



Using an Open Source Python Toolbox (Signac) to Manage High Dimensional Research Data

Chrysy Xiyu Du



Mechanical Engineering



My background

Assistant Professor in Mechanical Engineering (Jan. 2023)

PhD in the Glotzer Group at University of Michigan

Highly computation driven group with many open-source softwares

<https://github.com/glotzerlab>

Signac



Python based data-management software

Completely Open source

Can be installed with conda and pip

Slack community

Documentation: <https://docs.signac.io/en/latest/>

Why do I need it?

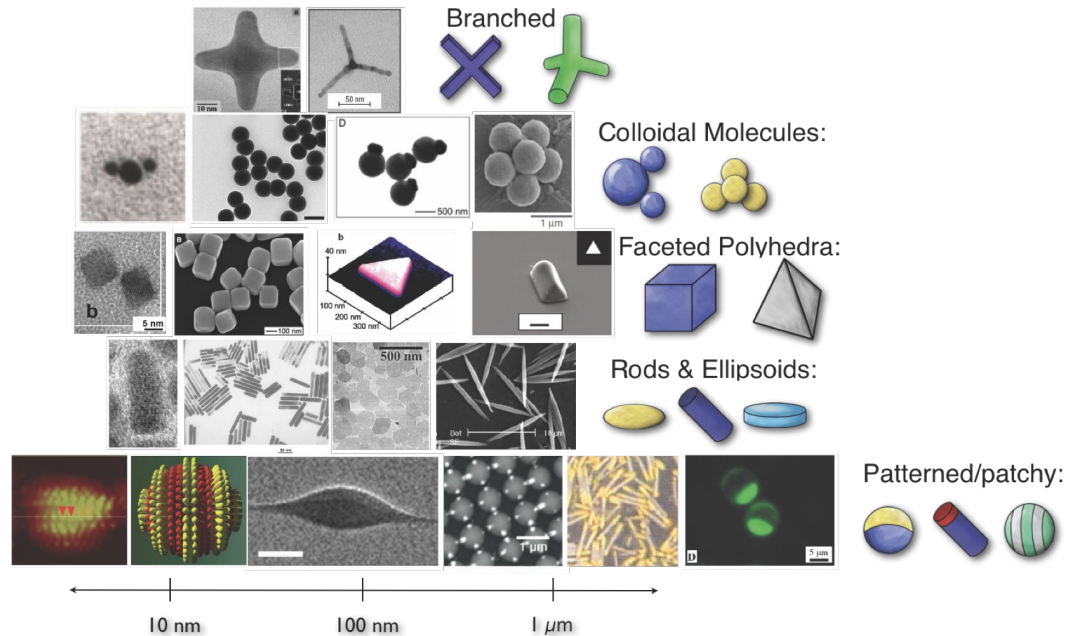
My research: use theory and computer simulation to design the next generation of soft materials.

Why do I need it?



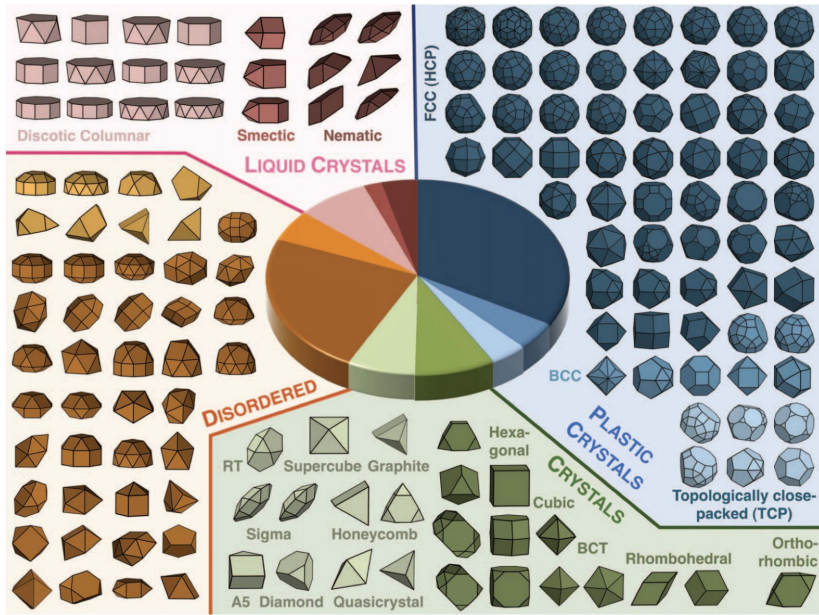
Movie: Big Hero 6 (gif by [aviscranio.tumblr.com](https://www.tumblr.com/aviscranio))

