# noWorkflow: Capturing, Analyzing, and Managing Provenance from Python Scripts

**João Felipe Nicolaci Pimentel (UFF),** Juliana Freire (NYU), Vanessa Braganholo (UFF), Leonardo Murta (UFF)

Instituto de **computação**

**NYU**
POLYTECHNIC SCHOOL
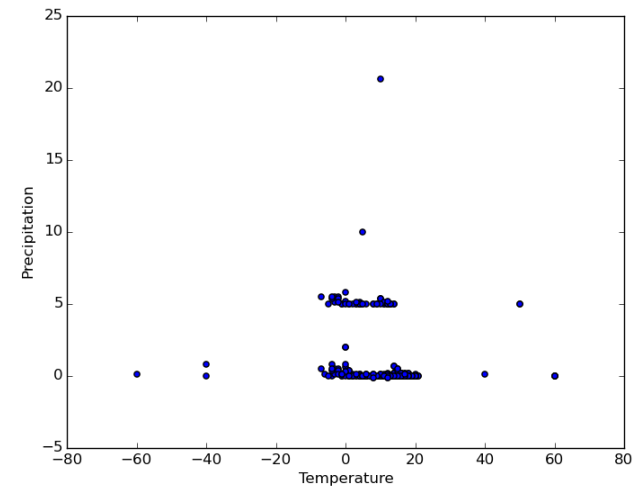OF ENGINEERING

# Outline

- Motivation
  - Scientific Experiments
  - Provenance
- noWorkflow
  - Collection
  - Management
  - Analysis
- IPython Notebook

# Scientific Experiments

noWorkflow: Capturing, Analyzing, and Managing
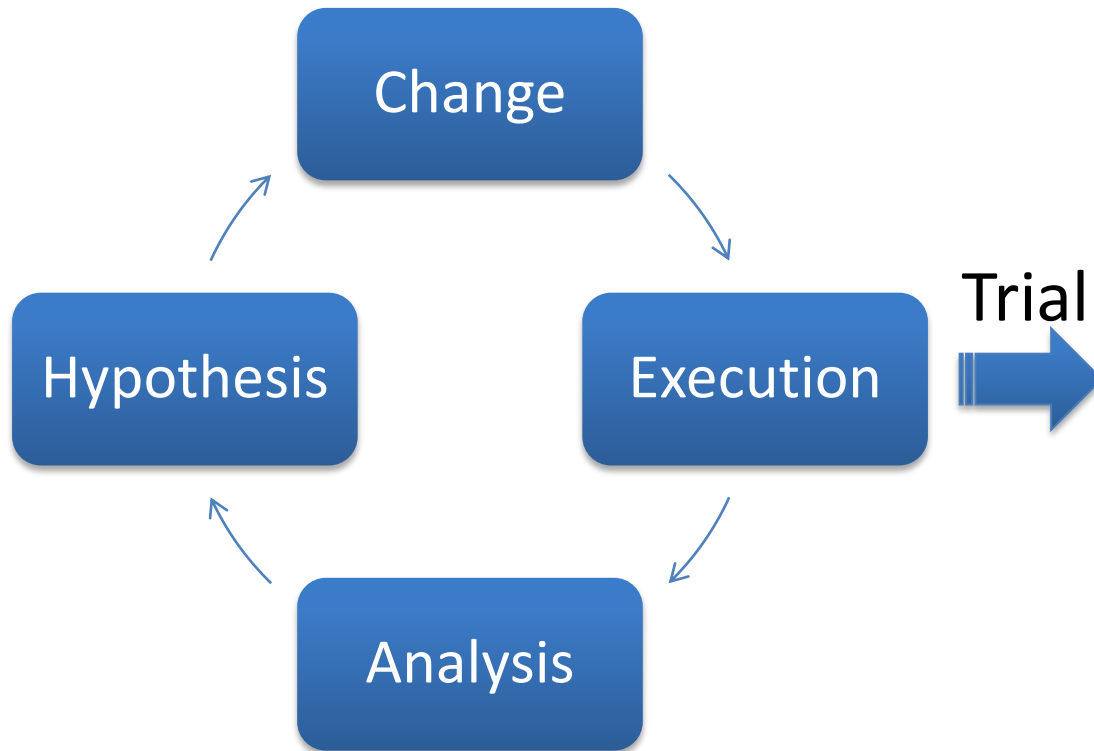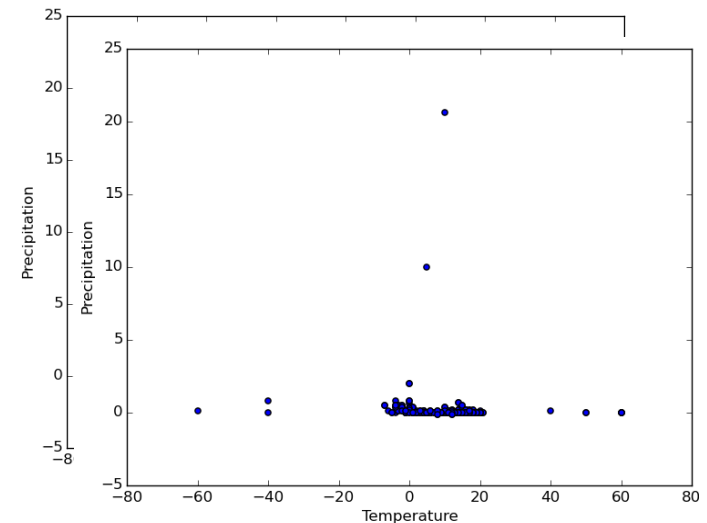Provenance from Python Scripts

