



Atlantisc  
2 · 0 · 2 · 0

Recherche,  
Formation  
& Innovation  
en PAYS de la LOIRE



# QuChemPedia

Quantum Chemistry Collaborative EncyclopediA

Formula ↕

Thomas CAUCHY  
Benoit DA MOTA

<thomas.cauchy@univ-angers.fr>  
<benoit.damota@univ-angers.fr>

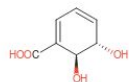
Intelligence artificielle pour la chimie des matériaux - 25 Sept. 2018

10 results for 1.



&lt; Page : 1 / 1 &gt;

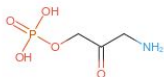
Results per page : 10



- OPT  
 FREQ  
 TD  
 SP  
 OPT\_ES  
 FREQ\_ES

**Formula :** C<sub>7</sub>H<sub>6</sub>O<sub>4</sub>  
**Charge :** 0  
**Multiplicity:** 1  
**Solvent :** gas / **Solvation method :** Unknown

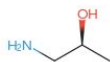
**IUPAC :** Unknown  
 GAMESS  
 B3LYPV1R / **Basis set :** 6-31G\* (181 functions)  
**Ending energy :** -572.4028395 a.u.



- OPT  
 FREQ  
 TD  
 SP  
 OPT\_ES  
 FREQ\_ES

**Formula :** C<sub>3</sub>H<sub>8</sub>NO<sub>5</sub>P  
**Charge :** 0  
**Multiplicity:** 1  
**Solvent :** gas / **Solvation method :** Unknown

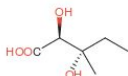
**IUPAC :** Unknown  
 GAMESS  
 B3LYPV1R / **Basis set :** 6-31G\* (170 functions)  
**Ending energy :** -891.4226034 a.u.



- OPT  
 FREQ  
 TD  
 SP  
 OPT\_ES  
 FREQ\_ES

**Formula :** C<sub>3</sub>H<sub>9</sub>NO  
**Charge :** 0  
**Multiplicity:** 1  
**Solvent :** gas / **Solvation method :** Unknown

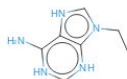
**IUPAC :** Unknown  
 GAMESS  
 B3LYPV1R / **Basis set :** 6-31G\* (93 functions)  
**Ending energy :** -249.6945983 a.u.



- OPT  
 FREQ  
 TD  
 SP  
 OPT\_ES  
 FREQ\_ES

**Formula :** C<sub>6</sub>H<sub>12</sub>O<sub>4</sub>  
**Charge :** 0  
**Multiplicity:** 1  
**Solvent :** gas / **Solvation method :** Unknown

**IUPAC :** Unknown  
 GAMESS  
 B3LYPV1R / **Basis set :** 6-31G\* (174 functions)  
**Ending energy :** -536.7636026 a.u.

