



Reproducible computational workflows with **signac**

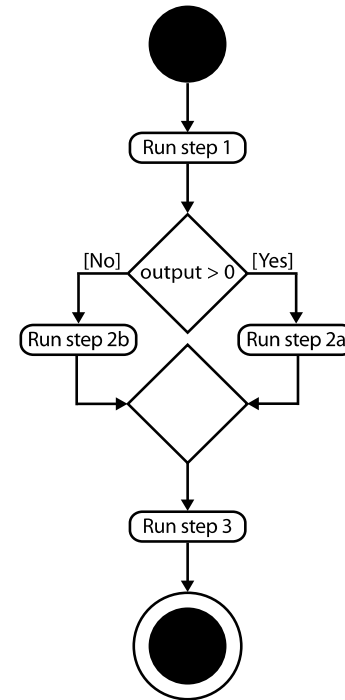
Bradley D. Dice, Carl S. Adorf, Vyas Ramasubramani, Sharon C. Glotzer

MICDE Symposium for the Center for Network and Storage Enabled
Collaborative Computational Science

October 15, 2018, Ann Arbor, MI

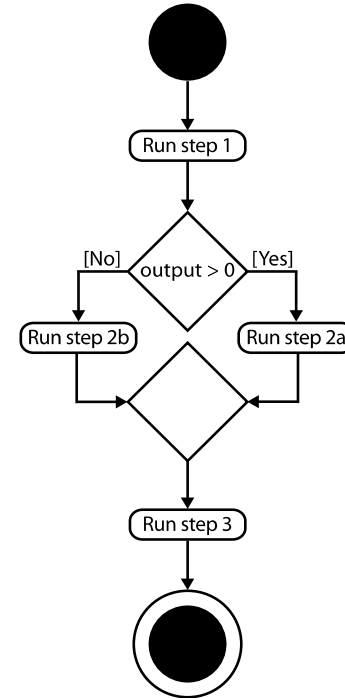
Planning a Parameter Study

concentration_A_0.25/
concentration_A_0.50/
concentration_A_0.75/



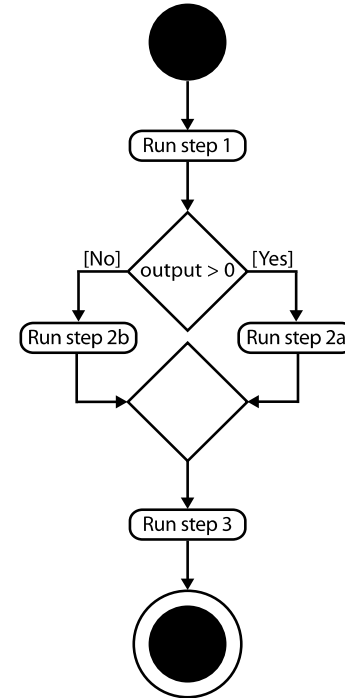
Planning a Parameter Study

concentration_A/0.25
concentration_A/0.50
concentration_A/0.75



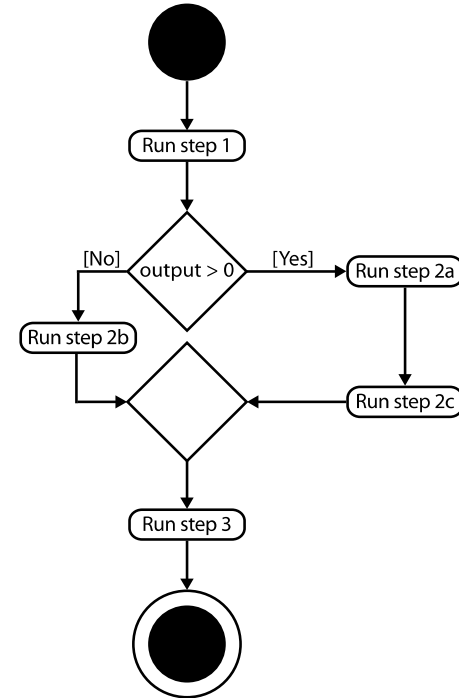
Planning a Parameter Study

conc_A/0.25
conc_A/0.50
conc_A/0.75



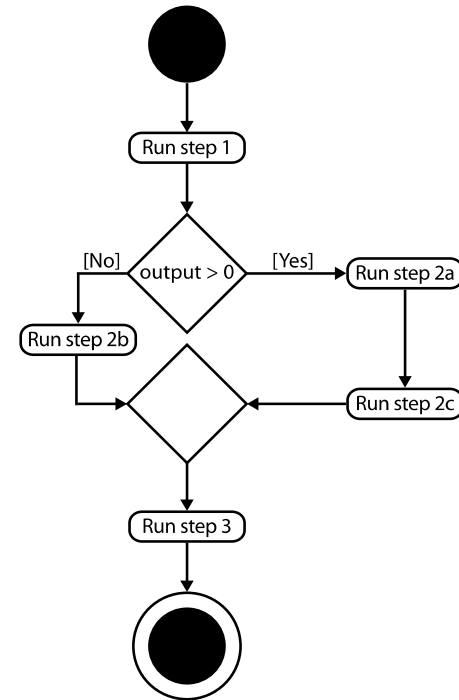