Large parameter space

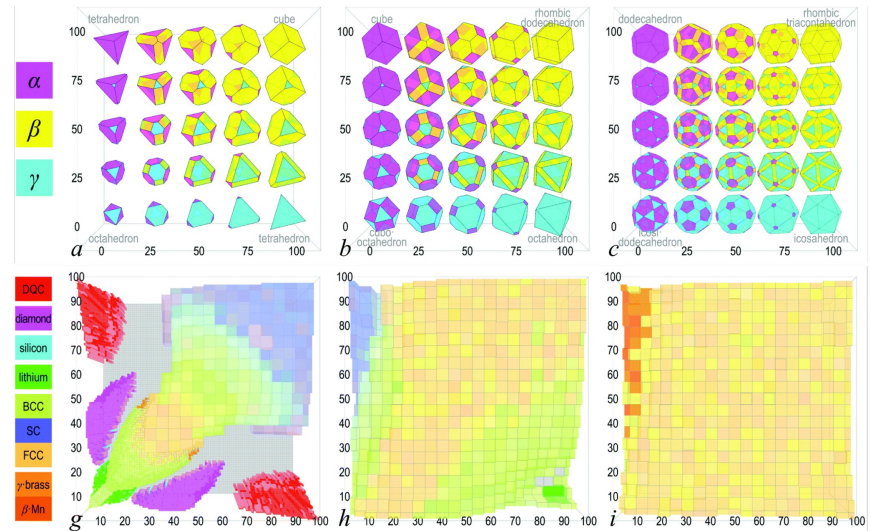


S. C. Glotzer, M. J. Solomon (2007)

Large parameter space

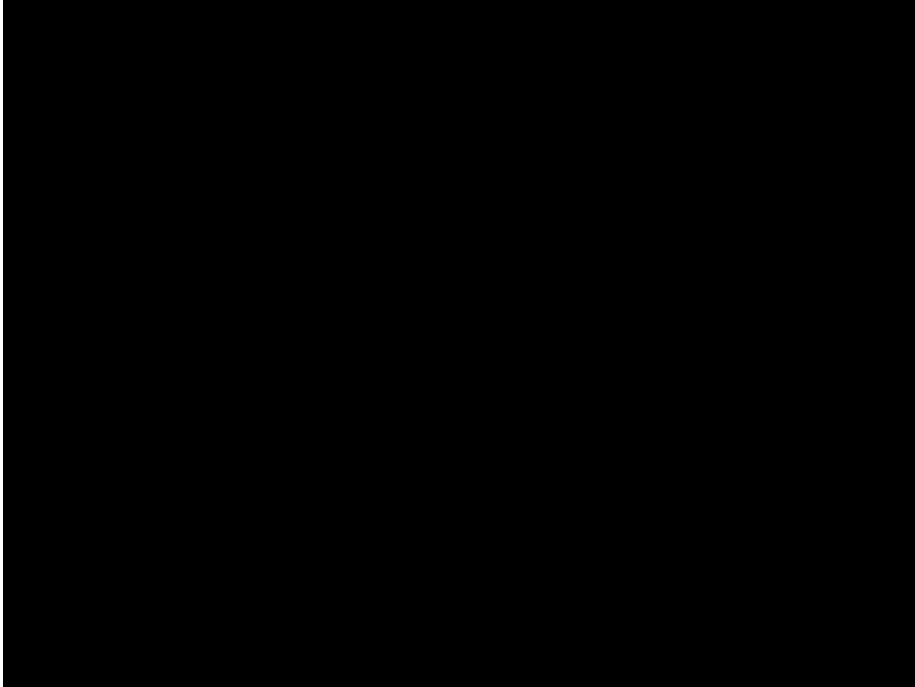


P. F. Damasceno, M. Engel, S. C. Glotzer, Science (2012)



D. Klotsa, E. R. Chen, M. Engel, and S. C. Glotzer, Soft Matter (2018)

What kind of data/metadata?



