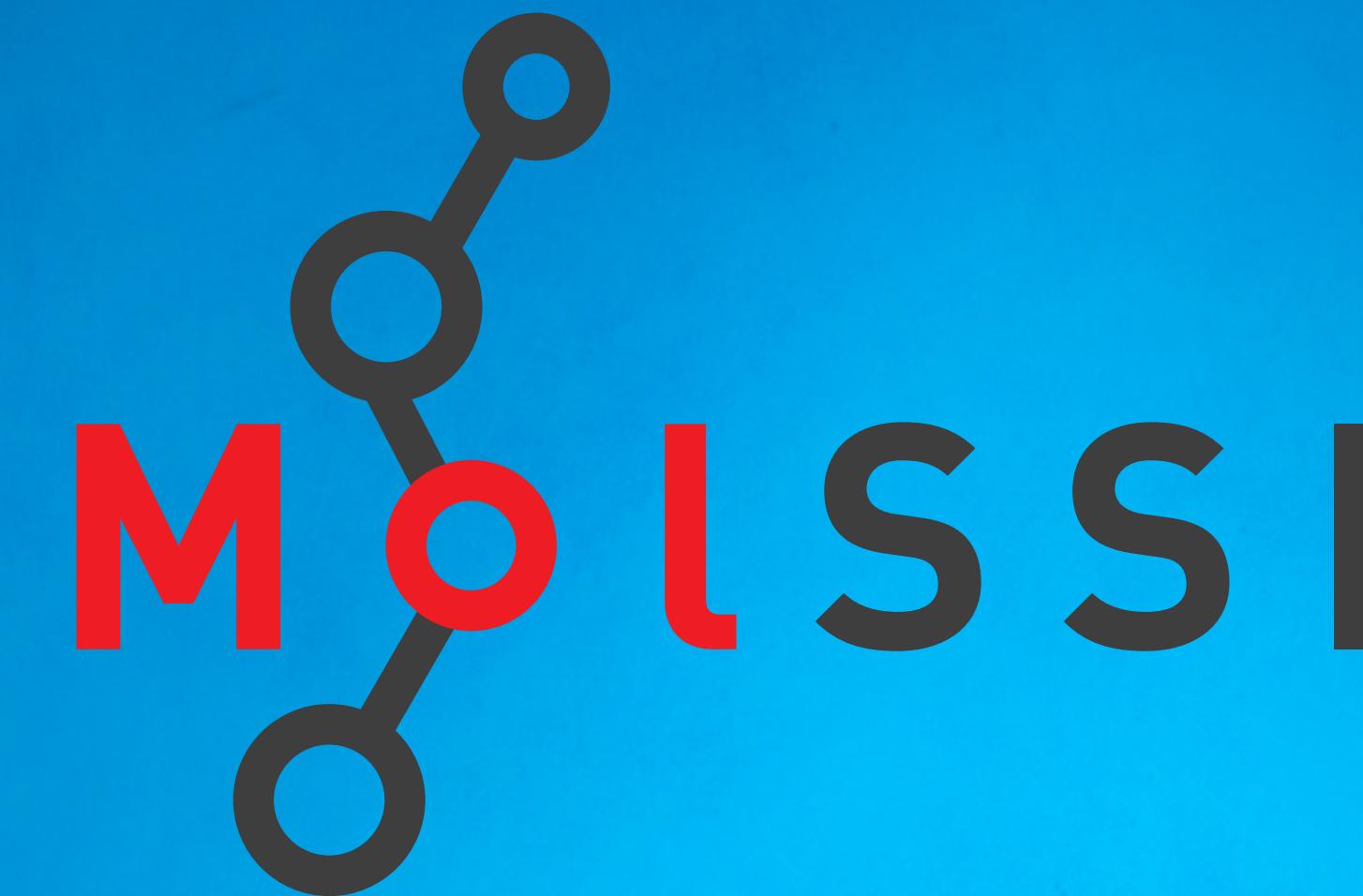


# PSICon 2018



Daniel G. A. Smith  
The Molecular Sciences Software Institute  
@dga\_smith

[molssi.org](http://molssi.org)

# MolSSI Education Initiatives

How do we change the software practices of an entire field?

- Primary objectives:

- Teamwork and collaborative projects
- MolSSI Best Practices
- Open-source paradigms
- Modern tools for modern software
- Teach content to increase both a students scientific capability and industrial marketability
- The “true” cost of ownership



- Inspiration:



- Core Components:

- Software Carpentry style courses
- Learning groups and hands on practice
- Taught by experts in the field

# Motivation

- Many terms thrown around:
  - Best Practice
  - DevOps
  - Open Source
  - Software Sustainability
  - “Incorrect Software”
  - Manage Code
  - Build and Package
  - “Just get the physics working”

“**Best**” is subjective and depends on your own community.

Lets discuss practical applications and how to make our lives easier!

# Open Source

Open Source software is software that can be freely accessed, used, changed, and shared (in modified or unmodified form) by anyone. Open source software is made by many people, and distributed under licenses that comply with the Open Source Definition. – Open Source Initiative



- **Free Redistribution** – Program is free of charge and distributable
- **Source Code** – Source code is freely available
- **Derived Works** – Allows modification of source in derived works