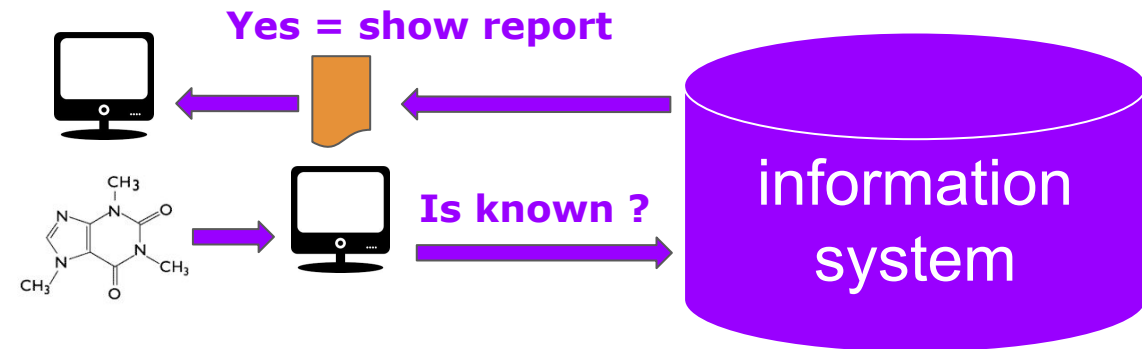


- OPT  
 FREQ  
 TD  
 SP  
 OPT\_ES  
 FREQ\_ES

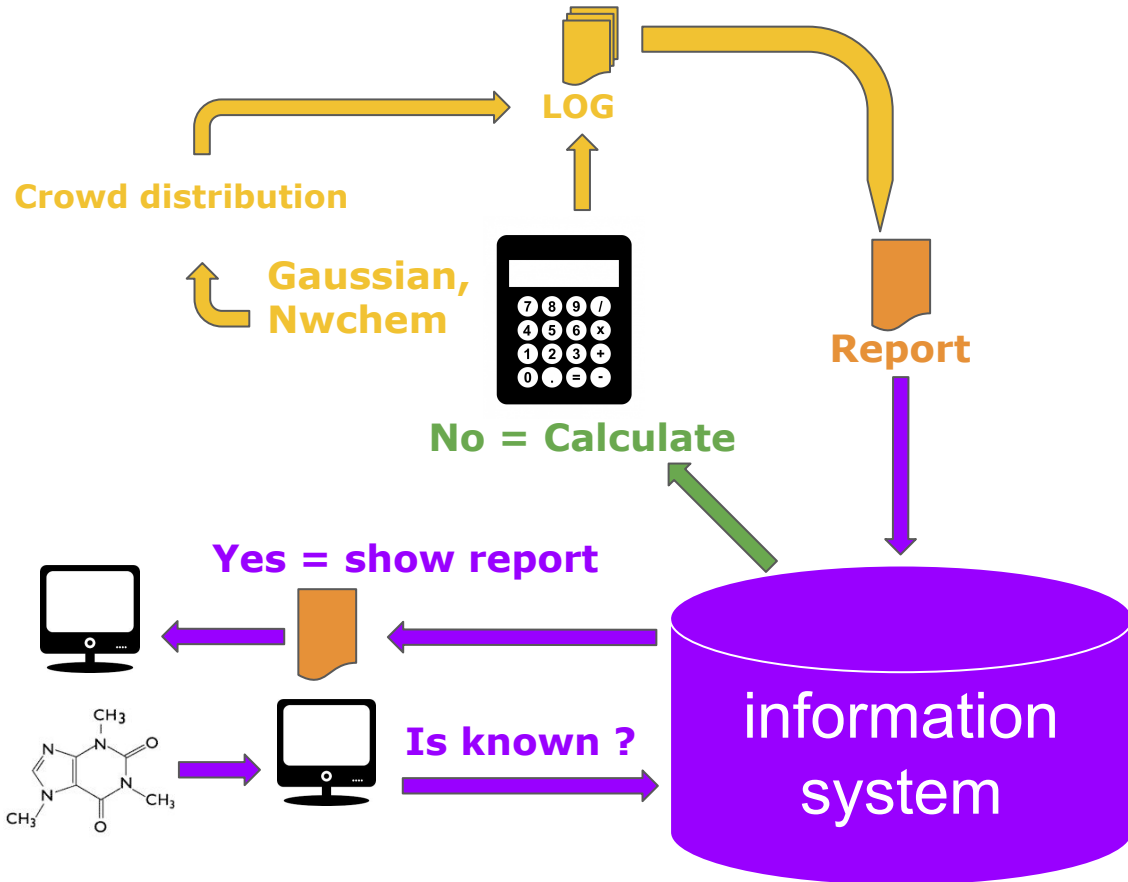
**Formula :** C<sub>7</sub>H<sub>9</sub>N<sub>5</sub>  
**Charge :** 0  
**Multiplicity:** 1  
**Solvent :** gas / **Solvation method :** Unknown

**IUPAC :** Unknown  
 GAMESS  
 B3LYPV1R / **Basis set :** 6-31G\* (198 functions)  
**Ending energy :** -545.9492504 a.u.