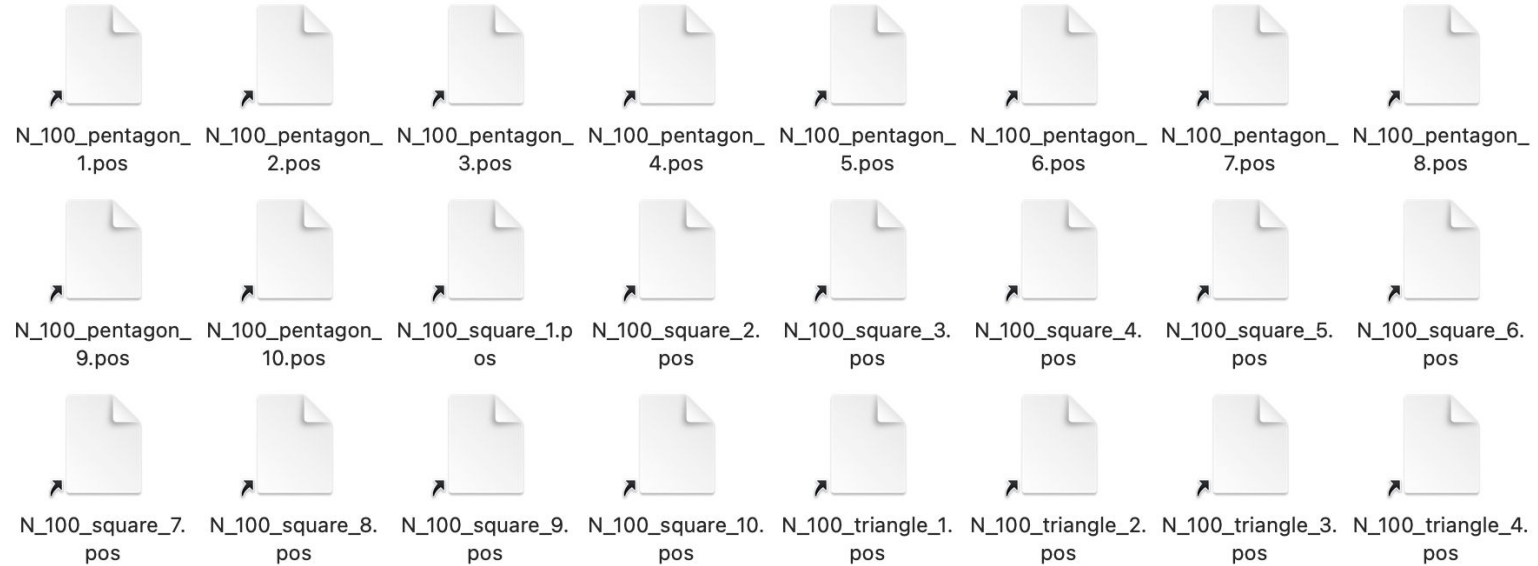
Metadata:

Number of particles, system density, temperature, how many steps, interaction type, etc

Data:

Positions of every particle, measured system information (energy, pressure, etc)

Before signac...



Metadata is a part of the file name for the data file.