# Open Source

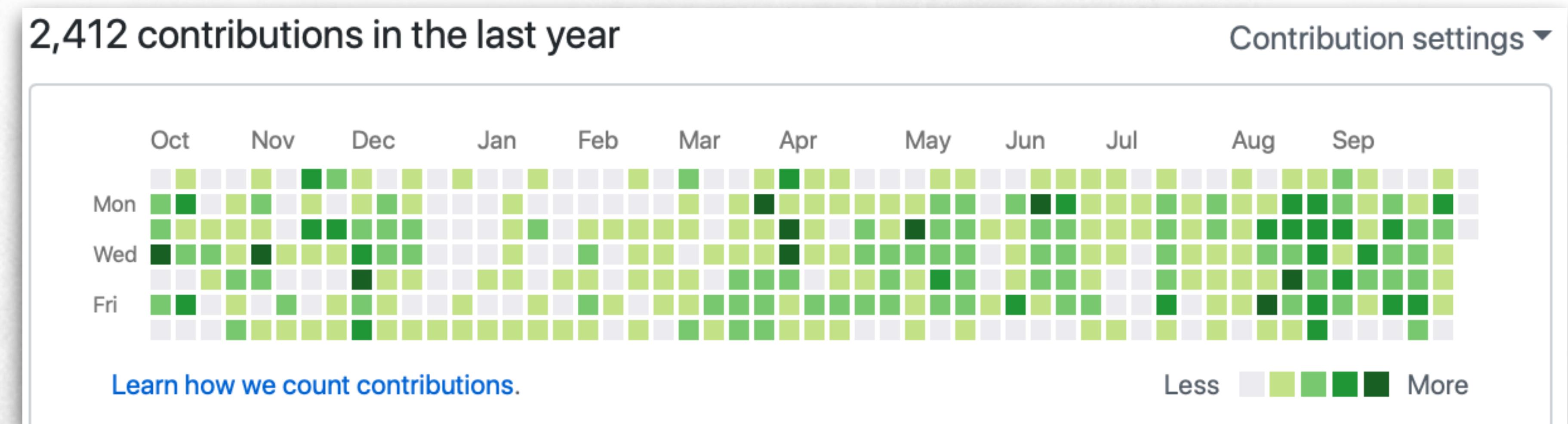
- **Community** - Building engaged cooperative communities
- **Open Collaboration**- Anyone is free to join the community
- **Development Costs** – Allows common goal execution



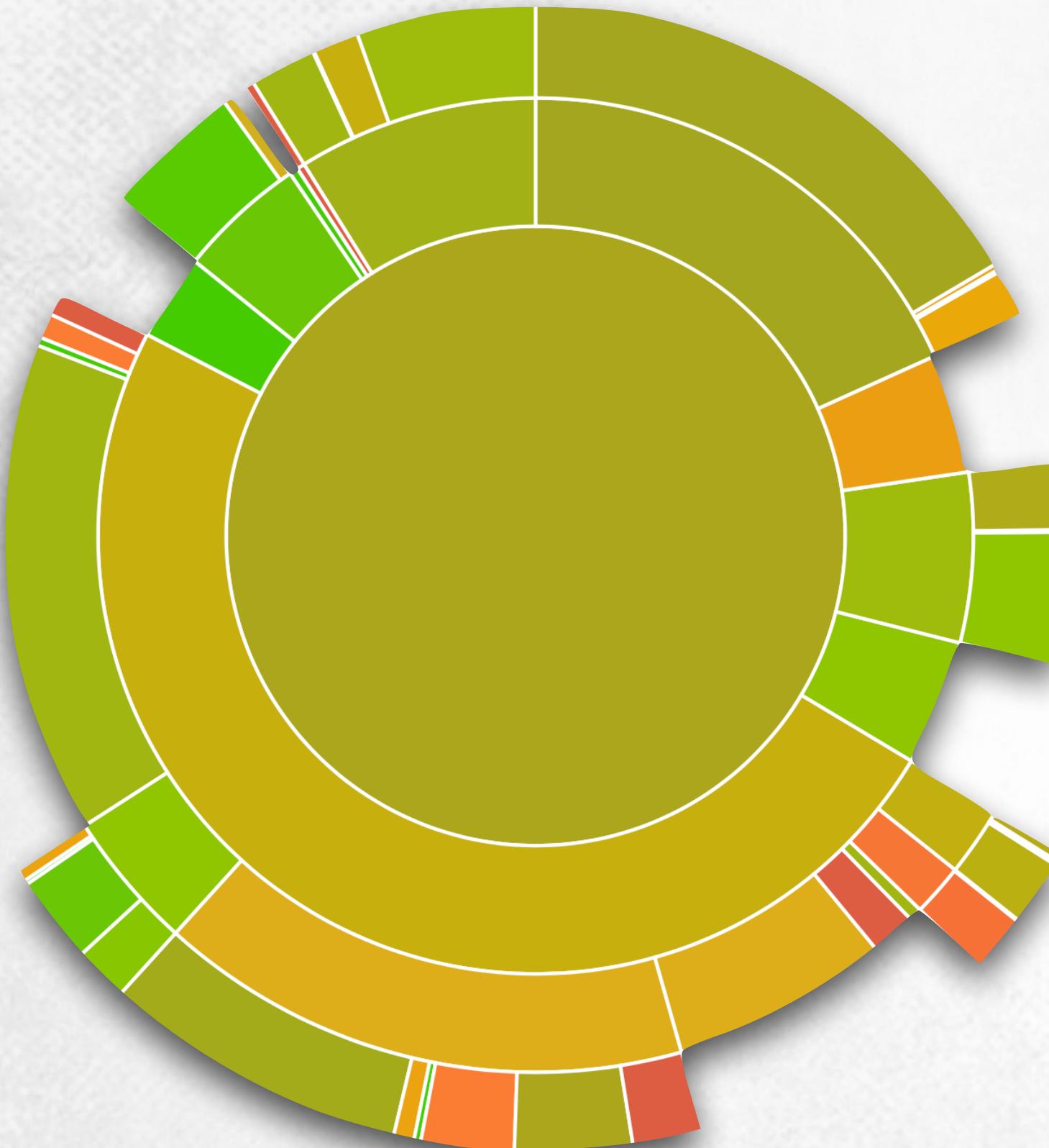
# Open Source

## Personal Reasons

- **Improves Coding Skills** – Code reviews and feedback
- **New Technologies** – Large base of users to provide unique experience
- **Peer Recognition** – Complete portfolio of what a person has done
- **Job Prospects** – Excellent for resumes



