



Outline

- Why care?
- Do and Don't data processing
- Simulation data processing
- Research coding basics (Code editor, Linting, Formating, Testing,
 Documentation, start a project, coding environments)
- Use AI productively for your research data workflows
- Sharing code (Webapps, Colab, packaging, Zenodo)

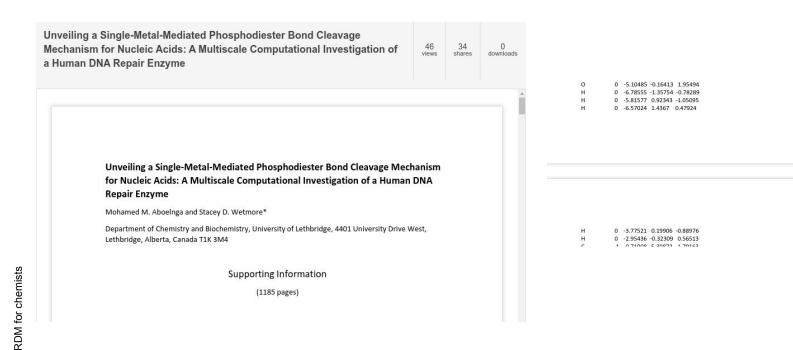


Why care?

S24



Chemists generate a lot of data and we're required to share it in Switzerland





RETURN TO ISSUE

EDITORIAL

NEXT >

Guidelines for Reporting Molecular Dynamics Simulations in JCIM Publications

Thereza A. Soares*, Zoe Cournia, Kevin Naidoo, Rommie Amaro, Habibah Wahab, and Kenneth M. Merz Jr.

Cite this: J. Chem. Inf. Model. 2023, 63, 11, 3227–3229

Publication Date: May 16, 2023 V

https://doi.org/10.1021/acs.jcim.3c00599

Copyright © Published 2023 by American Chemical Society

Request reuse permissions Sul

Subscribed

Article Views Altmetric Citations

10138 66 2

LEARN ABOUT THESE METRICS

Input/output files and data availability

ARTICLE SECTIONS

Jump To ∨

JCIM now requires depositing your input files (including topologies, parameters, input files, and initial configurations) and output trajectory files in a public repository providing a DOI and provide the accession link in the manuscript. Some available repositories are Zenodo (https://zenodo.org/), Dryad (https://datadryad.org/), Nomad (https://nomad-repository.eu/), and Figshare (https://figshare.com/). If the full trajectories are too large, a selection of clustered structure representing the respective trajectories should be made available in the public repository.

Real problems from my own PhD

People not sharing their data/models even after inquiring multiple times

Compiled binary only for 32bit systems and environment not specified

Simon Duerr

	Name	Website	Related Metals	Main Algorithm	Reported Performance	Ref.
	RDGB	http://www.cerm.unifi.it/home/research/genomebrowsing.html	Zn, Cu, Fe and other metals	Integration of tools for retrieval of protein domains and genome analysis	Accuracy: 89.6%, precision: 85.9%	[32]
review from 20221	Zincfinder	http://zincfinder.dsl.unifi.it	Zn	a SVM	^b AURPC: 0.590 (local predictor) and 0.633 (gating network)	[33]
	Zincpred	http://www.fos.su.se/~nanjiang/zincpred/download/	Zn	SVM- and homology-based algorithm	AURPC: 0.723 (local predictor) and 0.701 (gating network)	[34]
	TEMSP	http://netalign.ustc.edu.cn/temsp/	Zn	Structure-based algorithm with a range of geometric criteria	° AUC: 0.945	[35]
	Zincidentifier	http://protein.cau.edu.cn/zincidentifier/	Zn	A two-step feature selection method based on random forest algorithm	AUC: 0.955, AURPC: 0.829	[36]
	ZincExplorer	http://protein.cau.edu.cn/ZincExplorer/	Zn	A combination of SVM-, cluster- and template-based predictors	AURPC: 0.907	[37]
	ZincBinder	http://proteininformatics.org/mkumar/znbinder/	Zn	SVM model trained on PSSM-based input feature	AUC: 0.91	[38]
	ZINCCLUSTER	http://www.metalactive.in	Zn	SVM-based Ligand Finder and Cluster Finder algorithms	^d MCC: 0.798, F1-score: 0.801	[39]
	ZnMachine	http://bioinformatics.fzu.edu.cn/znMachine.html	Zn	A combination of several intensively-trained machine learning models	AUC: 0.933 (SVM) and 0.910 (neural network)	[40]
	HemeBIND	http://mleg.cse.sc.edu/hemeBIND/	Fe (heme)	SVM	MCC: 0.504, F1-score: 56.87%	[41]
	SCMHBP	http://iclab.life.nctu.edu.tw/SCMHBP/	Fe (heme)	Based on a newly-developed scoring card method for predicting heme-binding proteins	Accuracy: 85.90%	[42]
	Isph	http://biodev.extra.cea.fr/isph	Fe (Fe-S)	A penalized linear model based on machine learning approach	Precision: 87.9%, recall: 80.1% (extended model)	[43]
	MetalPredator	http://metalweb.cerm.unifi.it/tools/metalpredator/	Fe (Fe-S)	Integration of existing domain-based methodology with a new approach for discovering metal-binding motifs	Precision: 85.2%, recall: 88.6%	[44]
	MetSite	N/A	Fe, Zn, Cu, Mn, Ca, Mg	Artificial neural network	Mean accuracy: 94.5%	[45]
RDM for chemists	FINDSITE-metal	http://cssb.biology.gatech.edu/findsite-metal/	Fe, Zn, Cu, Mn, Ni, Co, Ca, Mg	Integration of structure/evolutionary information and machine learning approach (SVM)	Overall accuracy: 70–90%	[46]
	SeqCHED	http://ligin.weizmann.ac.il/seqched	Fe, Zn, Cu, Mn, Ni, Co, Ca, Mg	A modification of the CHED algorithm and machine learning filters (decision tree classifier and SVM)	Sensitivity: 84–85%, selectivity: 82–93% (stringent filtration)	[47]
	MetalDetector	http://metaldetector.dsi.uniff.it/v2.0/	Transition metals that use cysteine and histidine as ligands	A combination of different machine learning algorithms (SVM-HMM, structured-output SVM)	Precision: 60–79%, recall: 71–88%	[48]
	MIB	http://bioinfo.cmu.edu.tw/MIB/	Ca, Cu, Fe, Mg, Mn, Zn, Cd, Ni, Hg, Co	Fragment transformation method	Overall accuracy: 92.9– 95.1%	[49]

Table 1. Computational tools for metal-binding site and metalloprotein prediction.