Not sustainable

Hard to parse when doing batch analysis

Not robust to parameter expansion

Naming convention needs human knowledge

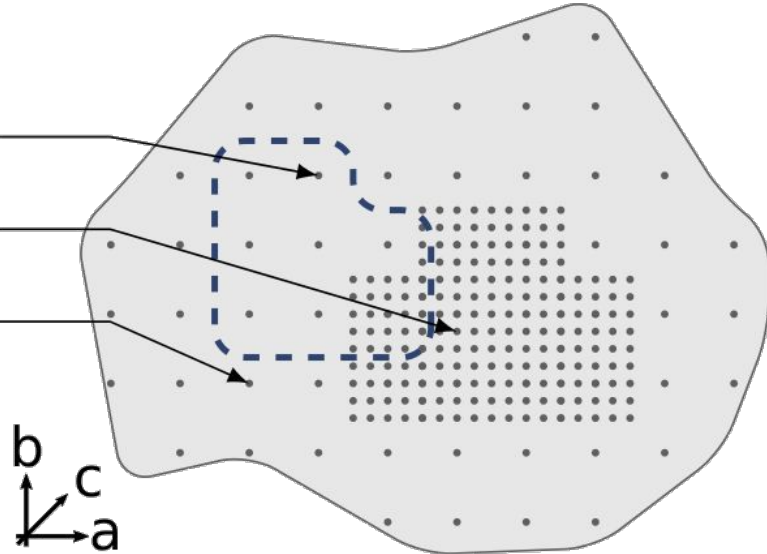
Hard to share data with others

Can we do better?

Introducing Signac

Index
$\{a = 3.0, b = 5.00, c = x\}$
$\{a = 5.0, b = 2.75, c = x\}$
$\{a = 2.0, b = 2.00, c = x\}$
\vdots