# Code Coverage

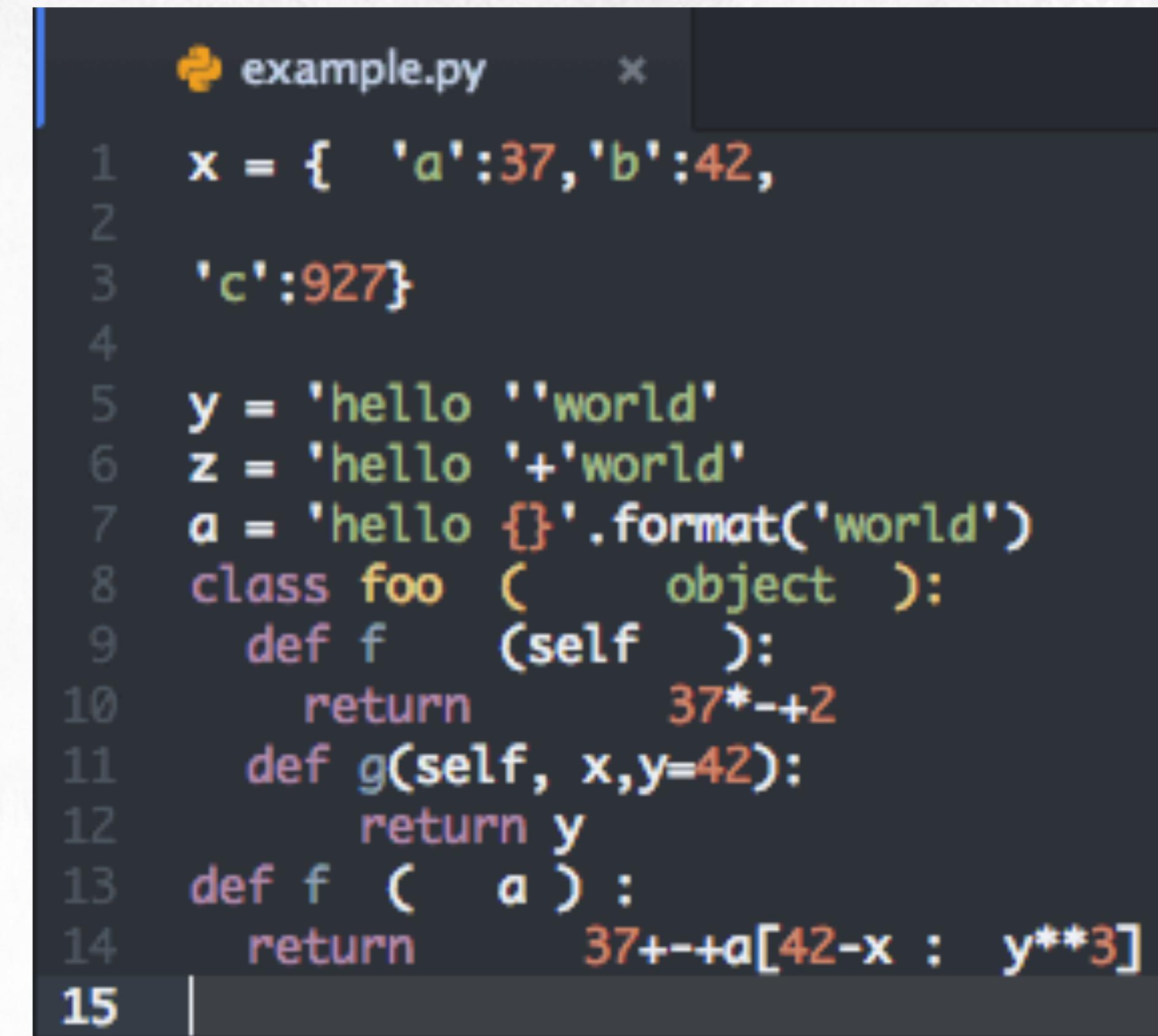


[codecov.io](https://codecov.io)

- Tendency to write many tests that cover a single area while completely neglecting others.
- 100% code coverage does not imply perfect code!
- General metric that informs testing competency that is difficult to fake.