0101010
011001110
01110101001

# Exploratory Development



Change

Hypothesis

Execution

Analysis

Trial

# Exploratory Development

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# Exploratory Development

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# Questions

- How long did it take to execute each trial?
- How was the source code for each trial?
- Which data were used?
- Which transformations were performed?
- Can I reproduce it?

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# Provenance



wasDerivedFrom

**Entity**

wasAttributedTo

wasGeneratedBy

**Agent**    used

wasAssociatedWith    **Activity**

http://www.w3.org/TR/prov-primer/

noWorkflow: Capturing, Analyzing, and Managing Provenance from Python Scripts

# Provenance

*"Refers to the documented history of an art object, or the documentation of processes in a digital object's life cycle"* [Moreau et al., 2008]

# How do I capture it?



- **Workflow Systems**
  - Transparent
  - Large start-up costs
  - Hard to integrate tools
- **OS-based solutions**
  - Transparent
  - Hard to connect to the semantics of the experiment
- **Script-based solutions**
  - Users must annotate their script with provenance capture directives

# noWorkflow

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# noWorkflow

- **Transparently** captures provenance of Python scripts
  - No changes required!
- Allows users to analyze provenance information



**Manage, Assess,** and **Reproduce**

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# Python

# vs

# noWorkflow

# simulation.py

```
 1| import csv
 2| import sys
 3| import matplotlib.pyplot as plt
 4| from simulator import simulate
 5|
 6| def run_simulation(data_a, data_b):
       ...
11|
12| def csv_read(f):
       ...
18|
19| def extract_column(data, column):
       ...
24|
25| def plot(data):
       ...
```

# simulation.py

main

```
36| data_a = sys.argv[1]
37| data_b = sys.argv[2]
38| data = run_simulation(data_a, data_b)
39| plot(data)
```

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# simulation.py

main

```
36| data_a = sys.argv[1]
37| data_b = sys.argv[2]
38| data = run_simulation(data_a, data_b)
39| plot(data)
```

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# simulation.py

```
 6| def run_simulation(data_a, data_b):
 7|     a = csv_read(data_a)
 8|     b = csv_read(data_b)
 9|     data = simulate(a, b)
10|     return data
```

```
12| def csv_read(f):
13|     reader = csv.reader(open(f, 'rU'), delimiter=':')
14|     data = []
15|     for row in reader:
16|         data.append(row)
17|     return data
```

# simulation.py

main

```
36| data_a = sys.argv[1]
37| data_b = sys.argv[2]
38| data = run_simulation(data_a, data_b)
39| plot(data)
```

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# simulation.py

```python
25| def plot(data):
26|     # Get Temperature
27|     t = extract_column(data, 0)
28|     # Get Precipitation
29|     p = extract_column(data, 1)
30|     plt.scatter(t, p, marker='o')
31|     plt.xlabel('Temperature')
32|     plt.ylabel('Precipitation')
33|     plt.savefig('output.png')
```