Name	Website	Related Metals	Main Algorithm	Reported Performance	Re
RDGB	http://www.corm.unifi.it/home/research/genomebrowsing.html	Zn, Cu, Fe and other metals	Integration of tools for retrieval of protein domains and genome analysis	Accuracy: 89.6%, precision: 85.9%	
Zincfinder	http://zincfinder.dsi.unifi.it	Zn	a SVM	b AURPC: 0.590 (local predictor) and 0.633 (gating network)	[3:
Zincpred	http://www.tos.su.se/-nanjjang/sinoprod/download/	Zn	SVM- and homology-based algorithm	AURPC: 0.723 (local predictor) and 0.701 (gating network)	[3
TEMSP	http://netalign.ustc.edu.cn/temsp/	Zn	Structure-based algorithm with a range of geometric criteria	^c AUC: 0.945	[3
Zincidentifier	http://protein.cau.edu.en/zhicidentifier/	Zn	A two-step feature selection method based on random forest algorithm	AUC: 0.955, AURPC: 0.829	[3
ZincExplorer	http://protein.eau.edu.en/ZineExplorer/	Zn	A combination of SVM-, cluster- and template-based predictors	AURPC: 0.907	[3
ZincBinder	http://proteininformatics.org/mkumar/znbinder/	Zn	SVM model trained on PSSM-based input feature	AUC: 0.91	[3
ZINCCLUSTER	-http://www.metalactive.in	Zn	SVM-based Ligand Finder and Cluster Finder algorithms	^d MCC: 0.798, F1-score: 0.801	[3
ZnMachine	http://bioinformatics.rzu.edu.cn/zniMachine.html	Zn	A combination of several intensively-trained machine learning models	AUC: 0.933 (SVM) and 0.910 (neural network)	[4
HemeBIND	http://mleg.cse.sc.edu/hemeBIND/	Fe (heme)	SVM	MCC: 0.504, F1-score: 56.87%	[4
SCMHBP	http://iclab.life.nctu.edu.tw/SCMHBP/	Fe (heme)	Based on a newly-developed scoring card method for predicting heme-binding proteins	Accuracy: 85.90%	[4
Isph	http://blodev.extra.cea.tr/isph	Fe (Fe-S)	A penalized linear model based on machine learning approach	Precision: 87.9%, recall: 80.1% (extended model)	[4
MetalPredator	http://metalweb.cerm.unifi.it/toois/metalpredator/	Fe (Fe-S)	Integration of existing domain-based methodology with a new approach for discovering metal-binding motifs	Precision: 85.2%, recall: 88.6%	[4
MetSite	N/A	Fe, Zn, Cu, Mn, Ca, Mg	Artificial neural network	Mean accuracy: 94.5%	[4
FINDSITE-metal	http://essb.biology.gatech.edu/findsite-metal/	Fe, Zn, Cu, Mn, Ni, Co, Ca, Mg	Integration of structure/evolutionary information and machine learning approach (SVM)	Overall accuracy: 70–90%	[4
SeqCHED	_http://ligin.weizmann.ac.il/seqchcd	Fe, Zn, Cu, Mn, Ni, Co, Ca, Mg	A modification of the CHED algorithm and machine learning filters (decision tree classifier and SVM)	Sensitivity: 84–85%, selectivity: 82–93% (stringent filtration)	[4
MetalDetector	<u>http://metaldetector.dsi.unifi.it/v2.0/</u>	Transition metals that use cysteine and histidine as ligands	A combination of different machine learning algorithms (SVM-HMM, structured-output SVM)	Precision: 60–79%, recall: 71–88%	[4
MIB	this version no longer available http://bioinfo.cmu.edu.tw/MIB/	Ca, Cu, Fe, Mg, Mn, Zn, Cd, Ni, Hg, Co	Fragment transformation method	Overall accuracy: 92.9- 95.1%	[4

Do and Don't data processing

No manual file editing

No manual file copying

Clear track record from raw data to plots/tables (e.g via scripting)

Go from Python to Microsoft Word

https://rowannicholls.github.io/py thon/data/export to word.html



Programming language













Simplicity is key

Multiply 2d arrays of size (100,100)









fast

RDM for chemists

Simulation data processing

Running a lot of computations?

Use workflow tools to keep track:

Computational Chemistry

- AIIDA
- atomate2/jobflow

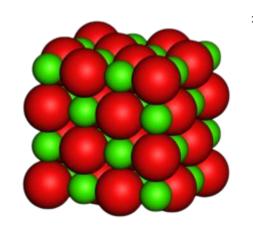
Machine Learning

Weights and Biases

Old school

Create structure in Avogadro or GaussView

Run pw.x < MGO.opt.in > MGO.opt.out



Extract energies grep ! MGO.opt.out > energies.dat

Plot them using gnuplot by entering the commands

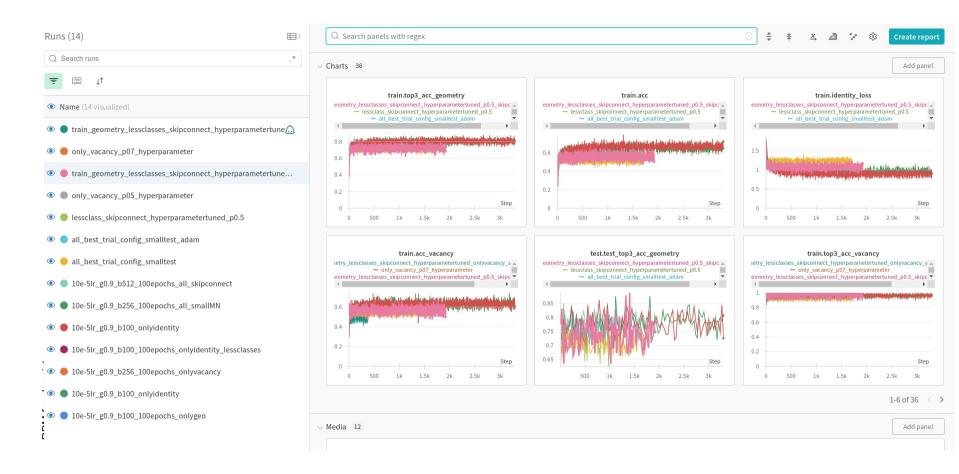
> plot 'energies.dat' using 1:2 with lines

struct_opt_band_calc.py

```
from atomate2.vasp.flows.core import RelaxBandStructureMaker
from jobflow import run_locally
from pymatgen.core import Structure
# construct a rock salt MgO structure
mgo_structure = Structure(
    lattice=[[0, 2.13, 2.13], [2.13, 0, 2.13], [2.13, 2.13, 0]],
    species=["Mg", "0"],
    coords=[[0, 0, 0], [0.5, 0.5, 0.5]],
# make a band structure flow to optimise the structure and obtain the band structure
bandstructure_flow = RelaxBandStructureMaker().make(mgo_structure)
# run the flow
run_locally(bandstructure_flow, create_folders=True)
```

EPFL

Weights and Biases



Prefer open software to closed software

Make everything explicit instead of running single commands

RDM for chemists

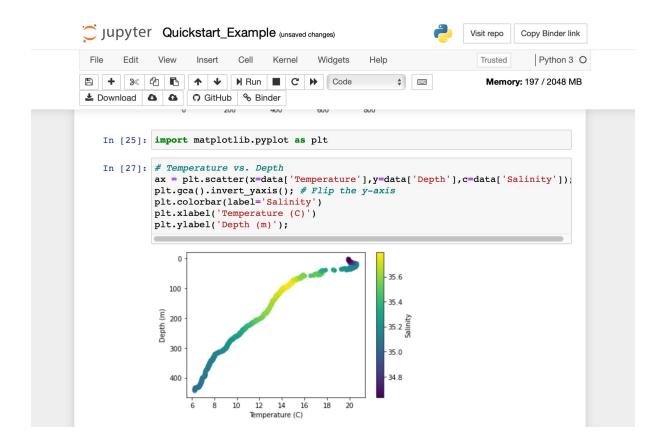
EPFL

Chances that you have to redo something multiple times are high.

Use automated scripts as much as possible to generate plots directly



Let's talk about Jupyter Notebooks







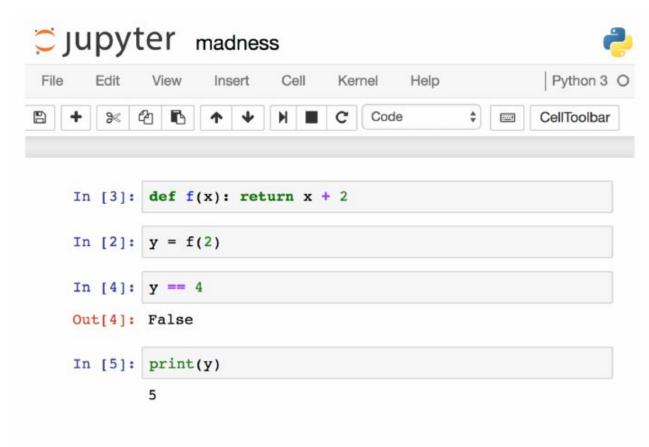


JAKE-CLARK.TUMBLR

In the same of the







if you look at the numbers, it's clear those cells weren't even executed in order!



@joelgrus #jupytercon









input functions

processing functions



main program

Notebook

VS.

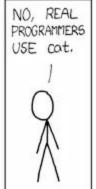
Script

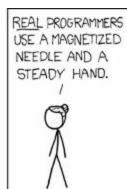
VS.

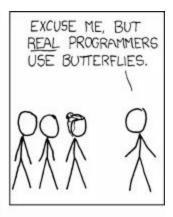
Code













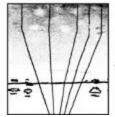
THE DISTURBANCE RIPPLES OUTWARD, CHANGING THE FLOW OF THE EDDY CURRENTS IN THE UPPER ATMOSPHERE.

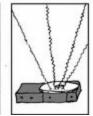


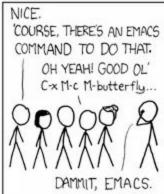


THESE CAUSE MOMENTARY POCKETS OF HIGHER-PRESSURE AIR TO FORM,

WHICH ACT AS LENSES THAT DEFLECT INCOMING COSMIC RAYS, FOCUSING THEM TO STRIKE THE DRIVE PLATTER AND FLIP THE DESIRED BIT.