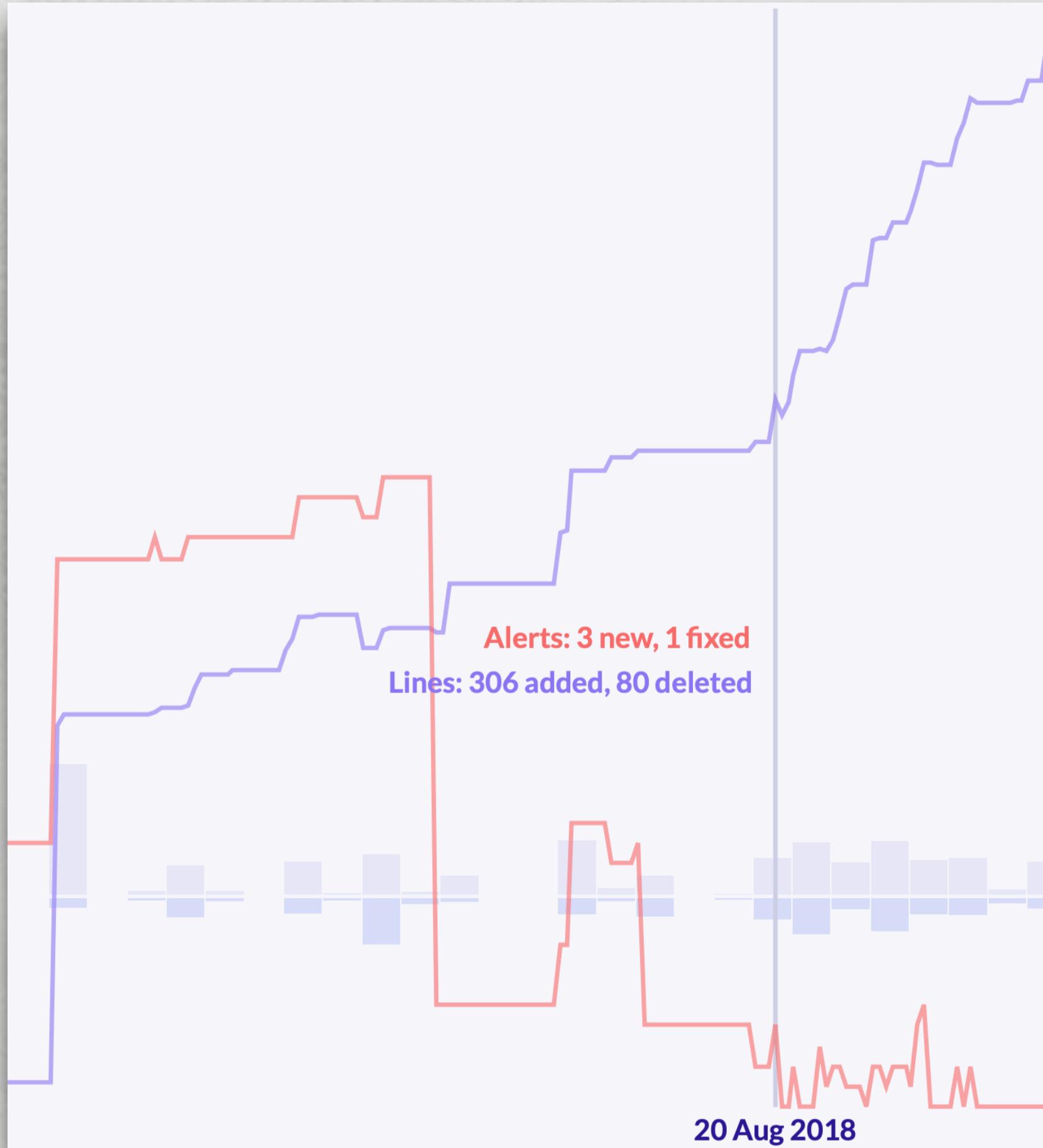
# Code Formatters

- Automatic application of a standard format
  - Most language have a single standard (Python)
  - Several language have multiple standard (C++)
- Provides consistent code
  - Within your own software stack
  - And amongst other software stacks
- Tools: clang-format (C/C++), yapf (Python)
- Most languages (C, C++, Python, Java, Julia, etc)



```
example.py
1 x = { 'a':37,'b':42,
2
3 'c':927}
4
5 y = 'hello ''world'
6 z = 'hello +'world'
7 a = 'hello {}'.format('world')
8 class foo ( object ):
9     def f ( self ):
10         return 37*-+2
11     def g(self, x,y=42):
12         return y
13     def f ( a ):
14         return 37+-+a[42-x : y**3]
15 |
```

# Code Quality



- Statically examine the code for a variety of errors and potential errors
- Catch bugs that do not have unit tests for them.
- Automation of code review:
  - Always catches “known” errors
  - Bad news is easier from a bot

lgtm.co

# Code Quality

## Actual errors

```
2150     for (iop = 0; iop < nop; ++iop) {  
2151         int axis, op_ndim, op_axis;
```

Local variable 'axis' hides a parameter of the same name.



## Possible errors

```
1440     printf("| IterIndex: %d\n", (int)NIT_ITERINDEX(iter));  
1441     printf("| Iterator SizeOf: %d\n",  
1442             (int)NIT_SIZEOF_ITERATOR(itflags, ndim, nop));
```

Multiplication result may overflow 'int' before it is converted to 'unsigned long'.