Planning a Parameter Study

conc_A/0.25/temp_08
conc_A/0.25/temp_1
conc_A/0.50
conc_A/0.75



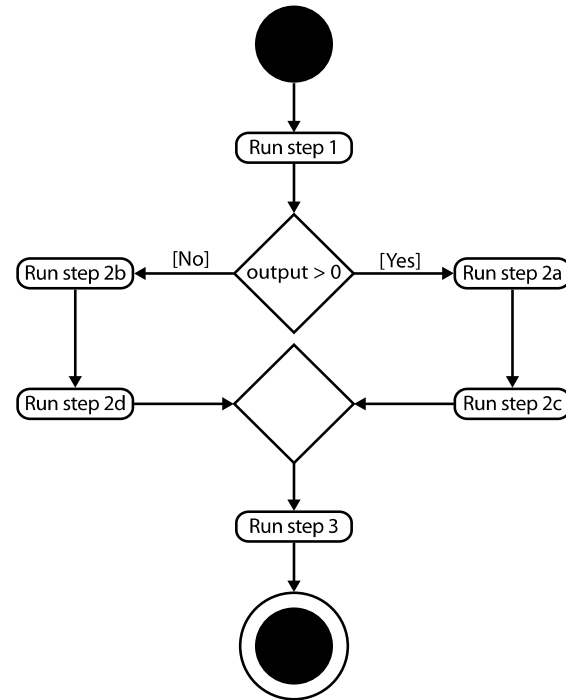
Planning a Parameter Study

temp_08/conc_A/0.25
temp_08/conc_A/0.50
temp_08/conc_A/0.75
temp_1/conc_A/0.25



Planning a Parameter Study

temp_08/conc_A/25/conc_B/05
temp_08/conc_A/50
temp_08/conc_A/75
temp_1/conc_A/25

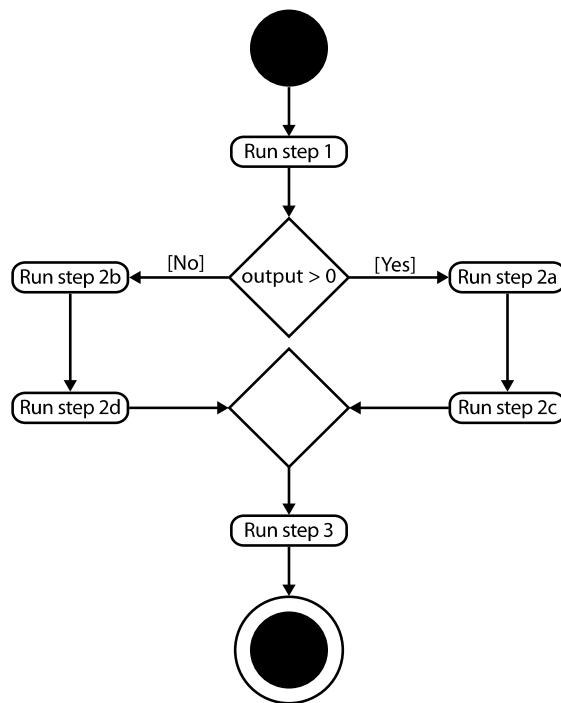


Planning a Parameter Study

```
#!/bin/bash
#SBATCH -J myproject
#SBATCH -N 16
#SBATCH -A ${MYACCOUNT}
#SBATCH -p ${QUEUE}
#SBATCH --ntasks-per-node 8
#SBATCH -t 12:00:00
```

```
cd ${WORKING_DIRECTORY}
```

```
mpirun -n 16 ./myscript.sh
```



The **signac** Framework