EPFL

Code editor

Many choices but if you don't have a strong preference for an editor:

VS Code offers:

- the right abstractions,
- version control,
- inbuilt terminal,
- linting,
- formatting,
- Al completion,
- Jupyter support and
- remote capabilities

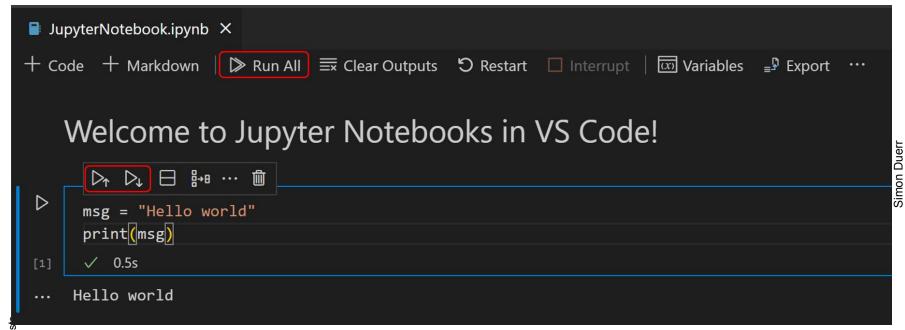
Why VS Code remains a developer favorite, year after year

Anastasija Uspenski



Notebooks in VS Code





Syntax Highlighting, Linting, Formatting

PDM for chami

EPFL

Tools for coding

- Linting
- Formatting
- Code environments
- Version control
- Continous integration
- Documention

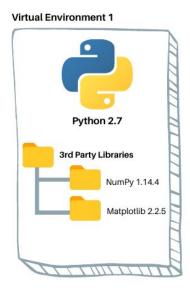
EPFL

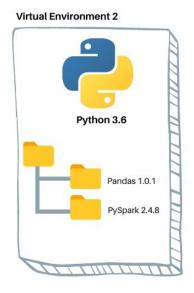
Linting: Automatically detect errors

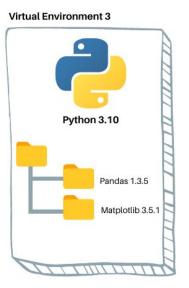
```
▷ ~ □ …
hello.py 2 •
 nello.py > ...
        msg = "Hello, Python"
        print msg
              flush: bool
         ) -> None: ...
         print(value, ..., sep=' ', end='\n', file=sys.stdout, flush=False)
         Prints the values to a stream, or to sys.stdout by default. Optional keyword arguments:
         file: a file-like object (stream); defaults to the current sys.stdout.
         sep: string inserted between values, default a space.
         end: string appended after the last value, default a newline.
         flush: whether to forcibly flush the stream.
         Parsing failed: 'Missing parentheses in call to 'print'. Did you mean
         print(...)? (<unknown>, line 3)' Pylint(E0001:syntax-error)
         View Problem (Alt+F8) No quick fixes available
                                                                                        PROBLEMS 2
                          TERMINAL
  hello.py 2
     🔞 Parsing failed: 'Missing parentheses in call to 'print'. Did you mean print(...)? (<unknown>, line 3)' Pylint(E0001:syntax-error) [Ln 3, Col 2]
    Statements must be separated by newlines or semicolons Pylance [Ln 3, Col 7]
```

Code environments

Do not install everything in one big environment









Different options:

Python version (e.g. Python 3.9.7)

(e.g. nltk 3.6.2)

hard

venv

conda

Package version venv/pip

- smaller/lightweight
- cannot change python version
- only python packages from pypi
- conda/mamba
 - change python version
 - heavier/slower
 - install also other scientific software

EPFL

Formatting



```
from seven dwwarfs import Grumpy, Happy, Sleepy, Bashful, Sneezy,
    x = \{ 'a':37, 'b':42, 
                                                                               x = {\text{"a": 37, "b": 42, "c": 927}}
     'c':927}
    x = 123456789.123456789E123456789
                                                                               if (
    if very long variable name is not None and \
 8
     very long variable name.field > 0 or \
     very long variable name.is debug:
     z = 'hello '+'world'
                                                                                  z = "hello " + "world"
12
    else:
                                                                          13 else:
                                                                          14
                                                                                  world = "world"
13
     world = 'world'
                                                                                  a = "hello {}".format(world)
     a = 'hello {}'.format(world)
                                                                          16
                                                                                  f = rf"hello {world}"
15
     f = rf'hello {world}'
                                                                          17 if this and that:
    if (this
     and that): y = 'hello ''world'#FIXME: https://github.com/psf/blac
    class Foo
                      object ):
                                                                          20
19
      def f
                (self ):
                                                                          21 class Foo(object):
20
        return
                      37*-2
                                                                          22
                                                                                  def f(self):
      def q(self, x, y=42):
                                                                          23
                                                                                      return 37 * -2
22
           return y
    def f ( a: List[ int ]) :
                                                                          25
                                                                                  def g(self, x, y=42):
24
                   37-a[42-u: v**3]
                                                                          26
       return
                                                                                      return y
```

```
from seven dwwarfs import Grumpy, Happy, Sleepy, Bashful, Sneezy, Do
    x = 123456789.123456789e123456789
        very long variable name is not None
        and very long variable name.field > 0
        or very long variable name is debug
        y = "hello " "world" # FIXME: https://github.com/psf/black/issu-
27
28
29 def f(a: List[int]):
        return 37 - a[42 - u : y**3]
```

format everytime on save and ensure consistent formatting with Python standards

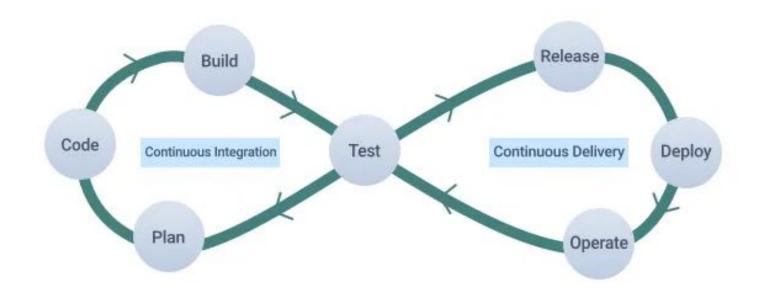
Testing

Make sure code does what it is supposed to do

```
# content of test_sample.py
def inc(x):
    return x + 1

def test_answer():
    assert inc(4) == 5
```