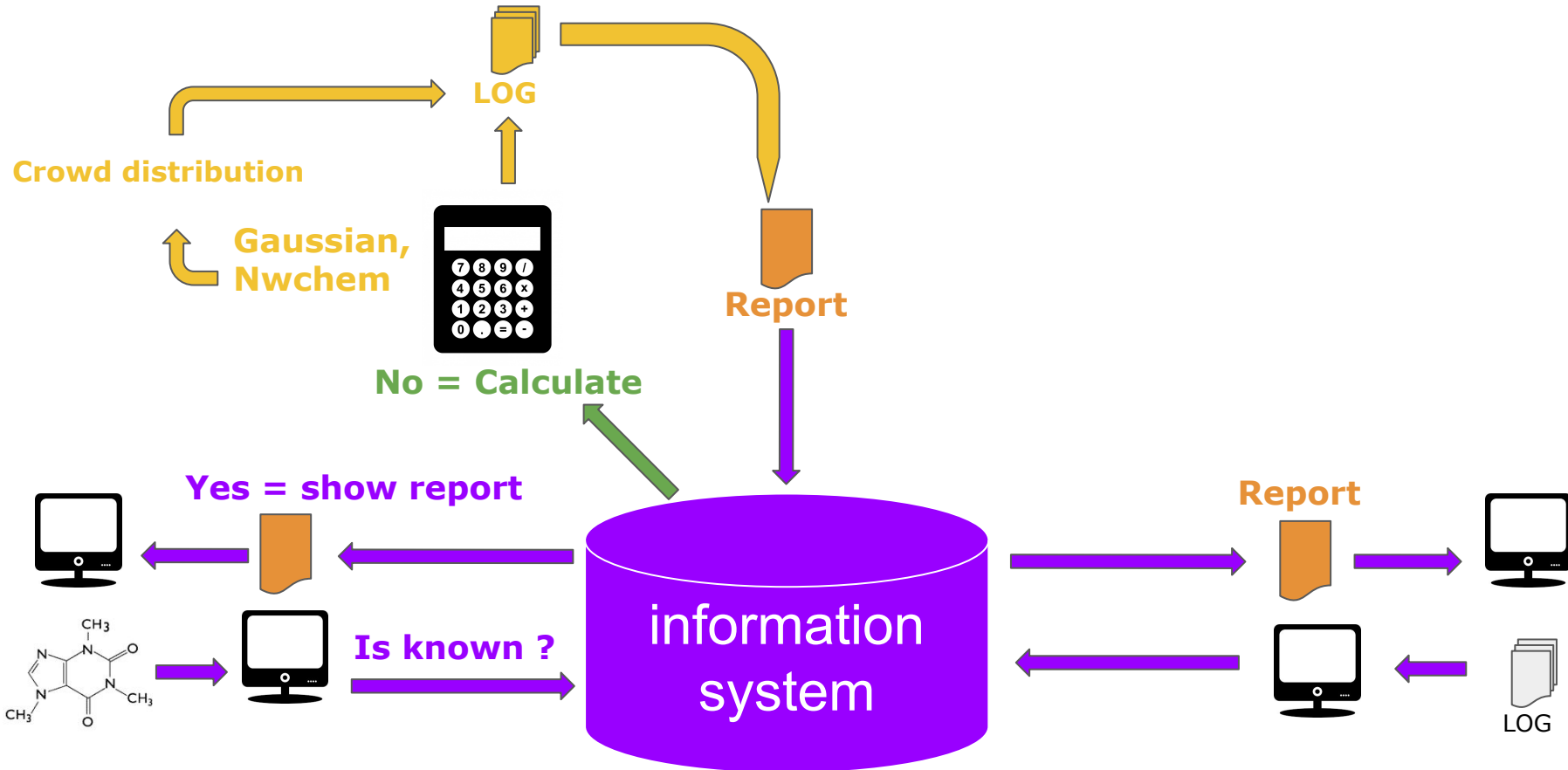
# The QuChemPedIA project



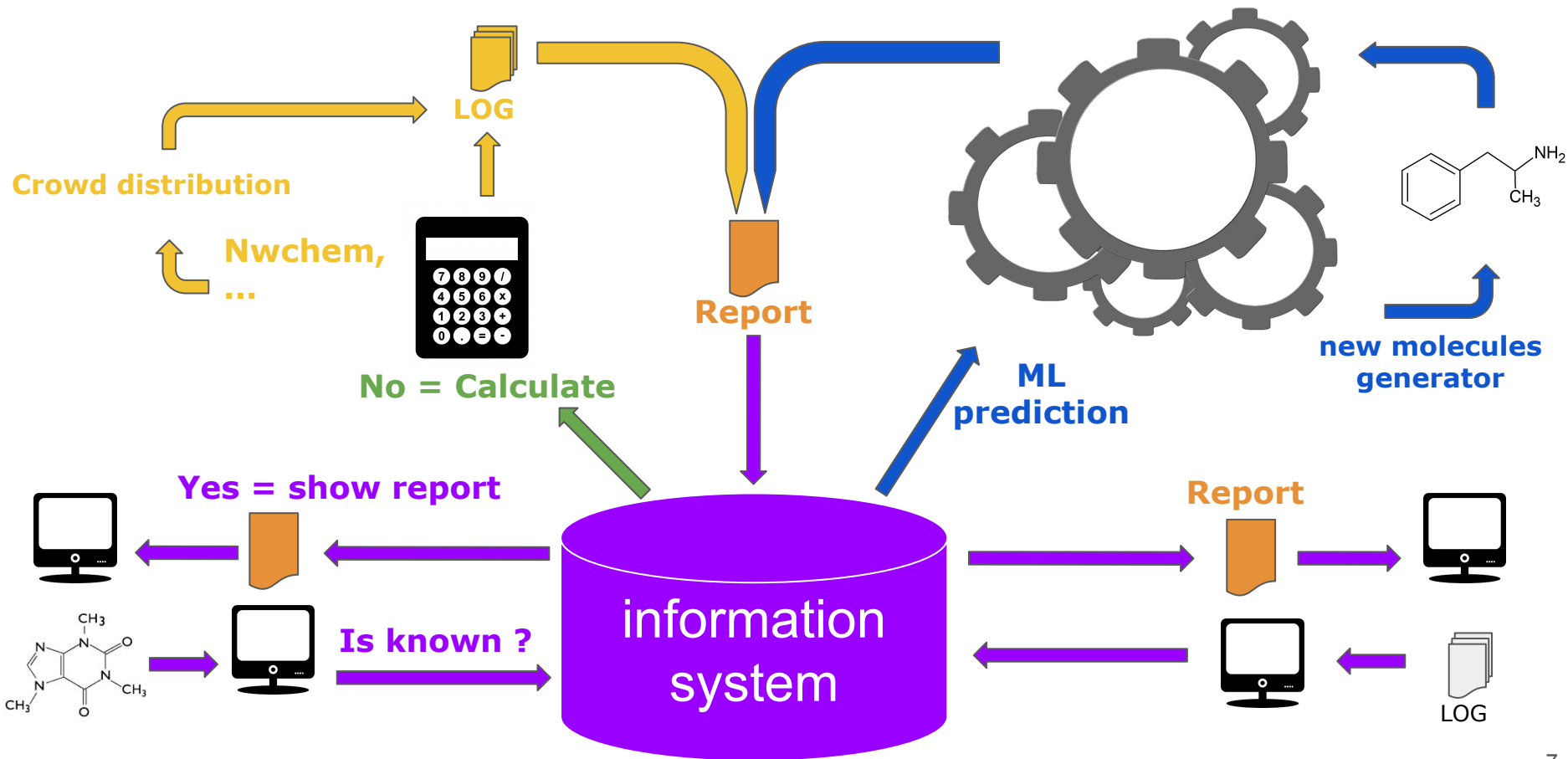
# The QuChemPedIA project



# The QuChemPedIA project



# The QuChemPedIA project





## Molecule

[Associated calculations](#)

[Authorship](#)

[Computational details](#)

[Results](#)

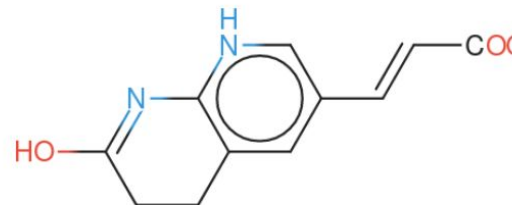
> [Geometry](#)

> [Thermochemistry](#)

> [Excited states](#)

## Molecule

<b>InChI</b> <a href="#">?</a>	1S/C11H10N2O3 /c14-9-3-2-8-5-7(1-4-10(15)16)6-12-11(8)13-9 /h1,4-6H,2-3H2,(H,15,16)(H,12,13,14)/b4-1+
<b>Canonical SMILES</b> <a href="#">?</a>	OC(=O)C=C /c1cnc2c(c1)CCC(=N2)O
<b>Monoisotopic mass</b>	218.06914219
<b>Formula</b>	C <sub>11</sub> H <sub>10</sub> N <sub>2</sub> O <sub>3</sub>
<b>Charge</b>	0
<b>Spin multiplicity</b>	1



## Associated calculations

Job type	Author	Description
OPT	Brice Harismendy	N/A

## Authorship

<b>Original log file</b>	<a href="#">Download</a>
<b>Primary author</b>	Brice Harismendy
<b>Affiliation</b>	None

## Computational details

<b>Software</b>	GAMESS (1MAY2013)
<b>Computational method</b>	DFT
<b>Functional</b>	B3LYP/6-31G*
<b>Basis set name</b>	6-31G*