# Comparison

$ **python** simulation.py \
> data1.dat data2.dat
$ display output.png

$ **now run** simulation.py \
> data1.dat data2.dat
$ display output.png

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# $ now vis

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# Collection

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# $ now run simulation.py data1.dat data2.dat

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# Deployment

- Environment

- Module dependencies

```
1| import csv
2| import sys
3| import matplotlib.pyplot as plt
4| from simulator import simulate
```

TABLE environment_attr

| id | name | value | trial_id |
|----|------|-------|----------|
| 1 | SC_REALTIME_SIGNALS | 200809 | 1 |
| 2 | rvm_version | 1.25.28 (stable) | 1 |
| 3 | SC_PII_OSI_COTS | -1 | 1 |
| 4 | SC_PII_OSI | -1 | 1 |
| 5 | SC_T_IOV_MAX | -1 | 1 |
| 6 | RUBY_VERSION | ruby-2.1.2 | 1 |
| 7 | SC_THREADS | 200809 | 1 |
| 8 | LC_PAPER | pt_BR.UTF-8 | 1 |
| 9 | SC_AIO_MAX | -1 | 1 |
| 10 | PROCESSOR | x86_64 | 1 |
| 11 | SC_USHRT_MAX | 65535 | 1 |
| 12 | SC_THREAD_KEYS_MAX | 1024 | 1 |

1 to 100 of 122

TABLE dependency

| trial_id | module_id |
|----------|-----------|
| 1 | 1 |
| 1 | 2 |
| 1 | 3 |
| 1 | 4 |
| 1 | 5 |
| 1 | 6 |
| 1 | 7 |
| 1 | 8 |
| 1 | 9 |
| 1 | 10 |
| 1 | 11 |
| 1 | 12 |

1 to 100 of 669

TABLE module

| id | name | version | path | code_hash |
|----|------|---------|------|-----------|
| 1 | BaseHT... | 0.3 | /usr/lib/... | 3fc68f6f19... |
| 2 | ConfigP... | | /usr/lib/... | 765dde108... |
| 3 | Cookie | | /usr/lib/... | e3a11a4d5... |
| 4 | FixTk | | /usr/lib/... | dfbe55683... |
| 5 | PIL | 1.1.7 | /home/j... | 5c969cc37... |
| 6 | PIL.Bmp... | 0.7 | /home/j... | 074c413f0... |
| 7 | PIL.Gifl... | 0.9 | /home/j... | 197e5bd77... |
| 8 | PIL.Gim... | | /home/j... | 671cae435... |
| 9 | PIL.Gim... | | /home/j... | 484b99960... |
| 10 | PIL.Image | 1.1.7 | /home/j... | e4e46dfff5... |
| 11 | PIL.Imag... | | /home/j... | 0958c7146... |
| 12 | PIL.Imag... | | /home/j... | fe1e169c2... |

1 to 100 of 669

# Definition

- Script

- Function definitions

- Arguments and Globals

```
 6| def run_simulation(data_a, data_b):
12| def csv_read(f):
19| def extract_column(data, column):
25| def plot(data):
```

Hash: 3461fe48c17619faeb81a6334567a47a162c7788

TABLE function_def

| id | name | code_hash | trial_id |
|----|------|-----------|----------|
| 1 | plot | bef07f4bbf... | 1 |
| 2 | run_simulation | 97894a102... | 1 |
| 3 | extract_column | 36fca5011c... | 1 |
| 4 | csv_read | d914038c9... | 1 |

1 to 4 of 4

TABLE object

| id | name | type | function_def_id |
|----|------|------|-----------------|
| 1 | data | ARGUMENT | 1 |
| 2 | extract_column | FUNCTION_CALL | 1 |
| 3 | data_a | ARGUMENT | 2 |
| 4 | data_b | ARGUMENT | 2 |
| 5 | simulate | FUNCTION_CALL | 2 |
| 6 | csv_read | FUNCTION_CALL | 2 |
| 7 | data | ARGUMENT | 3 |
| 8 | column | ARGUMENT | 3 |
| 9 | float | FUNCTION_CALL | 3 |
| 10 | f | ARGUMENT | 4 |
| 11 | open | FUNCTION_CALL | 4 |

1 to 11 of 11

noWorkflow: Capturing, Analyzing, and Managing Provenance from Python Scripts

# Execution

- Files content
- Function calls
- Parameter and global values
- Program arguments