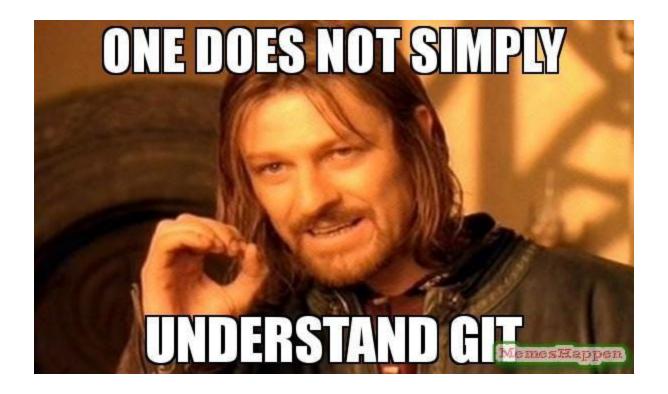




RDM for chemists

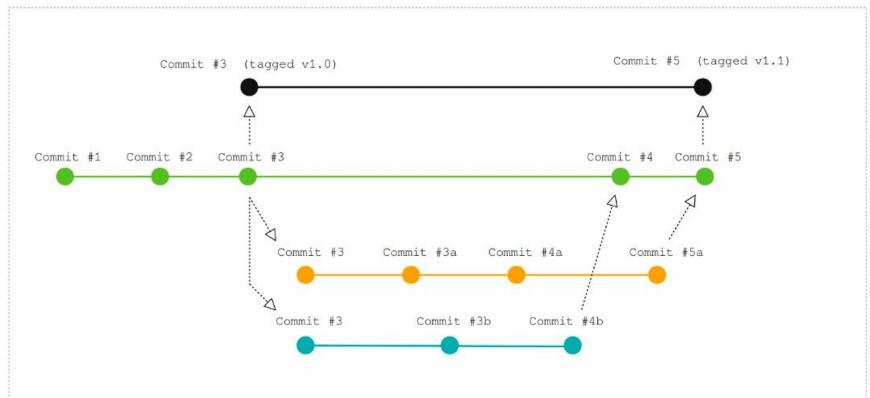
https://medium.com/digital-transformation-and-platform-engineering/a-quick-guide-to-continuous-integration-and-continuous-delivery-4df594 ae281

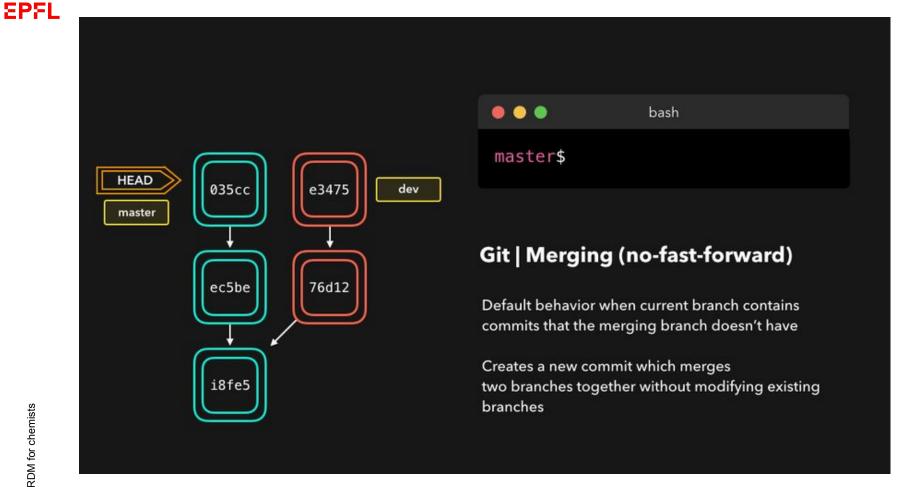
Version control

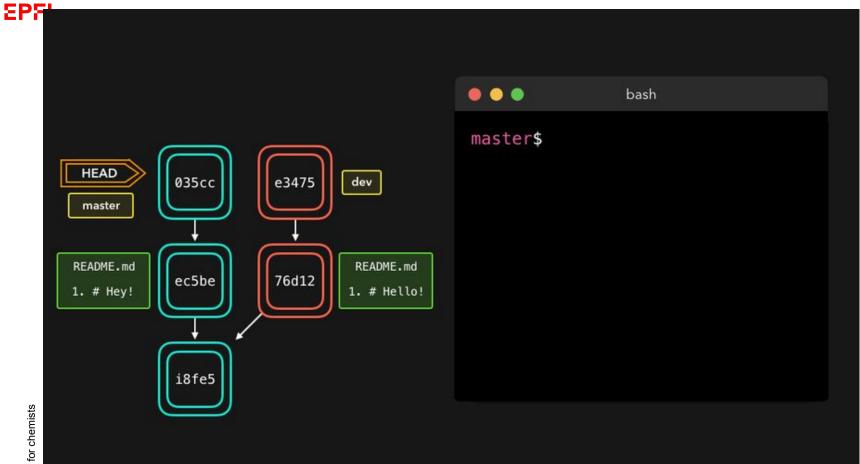


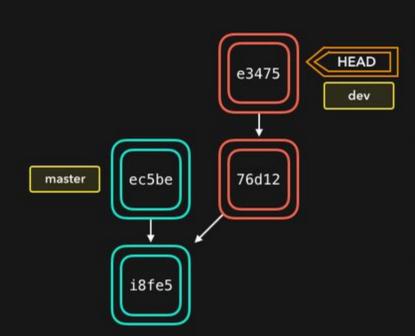
EL-.