Parameter & Data Space



Signac workspace



1fb813e07adf1a3
3517e7e...87cf4c0



2bc8c88ca0659d
09edfe0...073bdb



3ccf52224d9053
b96c68...60b6a19



4c8e3ebb5eb217
3952c5...0ab4413



5cc7b18c2f7ee71
a70a5d...3b929cf



9a343bd27b5ab9
fb9d3f1...a11d5c9



9c9c01be0e9839
519e81...b73dfdf



9c12279f0c6ca2
81a57b...6fef9862



14b3d4075439a8
117c859...2b96e4



22d0b46b713682
6c409c...de94b8a



899e66ab3076b7
341e94...9cbedb9



b47c8fcc2aa49ac
10d78ff...78ef102



bb470d7f4444e3
383926...3eb9a31



c97d278d4b6988
58ec7e...248e468



cde2aa99a73bb7f
2cce9f3...a9e418



dee000775878c4
9eda09...231c43ff



grad0.1.txt



grad0.01.txt



grad0.05.txt



loss0.1.txt



loss0.01.txt



loss0.05.txt



params0.1.txt



params0.01.txt



params0.05.txt



random_params.n
py



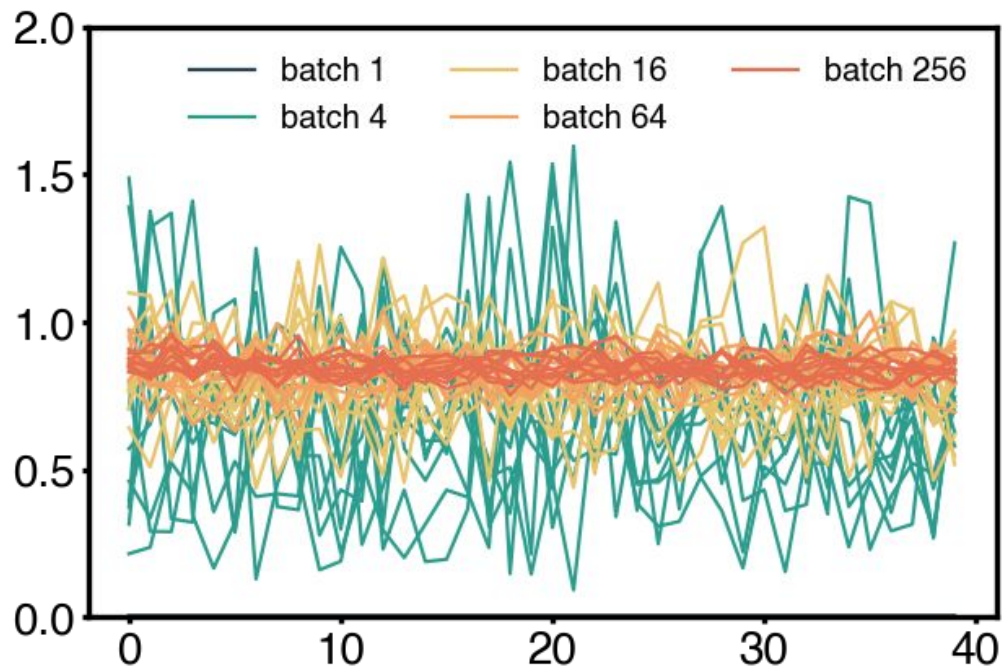
signac_statepoint
.json

Signac statepoint

```
{"N": 6, "s_type": "octahedra", "num_patch": 10, "type_patch": 2,  
"r": 2.5, "phi": 0.05, "num_steps": 40000, "scale": 1.3, "kT": 0.8,  
"D0": 4.0, "dt": 0.0001, "rand_search_count": 50, "batch_size": 16,  
"loop_batch": 4, "num_steps_opt": 1000, "opt_steps": 100,  
"learning_rate": 0.1, "closeness_penalty": 0.0,  
"closeness_penalty_nbrs": 1, "seed": 86711402, "replica": 3}
```

These metadata are being encoded into the hex string of the statepoint folders for quick search.

Post processing



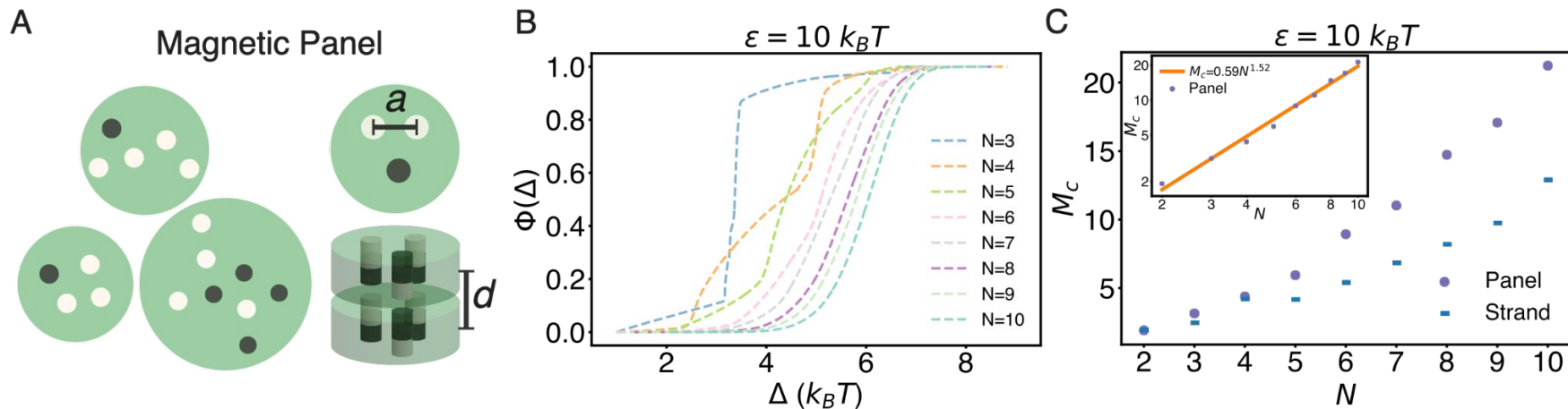
Post processing

```
colors = ['#264653', '#2a9d8f', '#e9c46a', '#f4a261', '#e76f51']
batch_size = [1, 4, 16, 64, 256]

fig = plt.figure(num = 1, figsize = (9, 6), dpi = 80, facecolor = None, edgecolor = 'k')

for i in range(len(batch_size)):
    for job in project.find_jobs({'batch_size': batch_size[i]}):
        ave_loss = np.load(job.fn("std_loss.npy"))
        if job.sp.replica == 1:
            plt.plot(ave_loss, color = colors[i], linewidth = 2, label = 'batch '+str(batch_size[i]))
        else:
            plt.plot(ave_loss, color = colors[i], linewidth = 2)
plt.legend(frameon = False, fontsize = 18, ncol=3)
plt.ylim([0, 2.0])
```


My publication using Signac



Over 1500 statepoint folders, other project has over 10k statepoints

Signac



Python based data-management software

Completely Open source

Can be installed with conda and pip

Slack community

Documentation: <https://docs.signac.io/en/latest/>