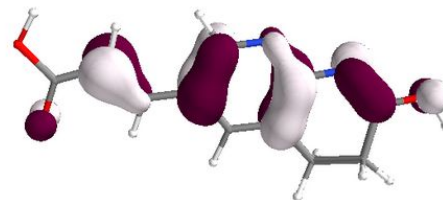
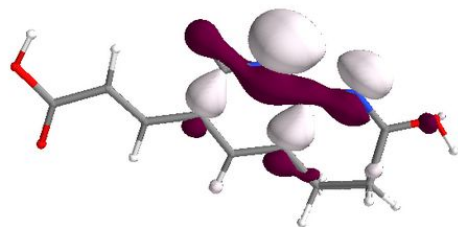
[Molecule](#)[Associated calculations](#)[Authorship](#)[Computational details](#)**🔧 Results**[> Geometry](#)[> Thermochemistry](#)[> Excited states](#)**Most intense Mulliken atomic charges**

mean = 0.000 e, std = 0.342

Atom	number	Mulliken partial charges
O	15	-0.585
O	13	-0.561
O	14	-0.463
N	12	-0.458
N	11	-0.457
C	2	-0.374
C	1	-0.345
C	10	0.364
H	24	0.414
C	8	0.502
C	9	0.610
H	25	0.407

**Geometry****Nuclear repulsion energy in atomic units**

988.60699 a.u.





# Molecular databases



## PubChemQC Project: A Large-Scale First-Principles Electronic Structure Database for Data-Driven Chemistry

Maho Nakata<sup>†</sup>  and Tomomi Shimazaki<sup>‡</sup> 

<sup>†</sup> Advanced Center for Computing and Communication, RIKEN, 2-1 Hirosawa, Wako, Saitama 351-0198 Japan

<sup>‡</sup> Advanced Institute for Computational Science, RIKEN, 7-1-26 Minatojima-minami-machi, Chuo-ku, Kobe, Hyogo 650-0047 Japan

*J. Chem. Inf. Model.*, 2017, 57 (6), pp 1300–1308

DOI: 10.1021/acs.jcim.7b00083

Publication Date (Web): May 8, 2017

Copyright © 2017 American Chemical Society

\*E-mail: maho@riken.jp.

 Cite this: *J. Chem. Inf. Model.* 57, 6, 1300-1308

 RIS Citation 

ioChem-BD Browse - Barcelona Supercomputing Center · BSC host

### PubChemDFT Collection home page



This is a **live** project. It uses spare computer time to compute, process, store, and publish open-access DFT results of molecules contained in the PubChem database. For each entry, we provide its optimized geometry, energies, charges, vibrational frequencies, cube files for electron density and electrostatic potential, etc ... Average rate is close to 1000 moles/day. Started in April 2017, in July 2017 70000 molecules have been completed.

Our Twitter bot @MolecuBot, born July 2017, tweets each time a new molecule is completed and published.

### Discover

#### Author

Central, ioChem-BD **182959**

#### Program name

Gaussian **182959**

#### Calculation type

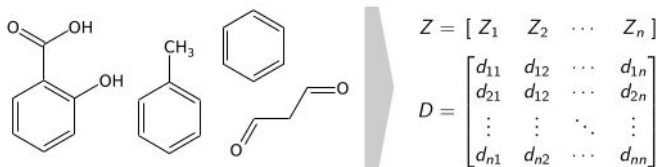
Geometry optimization Minimum **182959**

# Machine learning

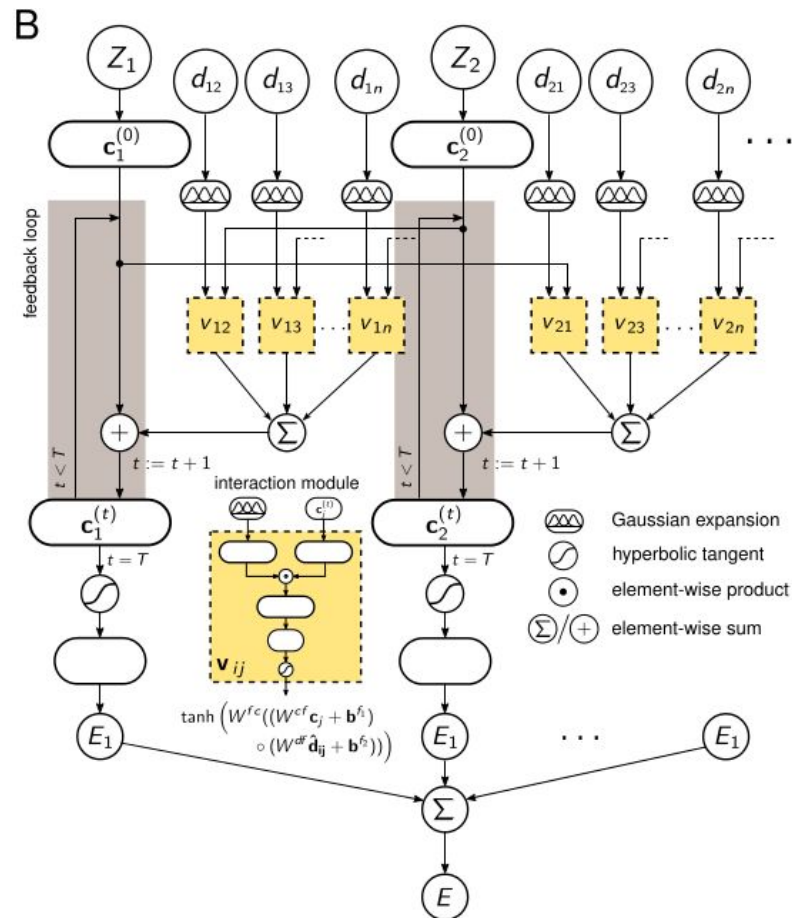
First step :  
Predicting interatomic distances

# Related works

[Schütt et al.] Quantum-Chemical Insights from Deep Tensor Neural Networks. *Nature Communications* 2017.