Master Branch (Production) Develop Branch Feature #1 Branch (Developer 1) Feature #2 Branch (Developer 2)

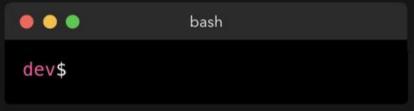








ΞŖ



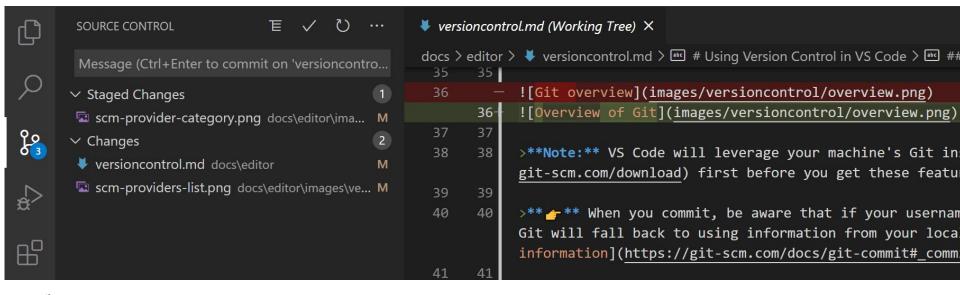
Git | Rebasing

Copies commits on top of another branch without creating a commit, which keeps a linear history

Changes the history as new hashes are created for the copied commits



VSCode can also help with git workflows



Simon Duerr

Documentation

```
using System;struct a{static int Main()
{object[]c={"\u0048e\x6c\x6co "+(C\u0068ar)
(86+1)+"or\x6c\x64"};typeof(Conso\u006ce).GetMet
\u0068o\u0064s()[101].Invoke(c,c);return 0;}}
```

worst

script.sh

```
. .
def simple_math(a, b):
    return a + b
def main():
    num1 = 10
    num2 = 5
    result = simple_math(num1, num2)
    print("Number 1:", num1)
    print("Number 2:", num2)
    print("Result:", sum_result)
if __name__ == "__main__":
    main()
```

bad

```
RDM for chemists
```

```
. .
                                 good
def sum(a, b):
   Perform summation two numbers.
    Parameters
   a : int or float
       The first number.
   b : int or float
       The second number.
    Returns
   int or float
       The sum of the two numbers.
   return a + b
def main():
   Entry point of the program.
   num1 = 10
   num2 = 5
   result = sum(num1, num2)
   print("Number 1:", num1)
   print("Number 2:", num2)
   print("Result:", result)
if __name__ == "__main__":
   main()
```

```
.
                                         excellent
import argparse
def sum(a, b):
    Calculate the sum of two numbers.
    Parameters
    a : int or float
        The first number.
    b : int or float
        The second number.
    Returns
    int or float
        The sum of the two numbers.
   return a + b
def main():
    Entry point of the program.
    parser = argparse.ArgumentParser(description="Calculate the
sum of two numbers.")
    parser.add argument("num1", type=int, help="First number")
    parser.add_argument("num2", type=int, help="Second number")
    args = parser.parse_args()
        result = sum(args.num1, args.num2)
    except Exception as e:
       print("An error occurred:", e)
    print("Number 1:", args.num1)
    print("Number 2:", args.num2)
    print("Result:", result)
if __name__ == "__main__":
    main()
```



```
def foo(arg1, arg2):
    """
    A method's docstring with parameters and return value.

    Use all the cool Sphinx capabilities in this description, e.g. to give usage examples ...

:Example:
    >>> another_class.foo('', AClass())

:param arg1: first argument
    :type arg1: string
    :param arg2: second argument
    :type arg2: :class:`module.AClass`
    :return: something
    :rtype: string
    :raises: TypeError
    """
    return '' + 1
```



automatically generate documentation website from code

foo(arg1, arg2)

A method's docstring with parameters and return value.

Use all the cool Sphinx capabilities in this description, e.g. to give usage examples ...

Example:

```
>>> another_class.foo('', AClass())
```

Parameters: • arg1 (string) - first argument

• arg2 (module.Aclass) - second argument

Return type: string
Raises: TypeError

Start a project

```
ERLEST STREET
scishers:girthde materis cookdecutter atcchedensiablesetiacutter-cooptee
project_name (project_name): Drug Minder Finder
first module same filting blander Finder(): drug search
author, name (hour name for year organization/company/name); Lovi N Madem
Geographian IX short description of the project is a book to convictor bor binames than are drugs or stop-time.
Select open morne Issuenie:
 - MIT
- 100-3-Classe
a - unwind
 - Note: made-large over the default assume manufactiff use factoric
 - Profer default eneconds channel with pilo fieldback
5 - Depositionalism from pap only the Limital
Chapter from to 1, 3 (LI):
```

Resources

https://github.com/choderalab/software-development

y master → y 4 Branches ♥ 0 Tags	Q Go to file	<> Code ▼	About	
pichodera Merge pull request #77 from chfw/i	master 4a8fce8 · 4 years ag	139 Commits	A primer on software development best practices for computational	
chapter_images	Fill in the Structure page	7 years ago	chemistry	
CONTINUOUS_INTEGRATION.md	Other minor edits.	5 years ago	python computational-chemistry development-practices	
DOCUMENTATION.md	adding back to top	7 years ago	☐ Readme	
LICENSING_GUIDELINES.md	Edits for readability	5 years ago	- ∕- Activity	
PACKAGING_AND_DEPLOYMENT.md	Link PyPI tutorial to PyPA user guide	6 years ago	© Custom properties ☆ 239 stars	
PHILOSOPHY_AND_SCOPE.md	Fix typo	5 years ago		
PYTHON_CODING.md	Update PYTHON_CODING.md	4 years ago	₹ 79 forks Report repository	
PYTHON_OPTIMIZATION.md	Fixing typos	6 years ago		
README.md	Add MolSSI Education Resources link	5 years ago	Releases	
STRUCTURING_YOUR_PROJECT.md	Add CMS cookiecutter info to structuring your project	5 years ago	No releases published	
UNIT_TESTING.md	UNIT_TESTING.md: mention codecov	4 years ago	Packages	
□ VERSION_CONTROL.md	Cleanup edits to Python Coding section.	5 years ago	No packages published	
□ README		≔	Contributors 16	

What do you really need?

