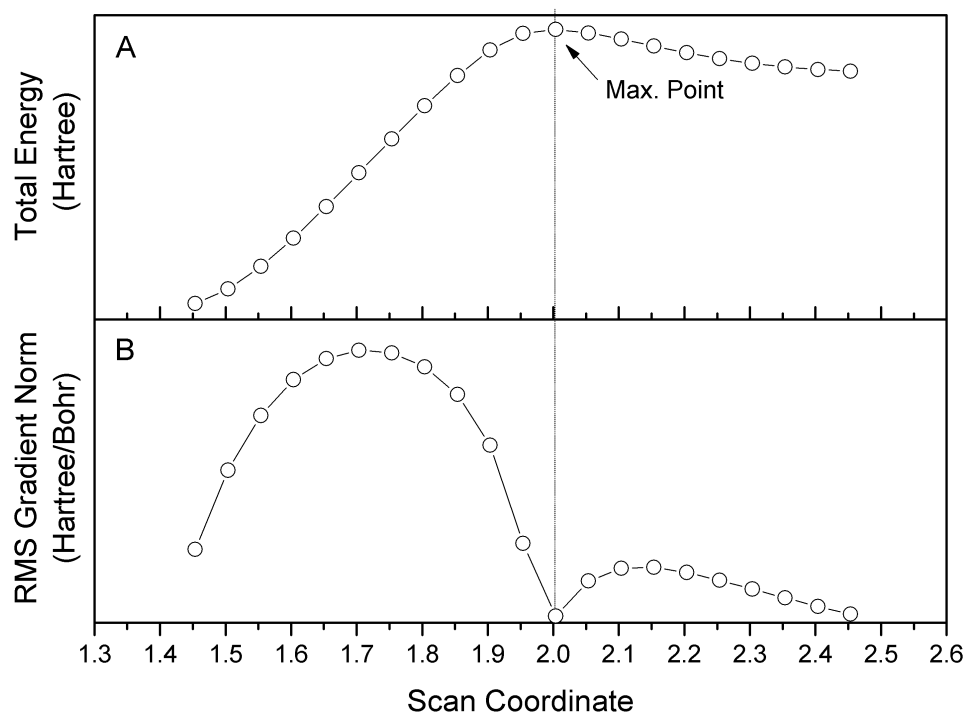


Figure 8: In A: The maximum energy point, where one can find a good approximation for the transition state geometry. In B: The first derivative norm of A.

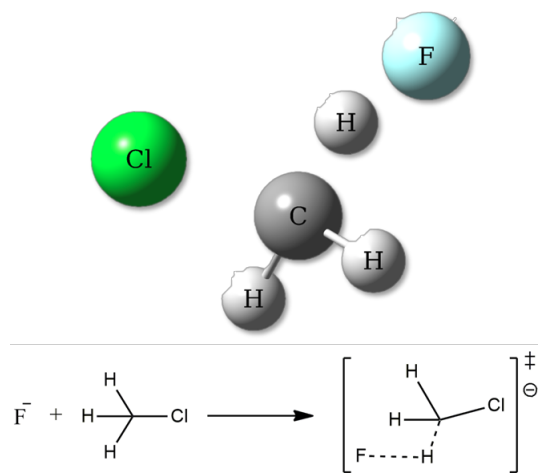


Figure 9: Transition state of methyl chloride deprotonation by fluoride ion.

for this step is:

```
%nprocshared=12
%mem=16GB
#opt=(TS,calcall,noeigentest) mp2 aug-cc-pvtz density=current
```

Final transition state optimization

-1 1			
C	0.504 296 00	0.000 114 00	-0.000 103 00
H	0.609 881 00	-0.505 729 00	0.936 557 00
H	0.609 855 00	-0.558 179 00	-0.906 487 00
H	0.609 836 00	1.064 214 00	-0.030 387 00
F	2.501 628 00	-0.000 057 00	0.000 054 00
Cl	-1.610 000 00	-0.000 028 00	0.000 027 00

Finally, it is necessary to check for the infrared frequencies. Since the TS represents a saddle point on the PES, i.e a maximum in one direction and a minimum in all others, it should have one, *and only one*, imaginary frequency. If more than one imaginary frequency is presented, it is necessary to repeat **Step 3** with different constraints. Figure 10 summarizes the results.

5.1.5 STEP 5 - Finding the IRC

Using the TS geometry as input, we request GAUSSIAN to find the intrinsic reaction coordinate in both **forward** and **reverse** directions.

```
%nprocshared=12
%mem=16GB
#irc=(reverse,maxpoints=15,calcall) mp2/aug-cc-pvtz
```

IRC

-1 1			
C	0.504 296 00	0.000 114 00	-0.000 103 00
H	0.609 881 00	-0.505 729 00	0.936 557 00
H	0.609 855 00	-0.558 179 00	-0.906 487 00
H	0.609 836 00	1.064 214 00	-0.030 387 00
F	2.501 628 00	-0.000 057 00	0.000 054 00
Cl	-1.610 000 00	-0.000 028 00	0.000 027 00

The output will contain 15 geometries in the forward direction and 15 geometries in the reverse direction. These geometries define the control coordinate for the REG analysis. For each one of these geometries, it is necessary to request GAUSSIAN to perform a single point calculation writing the wavefunction (.wfn or .wfx) file.