- GDB-9 dataset :
  - combinatorial molecular space with 9 heavy atoms: C, N, O and F.
  - ~134k small theoretical molecules
- full distances matrix (D): scaling issue
- “only” energy prediction



# Preliminary results

- PubChemQC dataset
  - > 3 millions of real molecules
  - general sampling of the real molecular space (organic chemistry)
- homogeneous data (DFT, B3LYP, 6-31G\*)
- simple neural networks (from 3 up to 9 fully connected layers)
- designed for strong scaling

Model	Objective	Given data
1	All interatomic distances	(partial) distance matrix
2	1 distance (CC, CH, OH)	all distances with other atoms
3	1 distance (CC, CH, OH)	distances with neighbors

# Partial distances matrix with trilateration

Distances better than coordinates: invariant by rotation and translation

But full distances matrix for  $n$  atoms  $\rightarrow n^2$  distances **bad scaling**

How to reconstruct converged geometry? 3D trilateration

Solution : distances from fixed points  $\rightarrow 4n$  distances **good scaling**

$d_{a_0, p_0}$	$d_{a_0, p_1}$	$d_{a_0, p_2}$	$d_{a_0, p_3}$
$d_{a_1, p_0}$	$d_{a_1, p_1}$	$d_{a_1, p_2}$	$d_{a_1, p_3}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$d_{a_n, p_0}$	$d_{a_n, p_1}$	$d_{a_n, p_2}$	$d_{a_n, p_3}$

+ Atomic masses [ $m_{a_0}, \dots, m_{a_n}$ ]

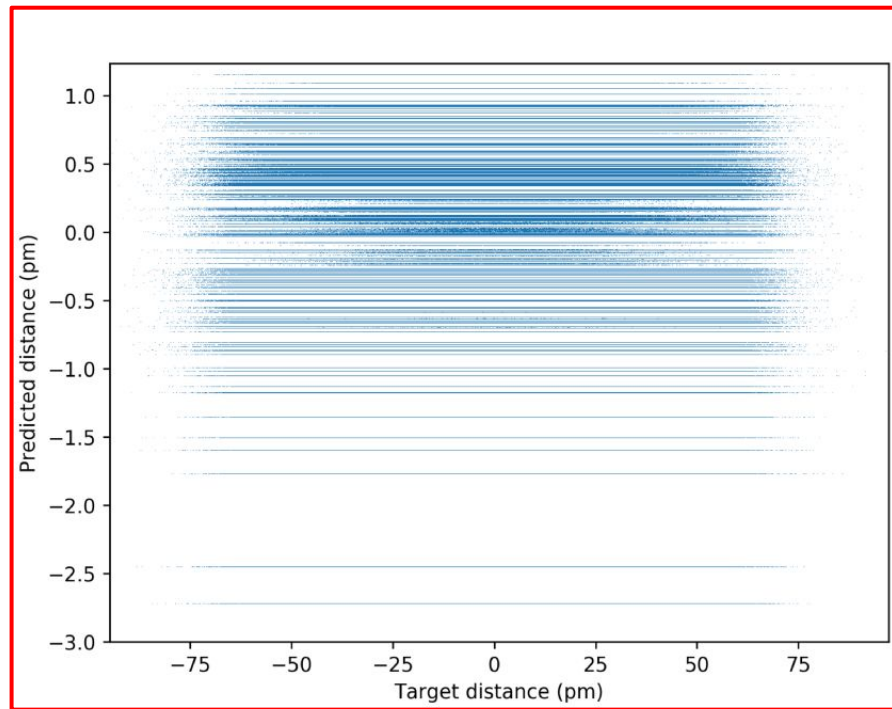
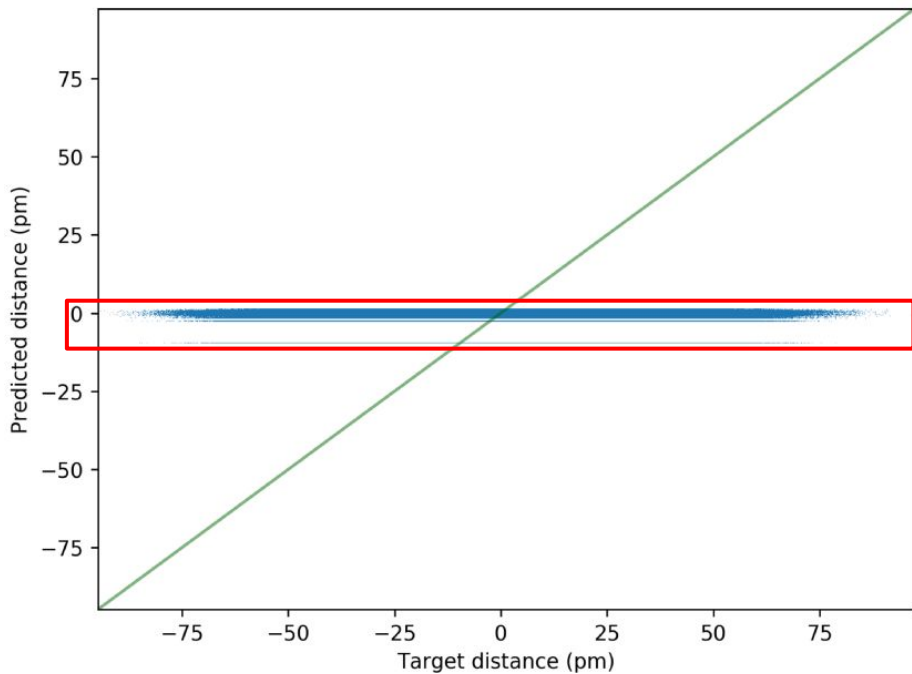
# Many tested models !