## Potential code duplication

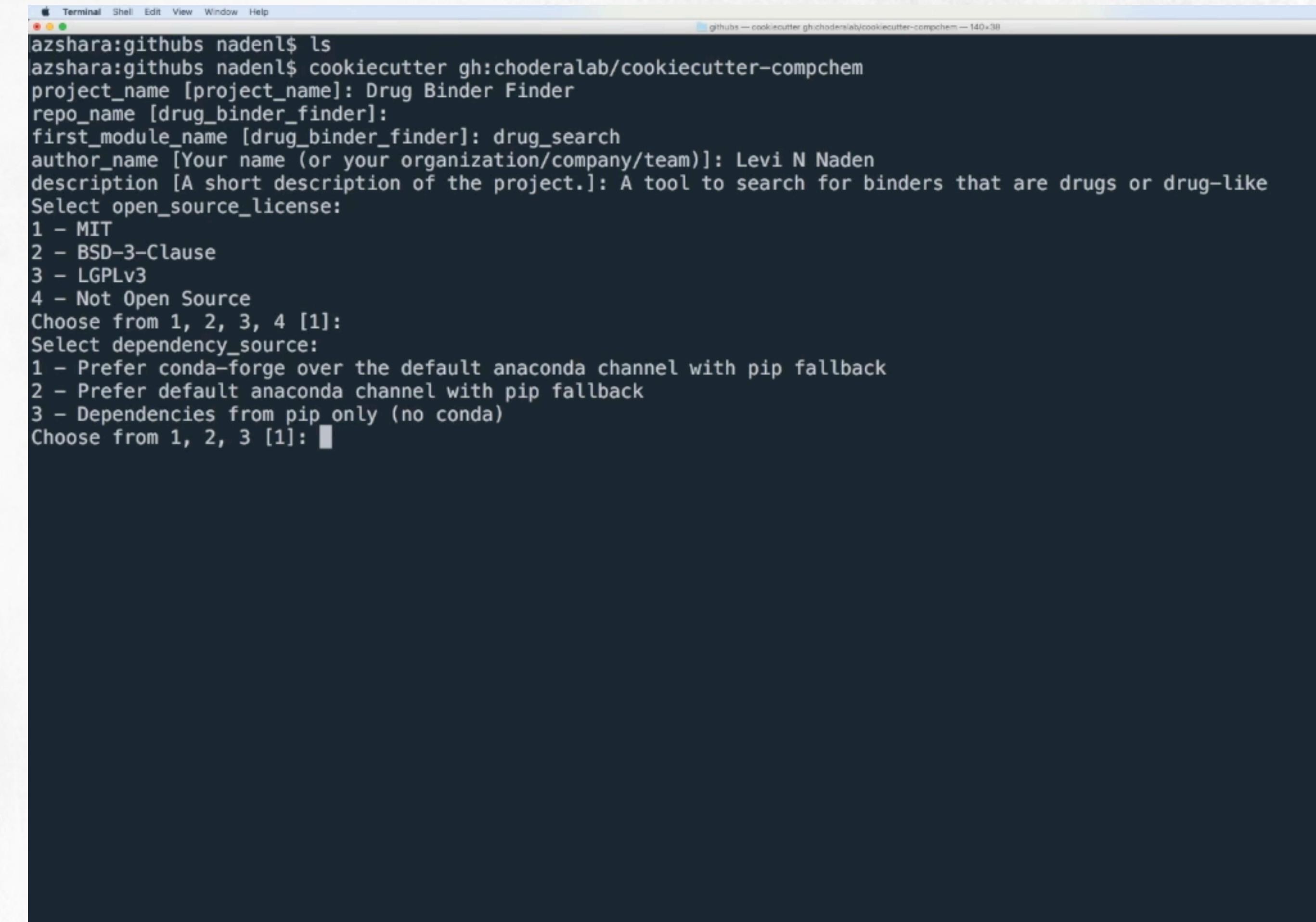
```
186     def list_final_energies(self, fragments=None, refresh_cache=False):
```

All statements in list\_final\_energies are similar in list\_final\_molecules.



# Cookiecutters [Python]

- Automatic setup:
  - Version Control (git)
  - License (MIT/BSD/LGPL/Other)
  - CI (Travis CI/Appveyor)
  - Testing (PyTest)
  - Code Coverage (Codecov)
  - Documentation (Sphinx)
- Takes ~5 minutes to integrate



```
azshara:githubs nadenl$ ls
azshara:githubs nadenl$ cookiecutter gh:choderalab/cookiecutter-compchem
project_name [project_name]: Drug Binder Finder
repo_name [drug_binder_finder]:
first_module_name [drug_binder_finder]: drug_search
author_name [Your name (or your organization/company/team)]: Levi N Naden
description [A short description of the project.]: A tool to search for binders that are drugs or drug-like
Select open_source_license:
1 - MIT
2 - BSD-3-Clause
3 - LGPLv3
4 - Not Open Source
Choose from 1, 2, 3, 4 [1]:
Select dependency_source:
1 - Prefer conda-forge over the default anaconda channel with pip fallback
2 - Prefer default anaconda channel with pip fallback
3 - Dependencies from pip only (no conda)
Choose from 1, 2, 3 [1]:
```

<https://github.com/MoSSI/cookiecutter-cms>

# Final Thoughts

- Make libraries with command line interfaces
  - Think multi-stage runs (workflows)
- Try not to use **stdout**
- If you are writing CUDA you are probably doing it wrong
- Recommend YAML over custom ASCII files
- Recommend JSON/msgpack/YAML/HDF5 for cross platform serialization
- Use the cookiecutter

<https://molssi.org/education/best-practices/>

# MolSSI Integral Reference Project

- Reference implementation and values
- Utilizes arbitrary-precision interval arithmetic (ball arithmetic)
- **Very slow**, but relatively simple implementation



4.78506540470550297026366517126315309034777632299183246390095  
52057465005515845927490470528135254482526 +/- 4.63e-101

<https://github.com/MolSSI/mirp>

# A New and Improved Basis-Set Exchange (BSE)

- Used by both end users and code developers
- Collaboration with original BSE developers (PNNL/EMSL)
- Improved provenance and reproducibility of calculations via unique identifiers
- Improved curation/reliability of the raw data
- Programmatic access via a public API

The screenshot shows the BASIS SET EXCHANGE version 2.0 interface. At the top, there is a navigation bar with links for Release Notes, Feedback, and Help. Below the navigation is a search bar with the placeholder "search basis sets...". To the right of the search bar is a list of basis sets, with "6-311G" highlighted. Further down, there is a detailed view of the "6-311G" basis set, including its description ("VTZ Valence Triple Zeta: 3 Funct.'s/Valence AO"), last update information, and latest version (0). There are also links for "More Information" and "Citations". On the right side of the page, there is a sidebar titled "Download basis set" with a "Format" dropdown set to "Gaussian94" and an "Optimize General Contractions" checkbox. Below the download section is a "Citation" section containing a link to a reference article.

BASIS SET EXCHANGE ver. 2.0

Release Notes Feedback Help ▾

Total found:

1	H												
3	Li	4	Be										
11	Na	12	Mg										
19	K	20	Ca	21	Sc	22	Ti	23	V	24	Cr	25	Mn
37	Rb	38	Sr	39	Y	40	Zr	41	Nb	42	Mo	43	Tc
55	Cs	56	Ba	72	Hf	73	Ta	74	W	75	Re		
87	Fr	88	Ra	104	Rf	105	Db	106	Sg	107	Bh		
57	La	58	Ce	59	Pr	60	Nd						
89	Ac	90	Th	91	Pa	92	U						

All

3-21G  
4-31G  
5-21G  
6-21G  
6-31++G  
6-31++G\*  
6-31++G\*\*  
6-31+G  
6-31+G\*  
6-31+G\*\*  
**6-311G**  
6-311G\*  
6-311G\*\*  
6-31G  
6-31G\*  
6-31G\*\*  
cc-pV5Z  
cc-pV6Z  
cc-pVDZ  
cc-pVQZ  
cc-pVTZ  
CRENBL ECP  
D-FS ECP

search basis sets...

