Université de Toulouse

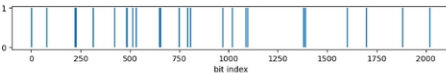# Python in the Physical Chemistry Lab

# [ pyPhysChem ]

**"Talktorials" in physical chemistry and data science / machine learning**

S. Christodoulou, Iann C. Gerber, F. Jolibois, R. Poteau
Python in the Physical Chemistry lab (pyPhysChem) github repository, release v. 1.8.0 (2023)

DOI 10.5281/zenodo.10198844

- interactive python
- images / videos
- mathematical equations
- enriched text (markdown)

Université de Toulouse

**integration of verbal explanations with numerical demonstrations or computer algebra system-based demonstrations** proves to be an influential pedagogical tool

**let's call them "talktorials"**



talktorials specifically tailored for computational chemistry and data science/machine learning

**such mixing is not new, but until recently it was restricted to rather simple applications that required a great deal of development effort**

**What is new is:**

 - **the combination of Python's popularity and libraries**

 - **the interactive nature of Jupyter Notebooks**

 - **personal computers performance; the prevalence of real-world applications that can quite easily be adapted for students thanks to Python libraries available in a lot of domains**

 - **the strong community support**

 - **the ease of reproducibility that makes tutorials more effective, as learners can directly use the code provided to experiment and build upon it**

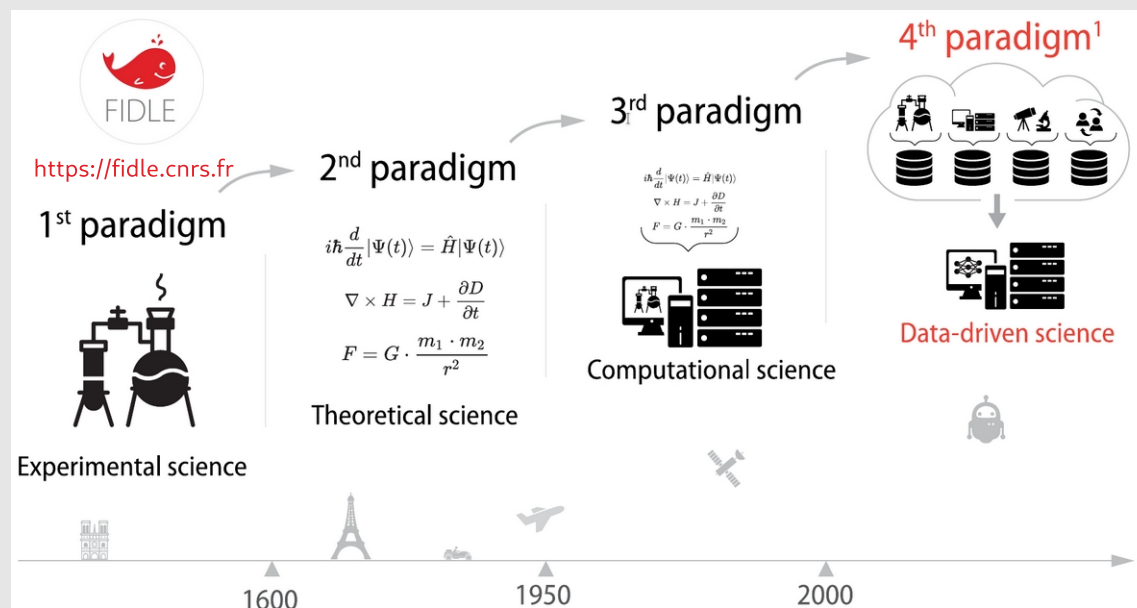**The emergence of such innovative approaches in the realm of computational chemistry is truly encouraging**
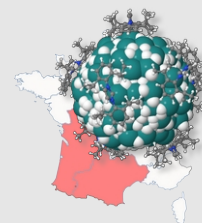
**It not only enables learners to grasp theoretical concepts but also offers a practical perspective on their application**

**For students specialising in computational chemistry who develop their own scripts, they acquire a dual skill set that could be sought after in various areas of research and industry**

**masters' (graduate) degrees**

**bachelor degree**

**all chemistry students, 2ⁿᵈ year**

Réseau Français de
Chimie Théorique

Traitement statistique de données
(data science pour débutants)

*Statistical treatment of data
(data science for beginners)*

Lecture et analyse de la base de données "iris" par la bibliothèque pandas
*Reading and analyzis of the "iris" database with the pandas library*

Ce sujet exploite une base de données souvent utilisée pour l'apprentissage de méthodes statistiques, la base **IRIS** :

- elle regroupe les caractéristiques de trois espèces de fleurs d'Iris : Setosa, Versicolor et Virginica
- la base regroupe 50 observations par espèce (soit 150 **individus**)
- chaque observation repose sur 4 caractéristiques (c'est-à-dire 4 **variables**): longueur et largeur de sépales ainsi que longueur et largeur de pétales

Un article wikipedia porte sur ce dataset, qui contient à la fois des données numériques (largeur & longueur de pétales et sépales) et descriptives (types d'iris).

pétale    sépale

setosa    virginica    versicolor

This subject uses a database often used for the training of statistical methods, the **IRIS** database:

Université
de Toulouse

pyPhysChem

Python in the
**Physical Chemistry Lab**
🐍 python™

**[ pyPhysChem ]**

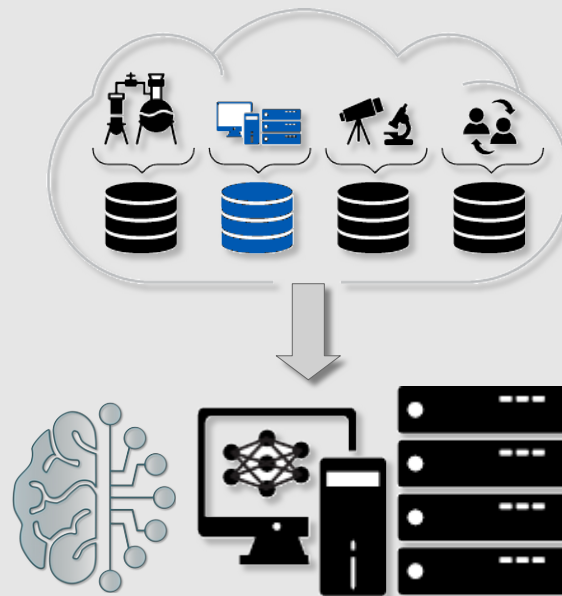https://github.com/rpoteau/pyPhysChem

DOI 10.5281/zenodo.10198844

- python for physicists and chemists in a nutshell

- Computer Algebra System

- Physical chemistry (incl. quantum chemistry)

- coding and use of representations of molecular structures and related data

- Data science and ML

can we expect a strong convergence between quantum and computational chemistry, data science and machine learning?

data-driven science

**increase in the number of students in the master's programmes in theoretical and computational chemistry?**

**... unless we do not really give them a dual skill set that could be sought after in various areas of research and industry**

Université
de Toulouse