Neural networks with fully connected layers:

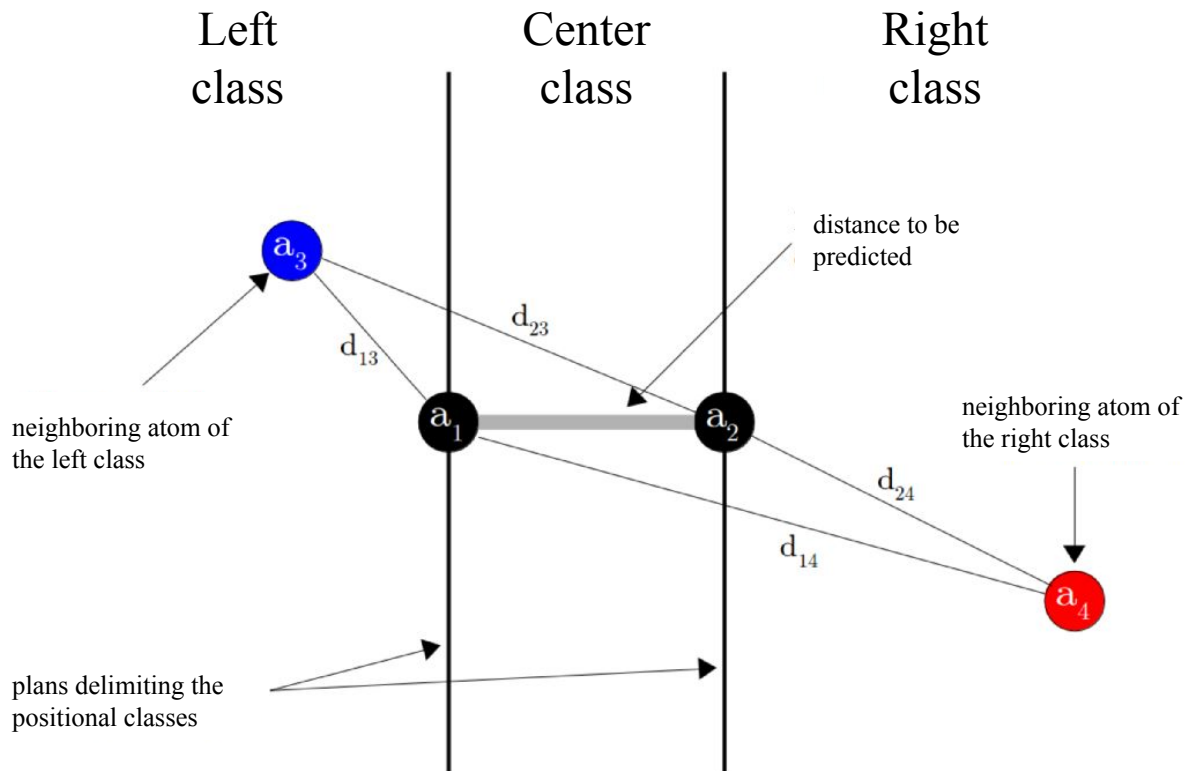
- Loss function: RMSE
- Learning rate: {0.1, 0.0001, 0.00001}
- Adam Optimizer  $\epsilon$ : {1000, 0.0001}
- Weights initialization : {0.2, 0.002}
- Hidden layers activation function: {elu, crelu}
- Exit layer activation function: linear
- Weight decay: {0.1, 0.01, 0.001}
- Layers width: 500 (up to 100 atoms per molecule)
- Networks depth: {3, 7}
- Batch size : {500, 200}
- Number of epochs : 3

→ Networks trained to predicted  $\Delta d$  from starting distances to converged distances

# Results

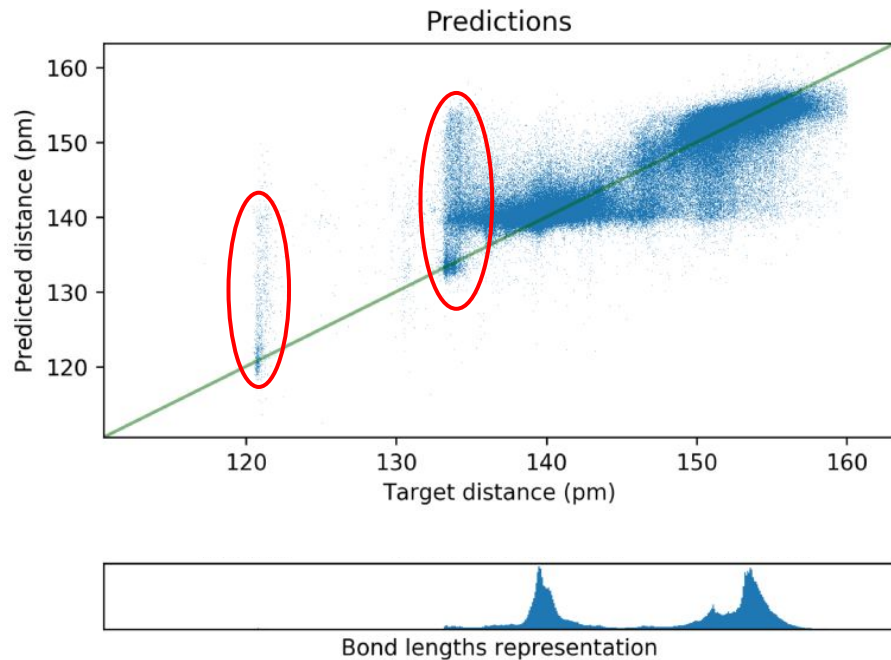
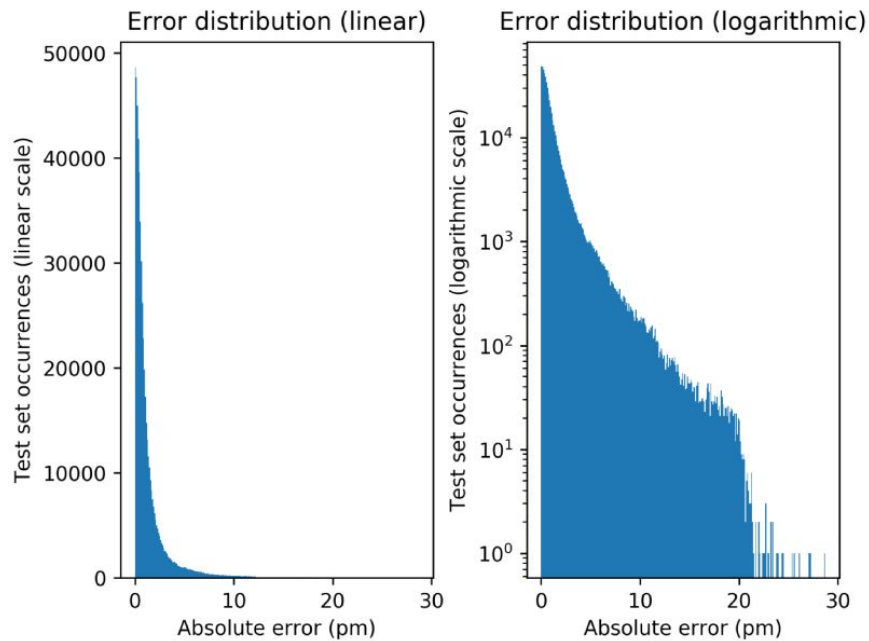