AIMALL will take the wavefunction files as input to obtain the IQA properties that will be gathered by REG.py.

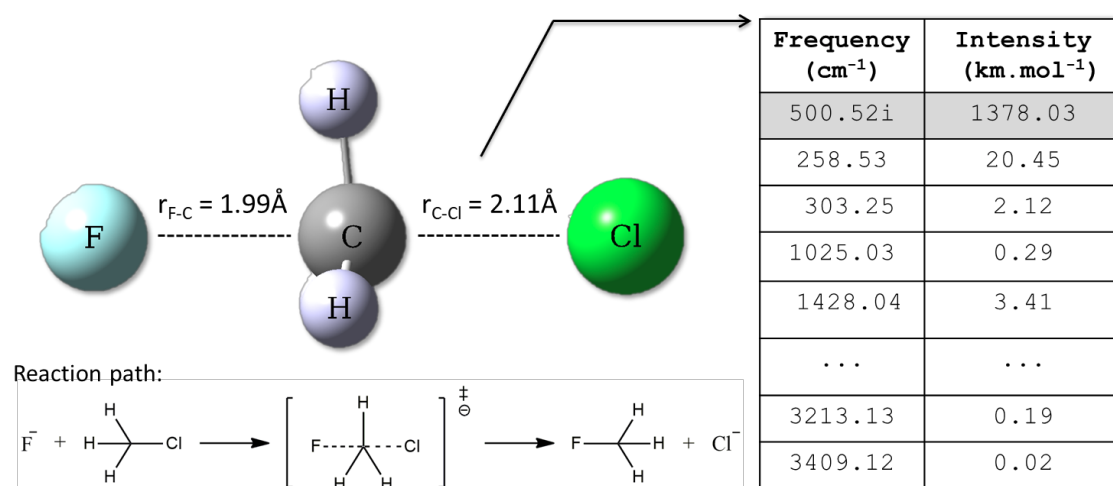


Figure 10: Transition state obtained after the proposed steps. Note that only one imaginary frequency exists, which means that we have found a saddle point in the PES.

6 Appendix B

6.1 Demonstrating $C = 0$

Consider the plot in figure 11 where the ordinate axis corresponds to a IQA term and the abscissa axis corresponds to the total energy of the system. The dashed line corresponds to the regression line obtained by equation 28.

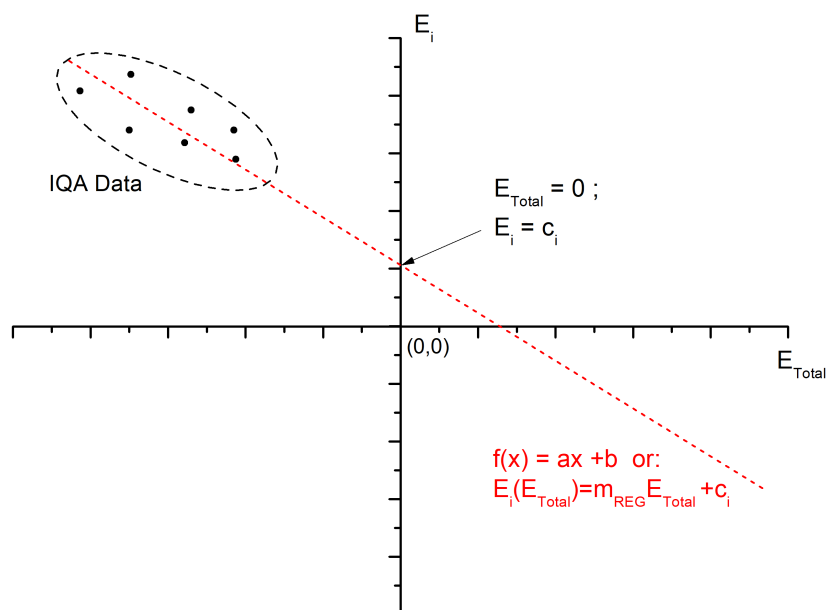


Figure 11: Example of a linear regression obtained by equation 28

Although the total energy of the system is never equal to zero, the regression line is a mathematical function and its domain corresponds to all real numbers, i.e. $(-\infty, +\infty)$.

The regression line is given by:

$$E_i(E_{total}) = m_{REG,i}E_{total} + c_i$$

extrapolating the curve to the point where $E_{total} = 0$:

$$E_i(E_{total} = 0) = m_{REG,i}0 + c_i$$

$$E_i(0) = c_i$$

consider now the sum of all regression lines, one for each IQA term, at the point where $E_{total} = 0$

$$\sum E_i(0) = \sum c_i$$

since $\sum E_i = E_{total}$, and $E_{total} = 0$ then:

$$E_{total} = \sum c_i$$

$$\sum c_i = 0 \implies C = 0; Q.E.D.$$

Since C is constant, $C = 0$ is true for the entire domain $(-\infty, +\infty)$.