A lightweight, application-agnostic software framework that unobtrusively helps users manage and scale file-based workflows, facilitating data reuse, sharing, and reproducibility.

Python 2/3 | open source BSD-3 | install with pip or conda

www.signac.io

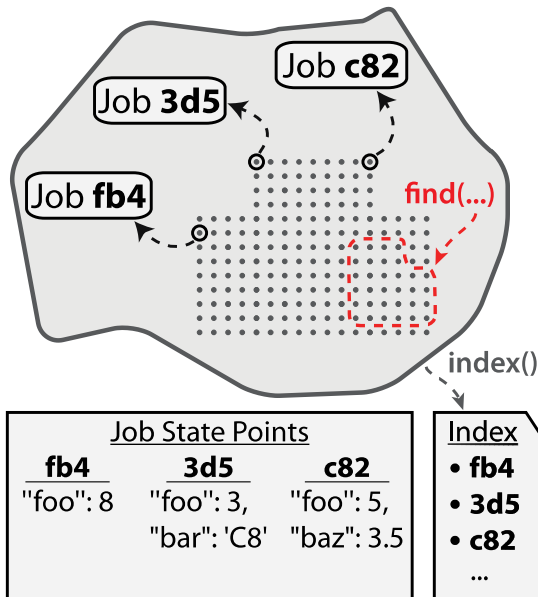
Topic Overview

1. Introduction
2. Projectile Demo
3. Recent Development and Future Goals

Data (**signac**) and Workflow (**signac-flow**) Management

signac

(a) Active Workspace

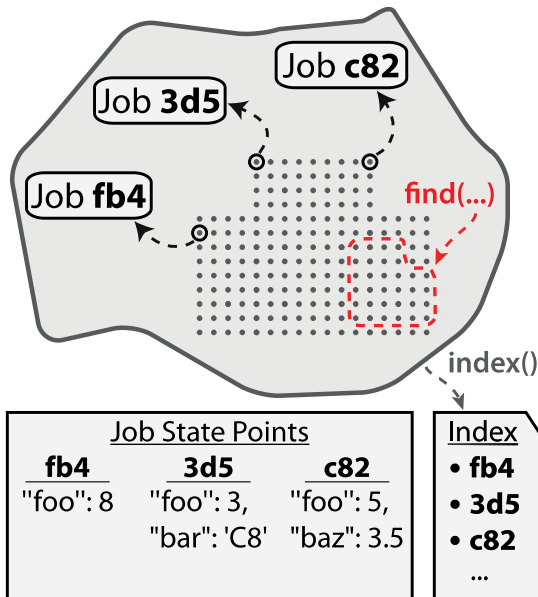


Data (**signac**) and Workflow (**signac-flow**) Management

signac

(a) Active Workspace

The workspace consists of **jobs**, data containers associated with distinct metadata mappings (called **state points**).