# One interatomic distance (CC, OH and CH)





# Results for Carbon-Carbon



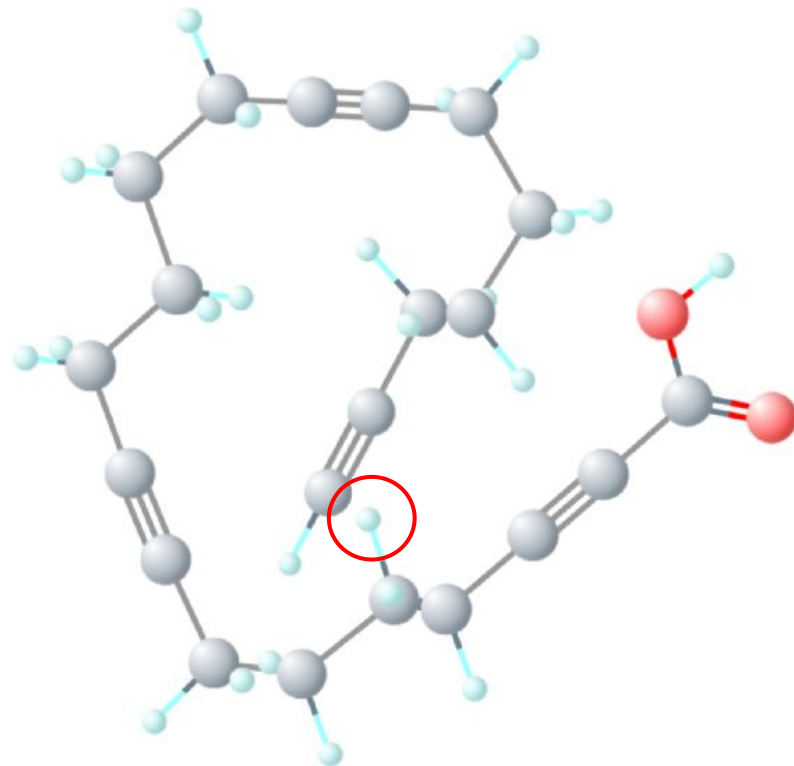
# Example of folded molecule (Pubchem CID 328310)

Unusually close atoms = “false” bonding.

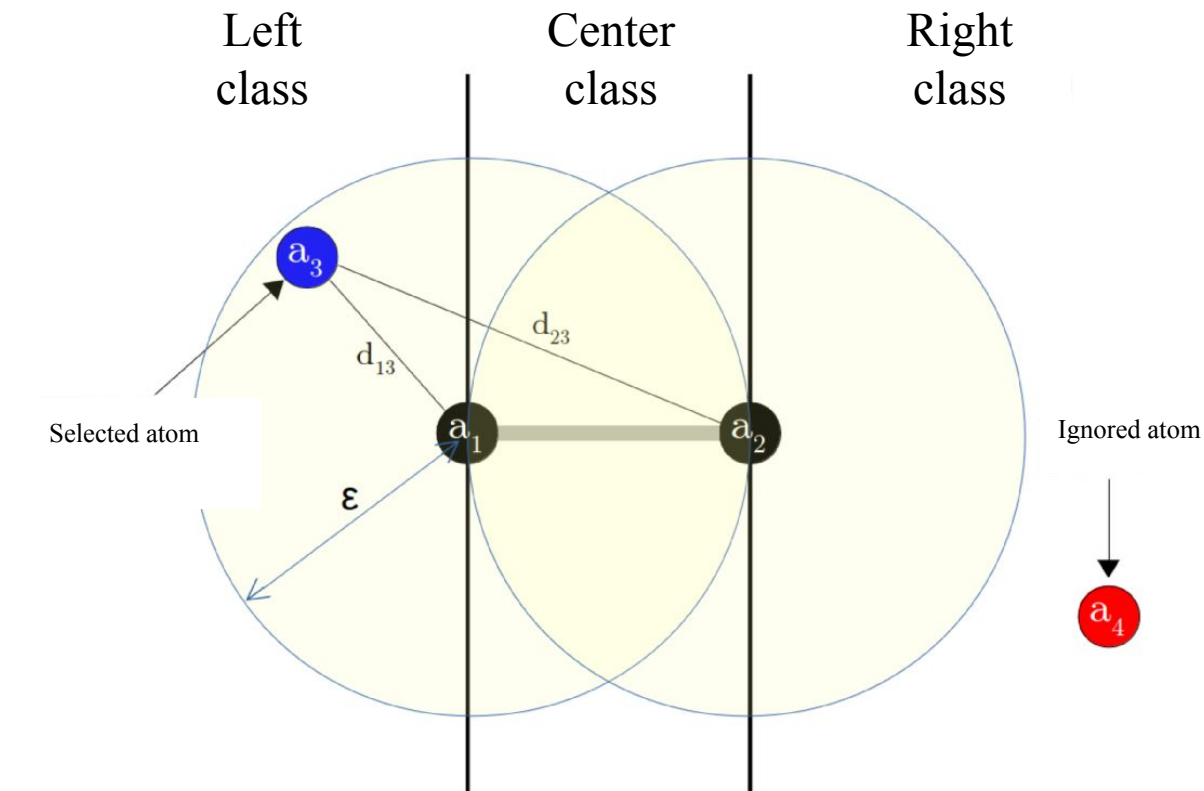
Bad surrounding predicted distances.