TABLE file_access

| id | name | mode | buffering | content_hash_before | content_hash_after | timestamp | functio... | trial_id |
|---|---|---|---|---|---|---|---|---|
| 1 | data1.dat | rU | default | 28f4192700d9e5d281... | 28f4192700d9e5d2... | 2015-05-14... | 4 | 1 |
| 2 | data2.dat | rU | default | 802a73cb49af95840b... | 802a73cb49af9584... | 2015-05-14... | 188 | 1 |
| 3 | /home/j... | rb | default | 1d7f6fa0c34e3d50be... | 1d7f6fa0c34e3d50... | 2015-05-14... | 1102 | 1 |
| 4 | output.p... | wb | default | 605d84723a48621a88... | 605d84723a48621a... | 2015-05-14... | 1102 | 1 |

1 to 4 of 4

TABLE function_activatior

| id | name | line | return | start | finish | caller_id | trial_id |
|---|---|---|---|---|---|---|---|
| 1 | /home/j... | 126 | None | 2015-05-... | 2015-05-... | | 1 |
| 2 | run_sim... | 38 | [['0.0', '0... | 2015-05-... | 2015-05-... | 1 | 1 |
| 3 | csv_read | 7 | [['0.0', '0... | 2015-05-... | 2015-05-... | 2 | 1 |
| 4 | open | 13 | <open fi... | 2015-05-... | 2015-05-... | 3 | 1 |
| 5 | reader | 13 | | 2015-05-... | 2015-05-... | 3 | 1 |
| 6 | list.appe... | 16 | | 2015-05-... | 2015-05-... | 3 | 1 |
| 7 | list.appe... | 16 | | 2015-05-... | 2015-05-... | 3 | 1 |
| 8 | list.appe... | 16 | | 2015-05-... | 2015-05-... | 3 | 1 |
| 9 | list.appe... | 16 | | 2015-05-... | 2015-05-... | 3 | 1 |
| 10 | list.appe... | 16 | | 2015-05-... | 2015-05-... | 3 | 1 |
| 11 | list.appe... | 16 | | 2015-05-... | 2015-05-... | 3 | 1 |

1 to 100 of 1102

TABLE object_value

| id | name | value | type | function_activation_id |
|---|---|---|---|---|
| 1 | data_b | 'data2.dat' | ARGUMENT | 2 |
| 2 | data_a | 'data1.dat' | ARGUMENT | 2 |
| 3 | f | 'data1.dat' | ARGUMENT | 3 |
| 4 | args | ('rU',) | ARGUMENT | 4 |
| 5 | name | 'data1.dat' | ARGUMENT | 4 |
| 6 | f | 'data2.dat' | ARGUMENT | 187 |
| 7 | args | ('rU',) | ARGUMENT | 188 |
| 8 | name | 'data2.dat' | ARGUMENT | 188 |

1 to 33 of 33

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# Trial

TABLE trial

| id | start | finish | script | code_hash | arguments | inherited_id | parent_id | run |
|---|---|---|---|---|---|---|---|---|
| 1 | 2015-05-... | 2015-05-... | simulati... | 3461fe48c1... | data1.dat data2.dat | | | 1 |

1   to   1   of   1

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# Fine-grained Collection

**$ now run –e Tracer simulation.py data1.dat data2.dat**

- Default provenance
- Variable assignments
- Dependencies

```
 6| def run_simulation(data_a, data_b):
 7|     a = csv_read(data_a)
 8|     b = csv_read(data_b)
 9|     data = simulate(a, b)
10|     return data
```