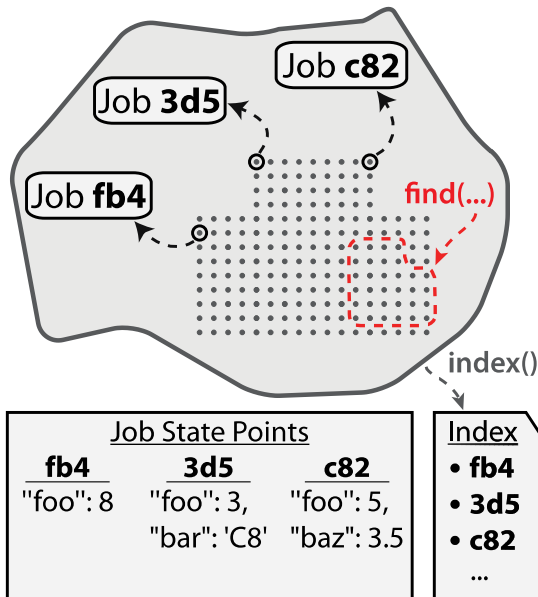


Data (**signac**) and Workflow (**signac-flow**) Management

signac

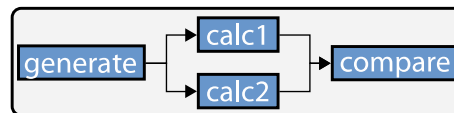
(a) Active Workspace

The workspace consists of **jobs**, data containers associated with distinct metadata mappings (called **state points**).



signac-flow

(b) Project Workflow



(c) Status Tracking

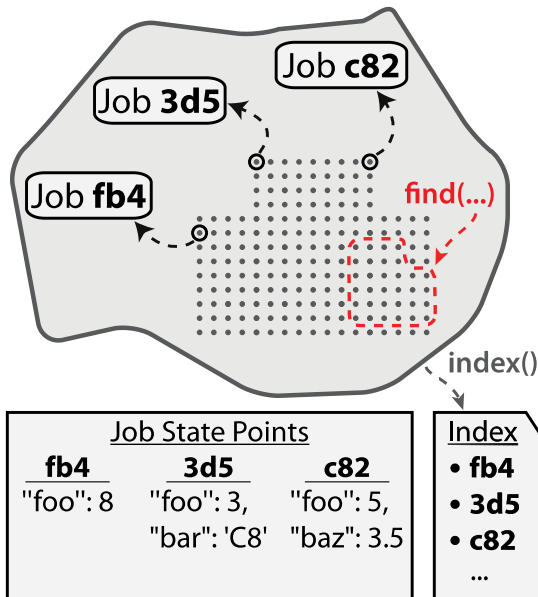
operation(job)	Status
generate(fb4)	✓
calc1(fb4)	✓
calc2(fb4)	✓
compare(fb4)	➔
generate(3d5)	✓
calc1(3d5)	✓
calc2(3d5)	➔
compare(3d5)	✗
...	...

Data (**signac**) and Workflow (**signac-flow**) Management

signac

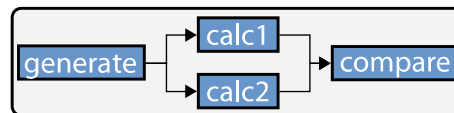
(a) Active Workspace

The workspace consists of **jobs**, data containers associated with distinct metadata mappings (called **state points**).



signac-flow

(b) Project Workflow



(c) Status Tracking

operation(job)	Status
generate(fb4)	✓
calc1(fb4)	✓
calc2(fb4)	✓
compare(fb4)	➔
generate(3d5)	✓
calc1(3d5)	✓
calc2(3d5)	➔
compare(3d5)	✗
...	...

The **workflow** consists of **operations**, linked through condition functions. This forms a directed graph.