**Neural network needs**  
**feature engineering**  
**domain specific knowledge**

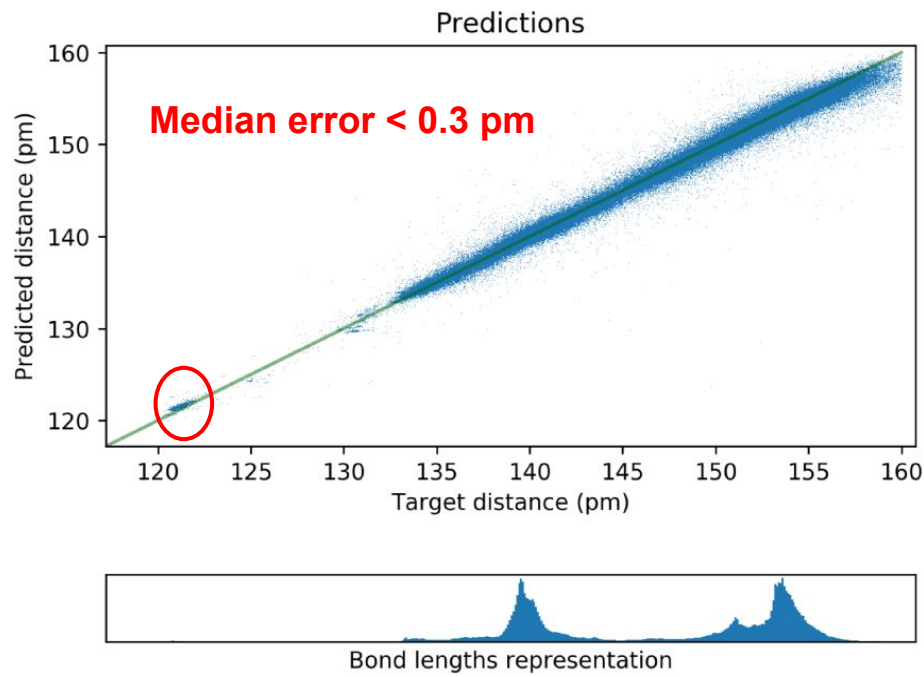
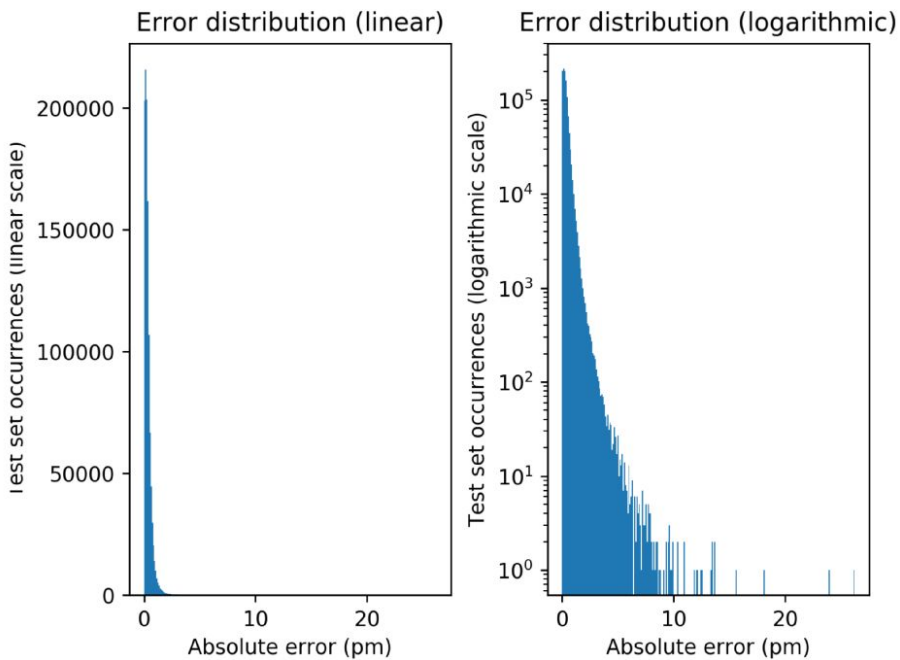
Limitation to the covalent neighboring region! (200 pm)



# One interatomic distance with neighborhood



# Results for Carbon-Carbon



**20 problematic CC bonds on non curated data!**

# Conclusion and perspectives

- Median error for CC < 0.3 pm; for CH < 0.2 pm and for OH < 0.1 pm
  - > 3 millions of real molecules (previously non curated)
  - general sampling of the real molecular space (organic chemistry)
- **Post-doc position available (now) on:**
  - iterative geometry optimization **coupling different models** (NN, KRR,...)
  - prediction of the **wavefunction** or UV-visible absorption/emission
  - **generative models** (GAN, autoencoders) or **combinatorial optimization** algorithms (objective and neighboring function) to efficiently explore the molecular space