	Research scripts/Notebooks	Research code			
Unit tests	-	yes			
Publish	Zenodo	Zenodo, PyPi, conda-forge, GitHub			
Linting/Formatti ng	yes	yes			
Documentation	inline/docstrings	Web documentation, tutorials (e.g using Sphinx)			
Format	e.g Notebook, scripts	e.g Python modules			
Environment	environment.yml	Docker container			
CI/CD	-	yes			
Versioning	yes, e.g git	yes, e.g git			



Example case:

I had a ML dataset of h5 files and corresponding csv files

I wanted to split the dataset into a smaller test set and validation set.

Simon Duerr

Simon Duerr

Use AI productively for your research data workflows

Be precise and split your commands



You

I have a folder with .h5 files with an identifier. In this folder there is a subfolder "envs" that contains corresponding csv files in the format identifier_envs.csv.

Take all of the files in the folder and move half of them to a folder test and half of them in vals

GPT-3.5 generates non working python code





You

Write a python programm that iterates over a pair of h5 files contained in a folder. Each h5 file has an associated _envs.csv file in the envs folder

generates working example that one can debug



You

Now add functionality to move half the h5/csv pairs to a new directory. Sample them pairs to choose randomly

Adds the required functionalities and results in working code

```
GITHUB COPILOT: CHAT

monalisa

Write unit tests for this function

GitHub Copilot

import unittest
```

The code assumes that the datetime module is imported. The test cases cover different scenarios such as valid input, empty input, input with comments, invalid date format, and invalid value format.

```
import datetime
     def parse expenses (expenses string):
         """Parse the list of expenses and return the list of triples (date, amount, currency
         Ignore lines starting with #.
         Parse the date using datetime.
         Example expenses_string:
             2023-01-02 -34.01 USD
             2023-01-03 2.59 DKK
             2023-01-03 -2.72 EUR
         expenses = []
         for line in expenses_string.splitlines():
             if line.startswith("#"):
             date, value, currency = line.split (" ")
             expenses.append((datetime.datetime.strptime (date, "%Y-%m-%d"),
                             float (value),
                             currency))
21
             return expenses
     expenses data = '''2023-01-02 -34.01 USD
                     2023-01-03 2.59 DKK
                     2023-01-03 -2.72 EUR'''
```

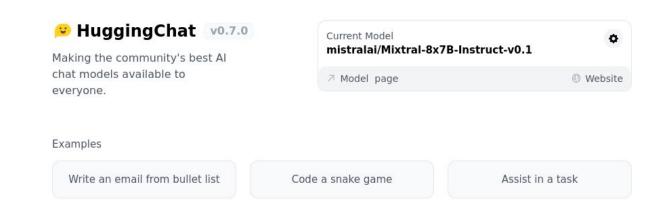
is sentiments.ts X

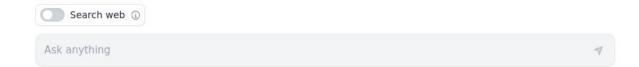
parse expenses.pv ×



RDM for chemists

apart from chatGPT





Simon Duerr

Sharing code

We need our code to be FAIR

indable Accessible Interoperable Reusable

How to package code

Just giving someone a code file does not mean the code will run as intended.

We need:

- data
- code
- environment specification

So how can we do this?

Simon Due

	?	?	?	?	?	?	?
	?	?	?	?	?	?	?
	?	?	?	?	?	?	?
	?	?	?	?	?	?	?
	?	?	?	?	?	?	?
000	?	?	?	?	?	?	?
	2						4

Solution: Intermodal Shipping Container

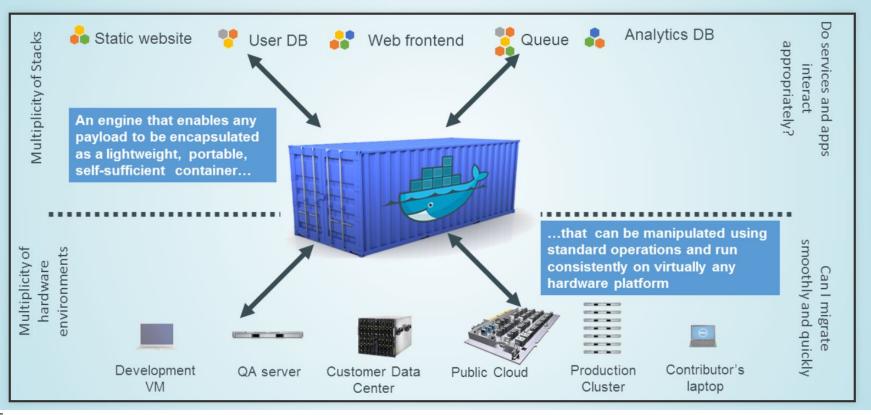




The Matrix from Hell

••	Static website	?	?	?	?	?	?	?
**	Web frontend	?	?	?	?	?	?	?
	Background workers	?	?	?	?	?	?	?
••	User DB	?	?	?	?	?	?	?
•	Analytics DB	?	?	?	?	?	?	?
**	Queue	?	?	?	?	?	?	?
		Development VM	QA Server	Single Prod Server	Onsite Cluster	Public Cloud	Contributor's laptop	Customer Servers
			1					111

Docker is a Container System for Code



Simon Duel

Build once... (finally) run anywhere*

- A clean, safe, hygienic, portable runtime environment for your app.
- No worries about missing dependencies, packages and other pain points during subsequent deployments.