The **signac** Team



Carl S. Adorf
Lead developer &
co-maintainer,
signac and signac-flow



Vyas Ramasubramani
Developer &
co-maintainer,
signac and signac-flow



Bradley Dice
Lead developer &
maintainer,
signac-dashboard

Top 6 Features Released In Last 12 Months

Advanced Searching and Aggregation



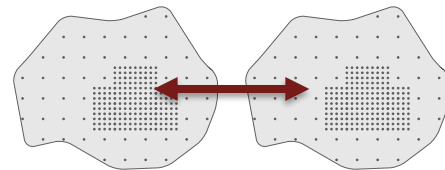
```
find({"T.$gt": 298})  
.groupby("P")
```

Schema Detection



```
'T': int([298, 300, 302])  
'P': float([0.1, 1.0])
```

Synchronization



One-script Projects



*Implement
workflows in
< 10 lines.*

Improved Container Support



Templating



- container by DPIcons from the Noun Project
- Relational Schema by Becris from the Noun Project

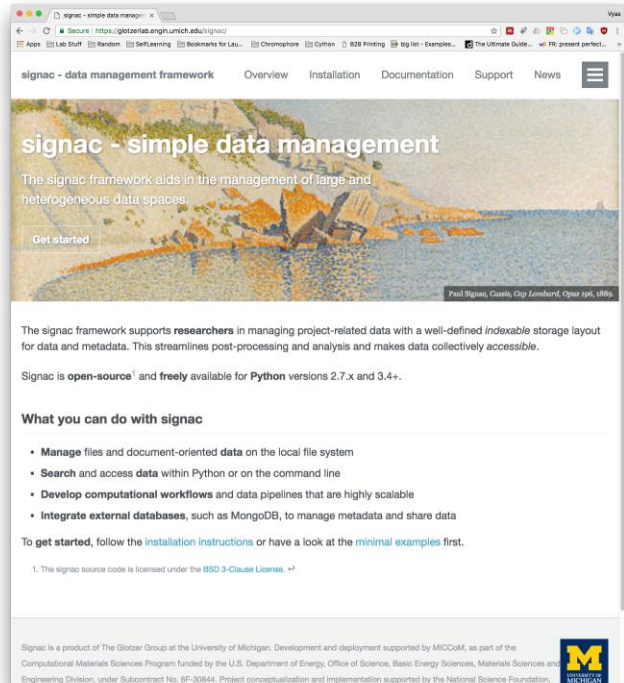
- search data by SBTS from the Noun Project
- Workflow by ProSymbols from the Noun Project