# Management

1
2
3

noWorkflow: Capturing, Analyzing, and Managing Provenance from Python Scripts

# Restore

(1) $ now run simulation.py data1.dat data2.dat

(2) $ now run simulation.py data1.dat data3.dat
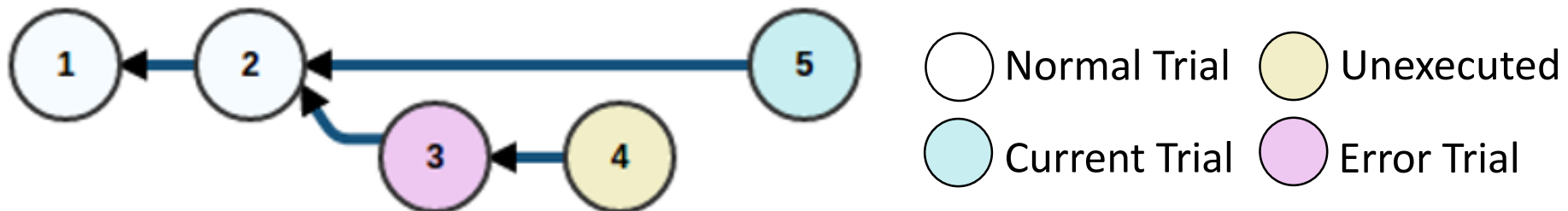
(3) $ now run simulation.py

Error

Try to fix simulation.py, save it, but do not run it.

(4) **$ now restore -li 2**

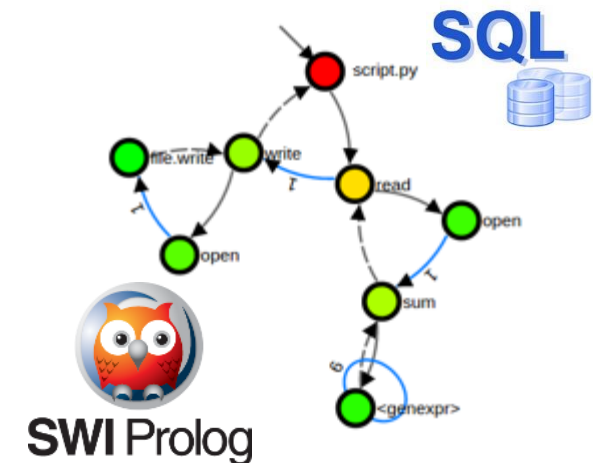Restores local modules, input and output files from trial 2

(5) $ now run simulation.py data1.dat data4.dat



Normal Trial    Unexecuted

Current Trial    Error Trial

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# $ now restore –li 2

```
[now] Backup Trial 4 created
[now] File simulation.py from trial 2 restored
[now] File
/home/joao/projects/nowissues/scipyla_weather/simulator.py from
trial 2 restored
[now] File output.png from trial 2 restored
[now] File data3.dat from trial 2 restored
[now] File data1.dat from trial 2 restored
```

noWorkflow: Capturing, Analyzing and Managing
Provenance from Python Scripts

# Analysis

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# Command line

- $ now list
- $ now show [trial]
- $ now diff [trial1] [trial2]
- $ now export –r [trial]
- $ now export –i [trial]
- $ now export –i history
- $ now export –i diff:[trial1]:[trial2]
- $ now vis

# $ now list

```
[now] trials available in the provenance store:
  Trial 1: simulation.py data1.dat data2.dat
          with code hash 3461fe48c17619faeb81a6334567a47a162c7788
          ran from 2015-05-15 00:45:09.908030 to 2015-05-15
00:45:20.602758
  Trial 2: simulation.py data1.dat data3.dat
          ...
  Trial 3: simulation.py
          ...
  Trial 4: simulation.py <restore 2>
          with code hash cf879bd94d8c5942de800a9a17aed2f98acd1300
          ran from 2015-05-15 00:47:31.103699 to None
  Trial 5: simulation.py data1.dat data4.dat
          with code hash 3461fe48c17619faeb81a6334567a47a162c7788
          ran from 2015-05-15 01:03:44.096991 to 2015-05-15
01:03:54.803461
```

# $ now show 5 -a

```
[now] trial information:
  Id: 5
  Inherited Id: None
  Script: simulation.py
  Code hash: 3461fe48c17619faeb81a6334567a47a162c7788
  Start: 2015-05-15 01:03:44.096991
  Finish: 2015-05-15 01:03:54.803461
[now] this trial has the following function activation graph:
  126:
/home/joao/projects/nowissues/scipyla_weather/simulation.py
(2015-05-15 01:03:52.165483 - 2015-05-15 01:03:54.803433)
      Globals, Arguments, Return value
    38: run_simulation (2015-05-15 01:03:52.165687 - 2015-05-15
01:03:54.176796)
        Globals:
        Arguments: data_b = 'data4.dat', data_a = 'data1.dat'
```