^{*} Where "anywhere" means an x86 server running a modern Linux kernel (3.2+ generally or 2.6.32+ for RHEL 6.5+, Fedora, & related)

```
FROM python: 3.7-slim-buster
COPY install_packages.sh .
RUN ./install_packages.sh
RUN useradd cheminfo
WORKDIR /home/cheminfo
COPY requirements.txt .
COPY dimensionality_reduction ./dimensionality_reduction
COPY README.md .
RUN pip install --no-cache-dir -r requirements.txt
```



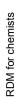
Webapps

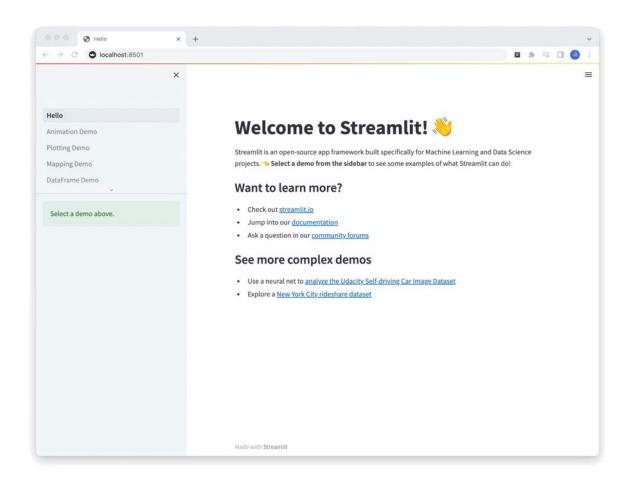
Gradio



RDM for chemists

https://www.gradio.app/guides/guickstart

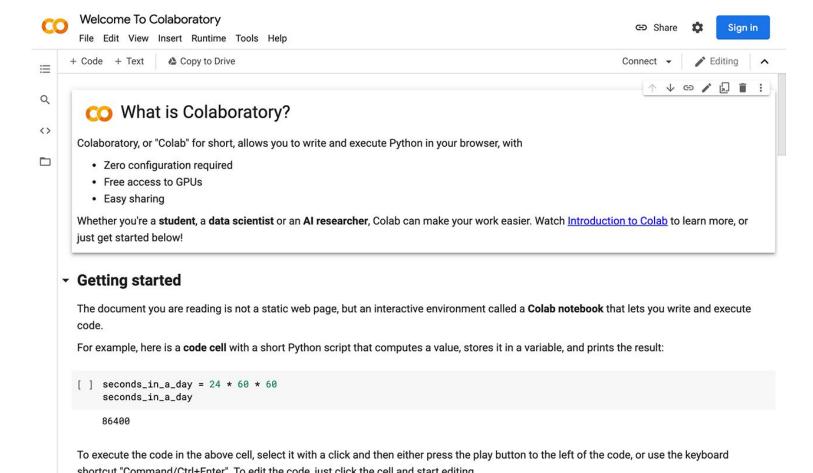




RDM for chemists

Simon Duerr

Google Colab



Q

 $\{x\}$

07

AlphaFold2.ipynb

File Edit View Insert Runtime Tools Help

+ Code + Text Copy to Drive

accessible to all. Nature Methods, 2022

jobname: "test

ColabFold v1.5.5: AlphaFold2 using MMsegs2

Easy to use protein structure and complex prediction using AlphaFold2 and Alphafold2-multimer. Sequence

alignments/templates are generated through MMsegs2 and HHsearch. For more details, see bottom of the notebook, checkout the ColabFold GitHub and read our manuscript. Old versions: v1.4, v1.5.1, v1.5.2, v1.5.3-patch

Mirdita M, Schütze K, Moriwaki Y, Heo L, Ovchinnikov S, Steinegger M. ColabFold: Making protein folding

Input protein sequence(s), then hit Runtime -> Run all

0

query_sequence: "PIAQIHILEGRSDEQKETLIREVSEAISRSLDAPLTSVRVIITEMAKGHFGIGGELASK

num relax: 0

specify how many of the top ranked structures to relax using amber

template_mode: none

• none = no template information is used. pdb100 = detect templates in pdb100 (see notes). custom - upload and search own templates (PDB or mmCIF format, see notes)

• Use: to specify inter-protein chainbreaks for modeling complexes (supports homo- and hetro-oligomers). For example PI...SK:PI...SK for a homodimer

Show code

Install dependencies

Show code

Simon Duerr

colab frequently updates the environment

not reliable enough for sharing research code

You will need to constantly maintain it

2023-12-18

- . Expanded access to Al coding has arrived in Colab across 175 locales for all tiers of Colab users
- Improvements to display of ML-based inline completions (for eligible Pro/Pro+ users)
- · Started a series of notebooks highlighting Gemini API capabilities
- Enable \(\mathbb{H} / \text{Ctrl+L to select the full line in an editor} \)
- · Fixed bug where we weren't correctly formatting output from multiple execution results
- · Python package upgrades
 - CUDA 11.8 to CUDA 12.2
 - tensorflow 2.14.0 -> 2.15.0
 - tensorboard 2.14.0 -> 2.15.0
 - keras 2.14.0 -> 2.15.0
 - Nvidia drivers 525 105 17 -> 535 104 05
 - tensorflow-gcs-config 2.14.0 -> 2.15.0
 - o bigframes 0.13.0 -> 0.17.0
 - o geemap 0.28.2 -> 0.29.6
 - o pyarrow 9.0.0 -> 10.0.1
 - o google-generativeai 0.2.2 -> 0.3.1
 - o jax 0.4.20 -> 0.4.23
 - o iaxlib 0.4.20 -> 0.4.23
- · Python package inclusions
 - kagglehub 0.1.4
 - o google-cloud-aiplatform 1.38.1

EPFL Packaging





Zenodo

Archive your code as is



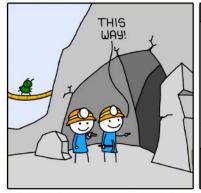
GitHub + Zenoco

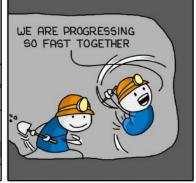
= Citable Code

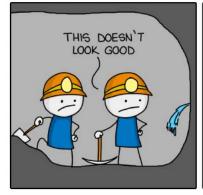
Simon Duerr

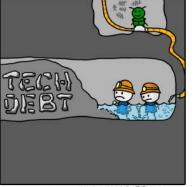
When should you start doing this?

TECH DEBT









MONKEYUSER.COM