News, Documentation & Chat Room Support

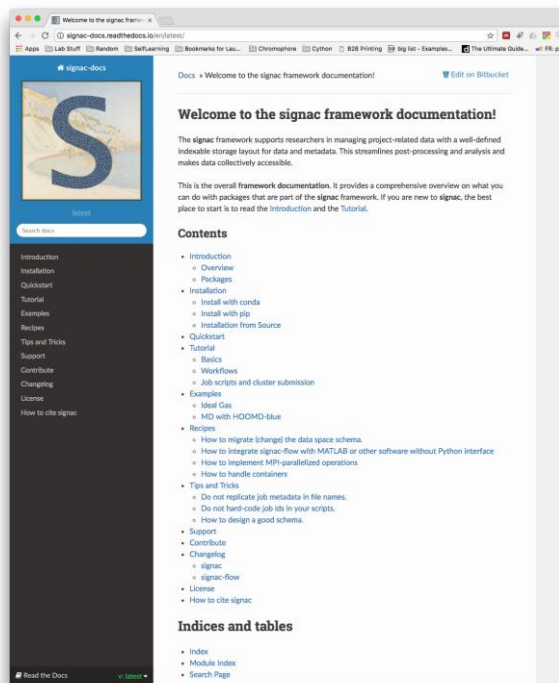
<http://www.signac.io>

<http://signac-docs.readthedocs.io>

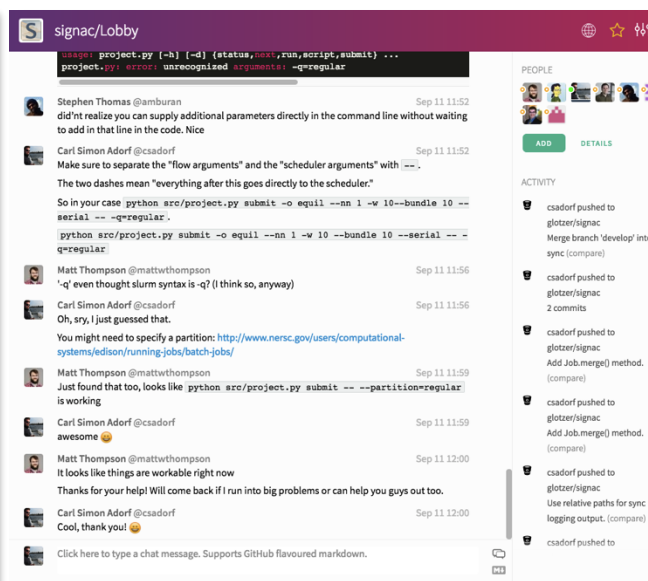
<https://gitter.im/signac/Lobby>



The screenshot shows the signac website homepage. At the top, there's a navigation bar with links for Overview, Installation, Documentation, Support, and News. The main heading is "signac - simple data management". Below it, a paragraph states: "The signac framework aids in the management of large and heterogeneous data spaces." A "Get started" button is visible. Further down, another paragraph explains: "The signac framework supports researchers in managing project-related data with a well-defined indexable storage layout for data and metadata. This streamlines post-processing and analysis and makes data collectively accessible." Below this, it says "Signac is open-source* and freely available for Python versions 2.7.x and 3.4+." A section titled "What you can do with signac" lists several capabilities: Manage files and document-oriented data on the local file system; Search and access data within Python or on the command line; Develop computational workflows and data pipelines that are highly scalable; Integrate external databases, such as MongoDB, to manage metadata and share data. At the bottom, it says "To get started, follow the installation instructions or have a look at the minimal examples first." and includes a footnote about the BSD 3-Clause License. The University of Michigan logo is in the bottom right corner.



The screenshot shows the signac documentation page. It features a large blue 'S' logo. The heading is "Welcome to the signac framework documentation!". A paragraph explains: "The signac framework supports researchers in managing project-related data with a well-defined indexable storage layout for data and metadata. This streamlines post-processing and analysis and makes data collectively accessible." Below this, it says "This is the overall framework documentation. It provides a comprehensive overview on what you can do with packages that are part of the signac framework. If you are new to signac, the best place to start is to read the Introduction and the Tutorial." A "Contents" section lists various topics: Introduction (Overview, Packages), Installation (Install with conda, Install with pip, Installation from Source), Quickstart, Tutorial (Basics, Workflow, Job scripts and cluster submission), Examples (Ideal Gas, MD with HOOMD-blue), Recipes (How to migrate (change) the data space schema, How to integrate signac-flow with MATLAB or other software without Python interface, How to implement MPI-parallelized operations, How to handle containers), Tips and Tricks (Do not replicate job metadata in file names, Do not hard-code job ids in your scripts, How to design a good schema), Support (Contribute, Changing, signac, signac-flow), and License. At the bottom, there's a section for "Indices and tables" with links to Index, Module Index, and Search Page.