Basis Set: 6-311G

Description: VTZ Valence Triple Zeta: 3 Funct.'s/Valence AO

Last Updated: Latest Version: 0

[More Information](#) [Citations](#)

Download basis set

Format Gaussian94  Optimize General Contractions

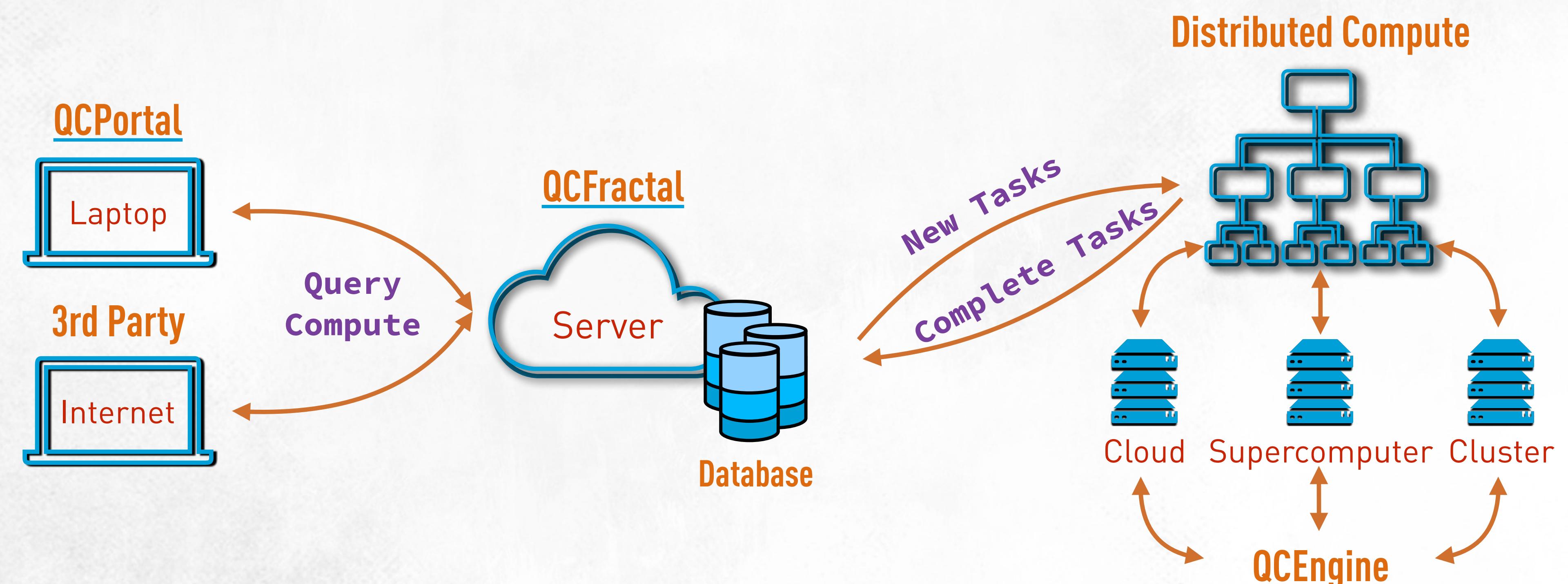
Citation

When publishing results obtained from use of the Basis Set Exchange, please cite:

- The Role of Databases in Support of Computational Chemistry
- Basis Set Exchange: A Community Database for Computational Chemistry, J., and Windus, T.L. J. Chem. Inf. Model., 47(3), 1045-1052,

# QCArchive Overview

Provide an open, community-wide quantum chemistry database to both facilitate and capture hundreds of millions of hours of computing time to enable large-scale forcefield construction, physical property prediction, new methodology assessment, and machine learning from data that would otherwise end up siloed or inaccessible.

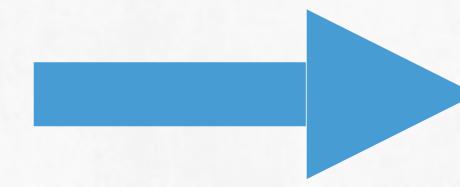


# Stack: QC Schema

[https://github.com/MoSSI/QC\\_JSON\\_Schema](https://github.com/MoSSI/QC_JSON_Schema)

- Communication channel between all piece of the ecosystem.
- *Community* project useful for many aspects of quantum chemistry.
- Not only JSON, but any key/value/array language (BSON/HDF5/XML/YAML)
- Molecule
- Input
- Output
- Optimization Trajectory

```
{  
  "molecule": {  
    "geometry": [0, 0, 0, 0, 0, 1],  
    "atoms": ["He", "He"]  
  },  
  "driver": "energy",  
  "model": {  
    "method": "SCF",  
    "basis": "sto-3g",  
  },  
}
```



```
{  
  ...Input  
  "provenance": {  
    "creator": "My QM Program",  
    "version": "1.1rc1",  
  },  
  "properties": {  
    "scf_n_iterations": 2.0,  
    "scf_total_energy": -5.433191881443323,  
    "nuclear_repulsion_energy": 2.11670883436,  
    "one_electron_energy": -11.67399006298957,  
    ...  
  },  
  "error": "",  
  "success": true,  
  "raw_output": "Output storing was not requested."  
}
```