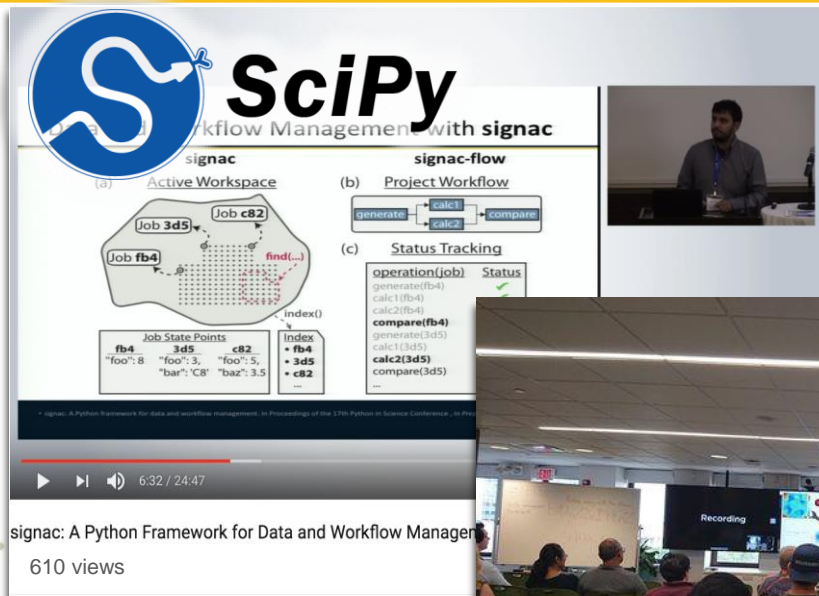


The screenshot shows the signac Lobby chat room. The header says "signac/Lobby". The chat area displays a conversation between several users. Stephen Thomas (@amburan) mentions not realizing that command-line parameters can be supplied directly. Carl Simon Adorf (@csadorf) explains the meaning of dashes in command-line arguments and provides a code snippet: `python src/project.py submit -o equil --nn 1 -w 10 --bundle 10 --serial -- -qregular`. Matt Thompson (@mattwthompson) and Carl Simon Adorf (@csadorf) discuss the use of the `--partition=regular` option. The chat also shows a list of recent activity, including pushes to the signac repository and commits. The bottom of the chat window has a text input field and a "Click here to type a chat message. Supports GitHub flavoured markdown." prompt.

Impact of signac

- In use by several research groups across the country, especially in materials science and molecular simulations
 - Princeton, Vanderbilt, Boise State, Air Force Research Laboratory, ...
- 10,000+ downloads of signac & signac-flow conda packages
- Actively used on projects with thousands of jobs and terabytes of data
- signac-flow has built-in support for Flux and XSEDE clusters, and adaptable templates for any SLURM/TORQUE scheduler

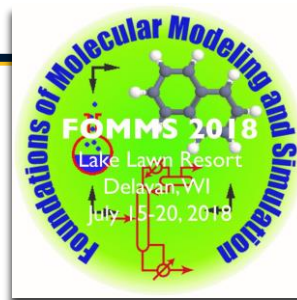
Community Building



The slide features the SciPy logo and the title "Workflow Management with signac". It is divided into three parts: (a) Active Workspace showing a directed graph of jobs (fb4, 3d5, c82) and their state points; (b) Project Workflow showing a linear sequence of operations (generate, calc1, calc2, compare); and (c) Status Tracking showing a table of operations and their statuses. A video player interface at the bottom indicates the video is 6:32 / 24:47 long and has 610 views.

signac: A Python Framework for Data and Workflow Management

610 views

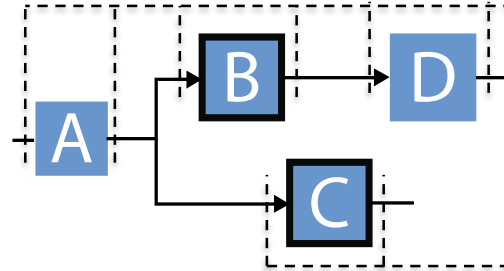
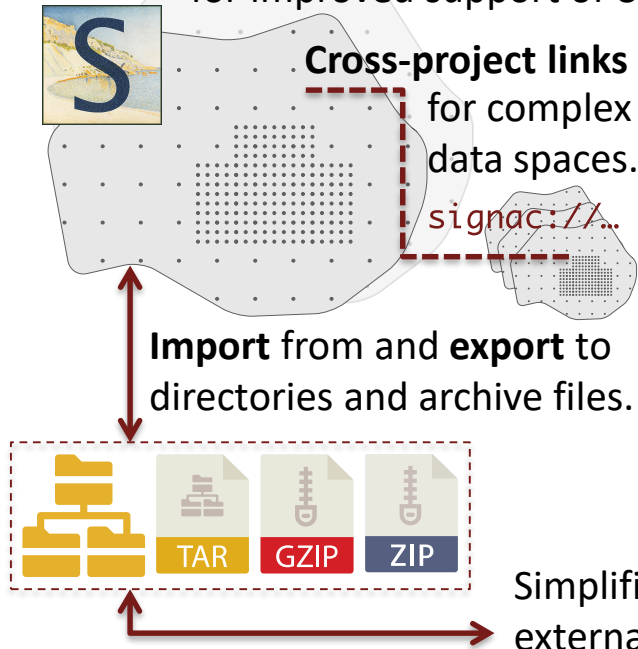


[signac Projectile Demo]

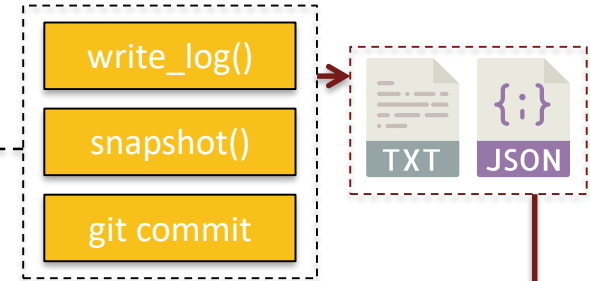
<https://github.com/bdice/signac-micde-cnsccs-2018>

Latest Development

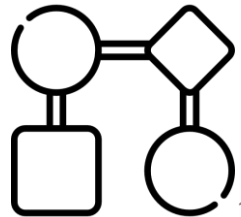
Better handling of **dynamically** growing data spaces for improved support of **optimization** workflows.



Execution hooks allow consistent tracking of all operations that manipulate the data space.



An **operations log** can be used to generate a graph of actions.

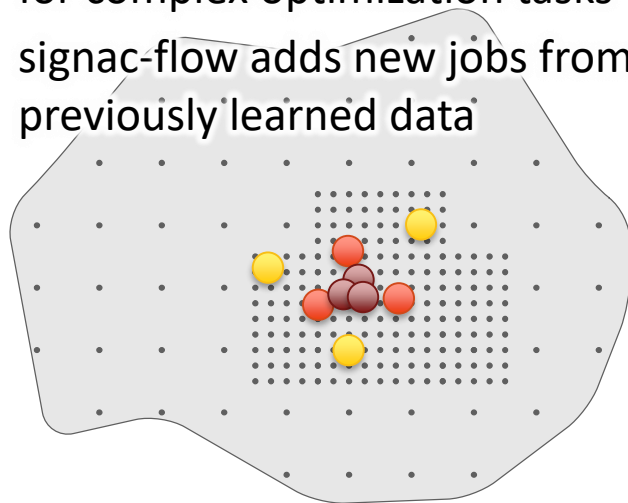


Future goal:

Automated Machine Learning Workflows

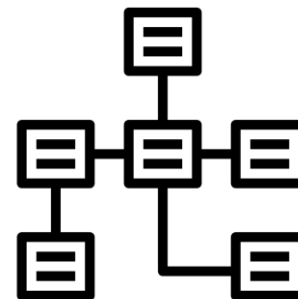
From **dynamic** workflows...

- signac data spaces are suited for complex optimization tasks
- signac-flow adds new jobs from previously learned data



... to distributed, iterative **exploration**

- prioritized exploration in high dimensional spaces
- distribution across systems



Thank you!

signac is a product of The Glotzer Group at the University of Michigan. Development and deployment supported by MICCoM, as part of the Computational Materials Sciences Program funded by the U.S. Department of Energy, Office of Science, Basic Energy Sciences, Materials Sciences and Engineering Division, under Subcontract No. 6F-30844. Project conceptualization and implementation supported by the National Science Foundation, Award # DMR 1409620.

www.signac.io