# Development Efforts

## ► Methods:

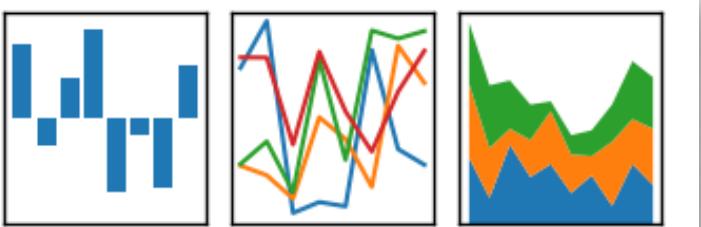
Collections: Datasets, OpenFFWorkflow

Services: Torsiondrives, geometry optimizations

Computation: Psi4, RDKit, TorchANI

## Front End

pandas  
 $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$



## ► Format:

A Quantum Chemistry schema so that different QC programs, AI evaluators, and forcefields can be plugged into the backend.

## ► Front-end:

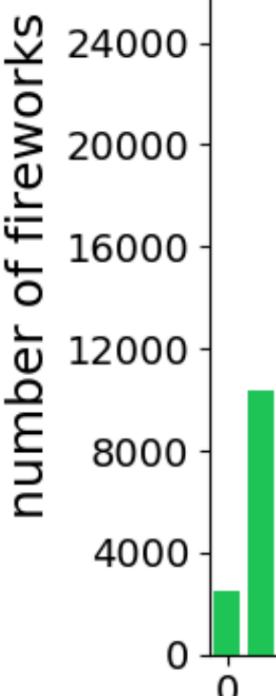
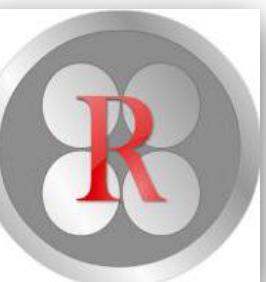
Python front end for querying, computing, and visualization. REST API for power users and web front ends.

## Database Technologies



Apache Arrow

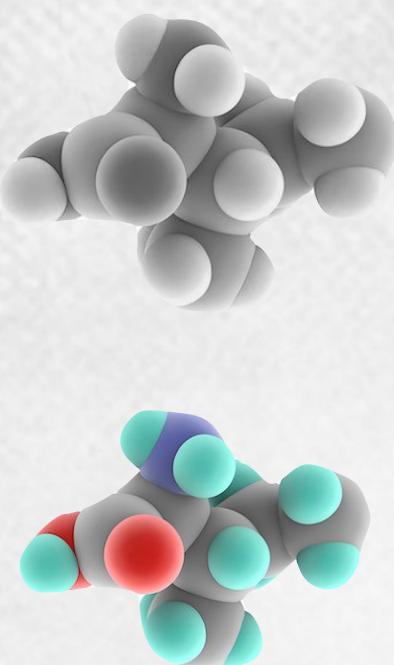
## Distributed Compute Engines



# Production Runs

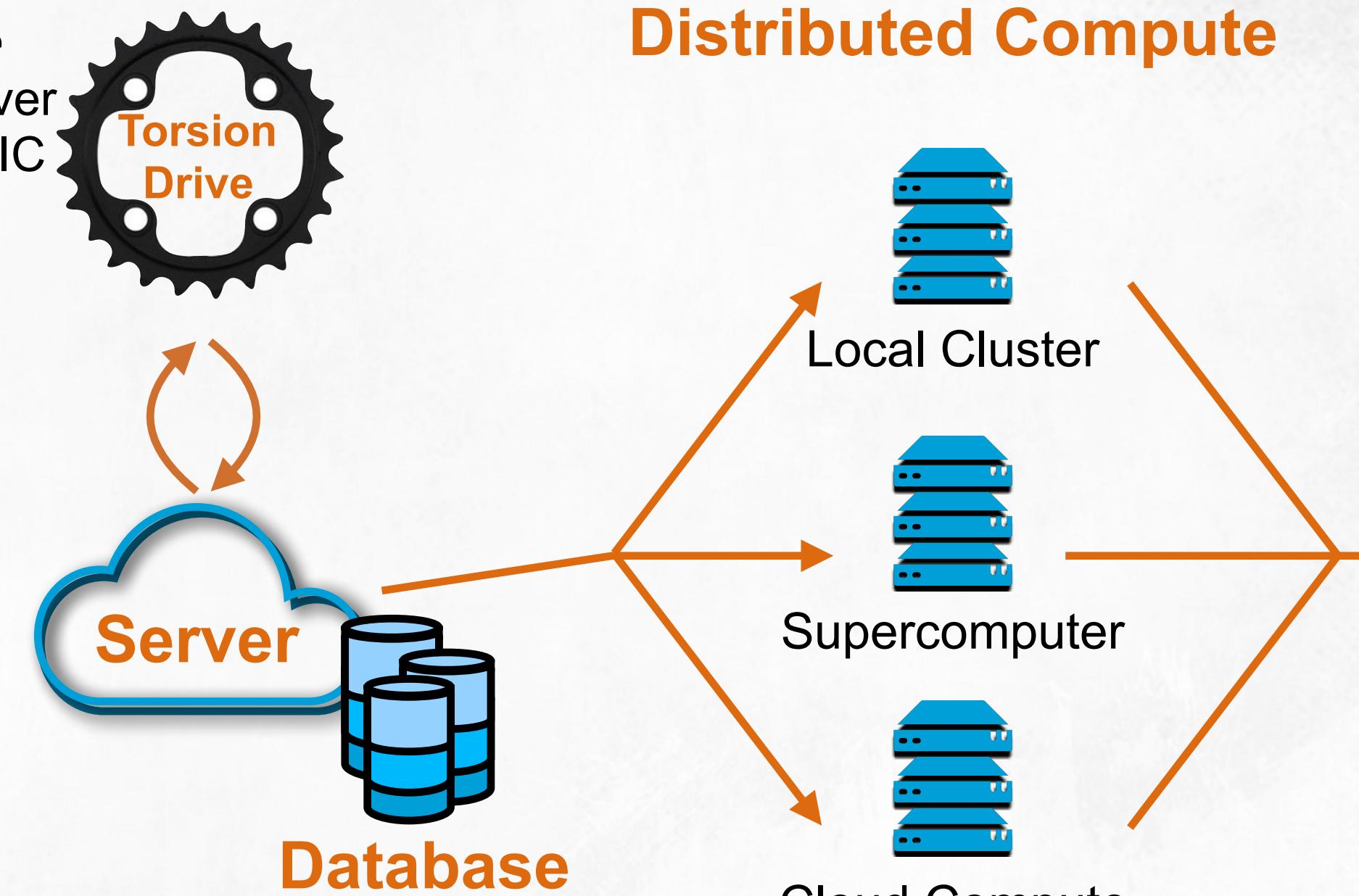
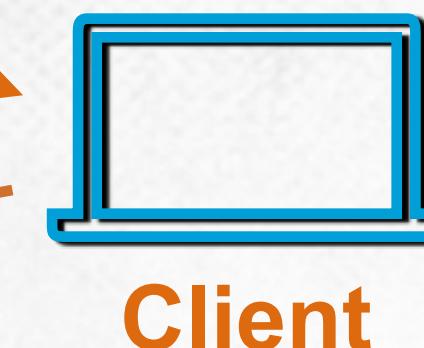
## Molecular Torsion Drives

A. The user runs a client that requests quantum chemical computations for lists of molecules.



B. Clients request results from the server. These are either returned from storage in the server or new computations are initiated.

D. The iterative *torsiondrive* procedure is run on the server and spawns new geomeTRIC computations as previous iterations complete.

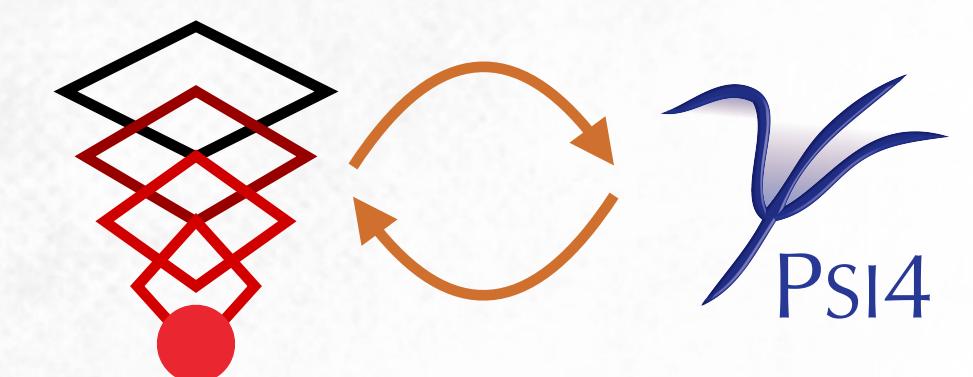


C. A central sever stores a record of complete computations and can execute “services” such as *torsiondrive*.

H. The molecules with requested quantum chemical data are returned to the user.

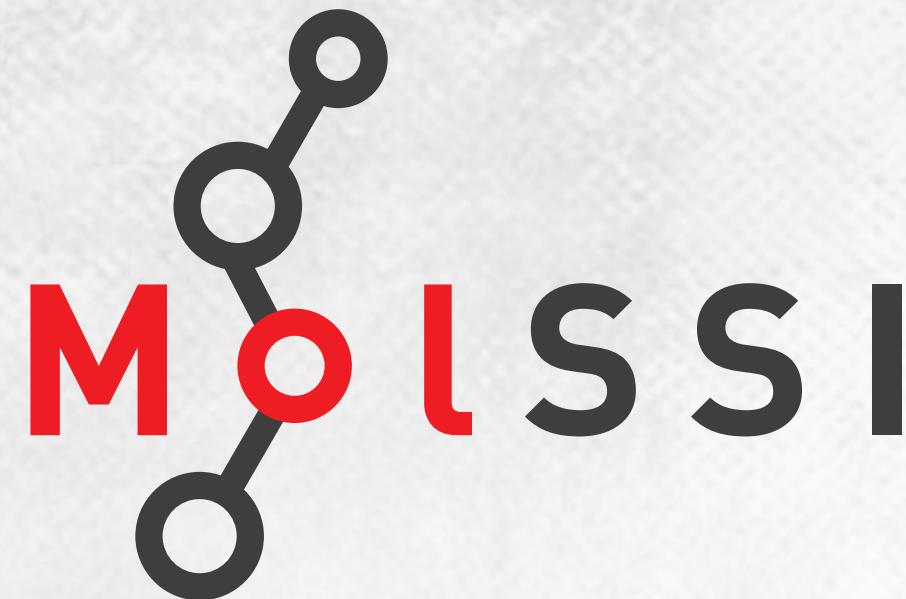
E. Compute clients can be run on any resource connected to the internet. These clients request computations and return results.

G. The results of the geomeTRIC / Psi4 quantum computations are formatted via the QC Schema and returned to the database.



## Geometry Optimization

# Challenges and Directions



## ► Challenges:

- Intricate software stack, focusing on simplifying user experience and interaction.
- Deployment of complex programs on large scale compute resources.
- Project is ~6 months old at the moment, still early days!

## ► Future Directions:

- A central database to hold not only Open Force Field data, but data from many molecular science subfields with a web-portal front-end.
- Enhanced “workflows” or chained quantum chemistry computations.
- Array-based storage for storing arrays like density matrices.