# $ now diff 1 2

```
[now] trial diff:
  duration changed from 10694728 to 10662788
  start changed from 2015-05-15 00:45:09.908030 to 2015-05-15
00:45:37.045604
  finish changed from 2015-05-15 00:45:20.602758 to 2015-05-15
00:45:47.708392
  arguments changed from data1.dat data2.dat to data1.dat
data3.dat
  parent_id changed from None to 1
```
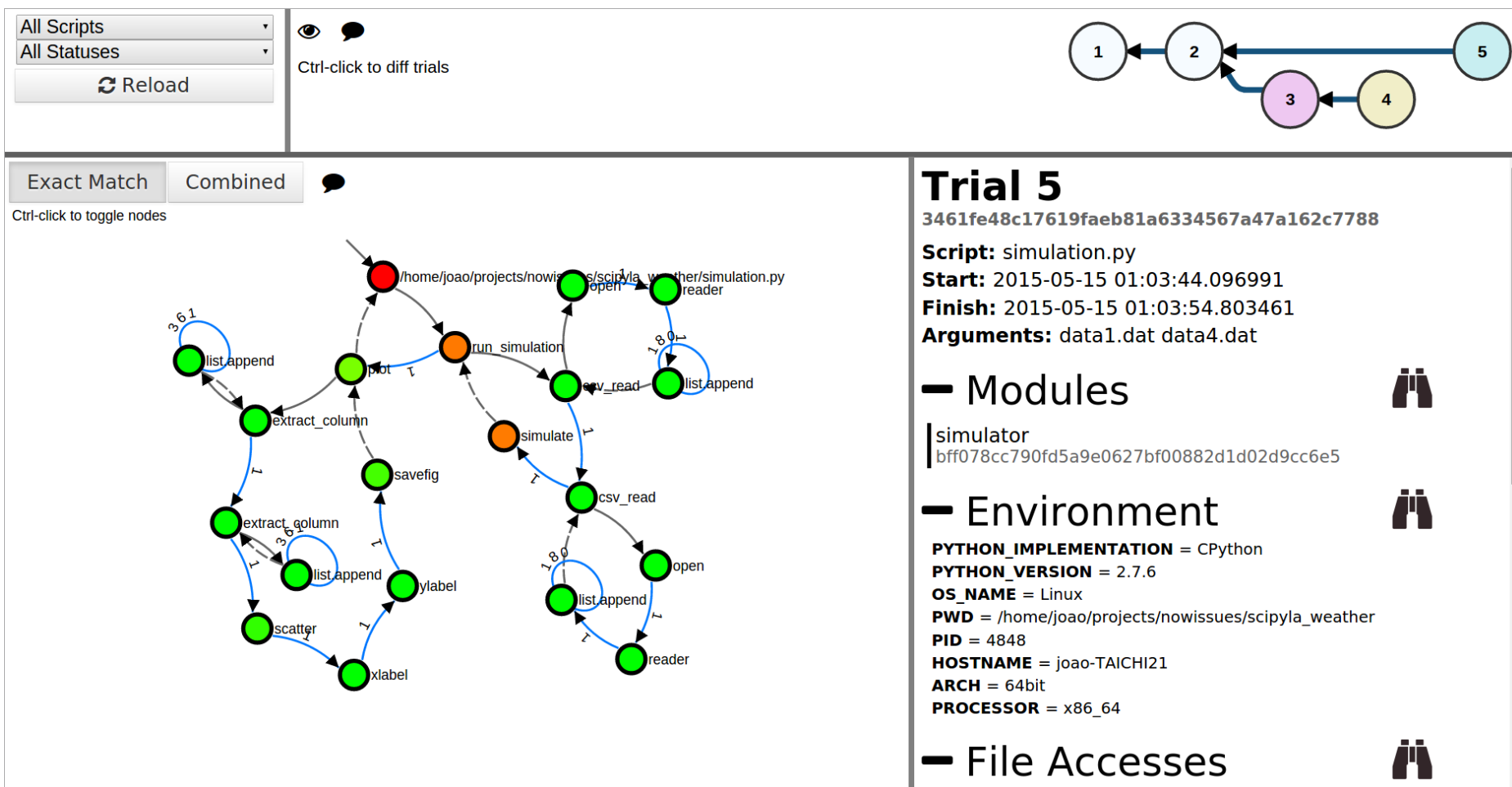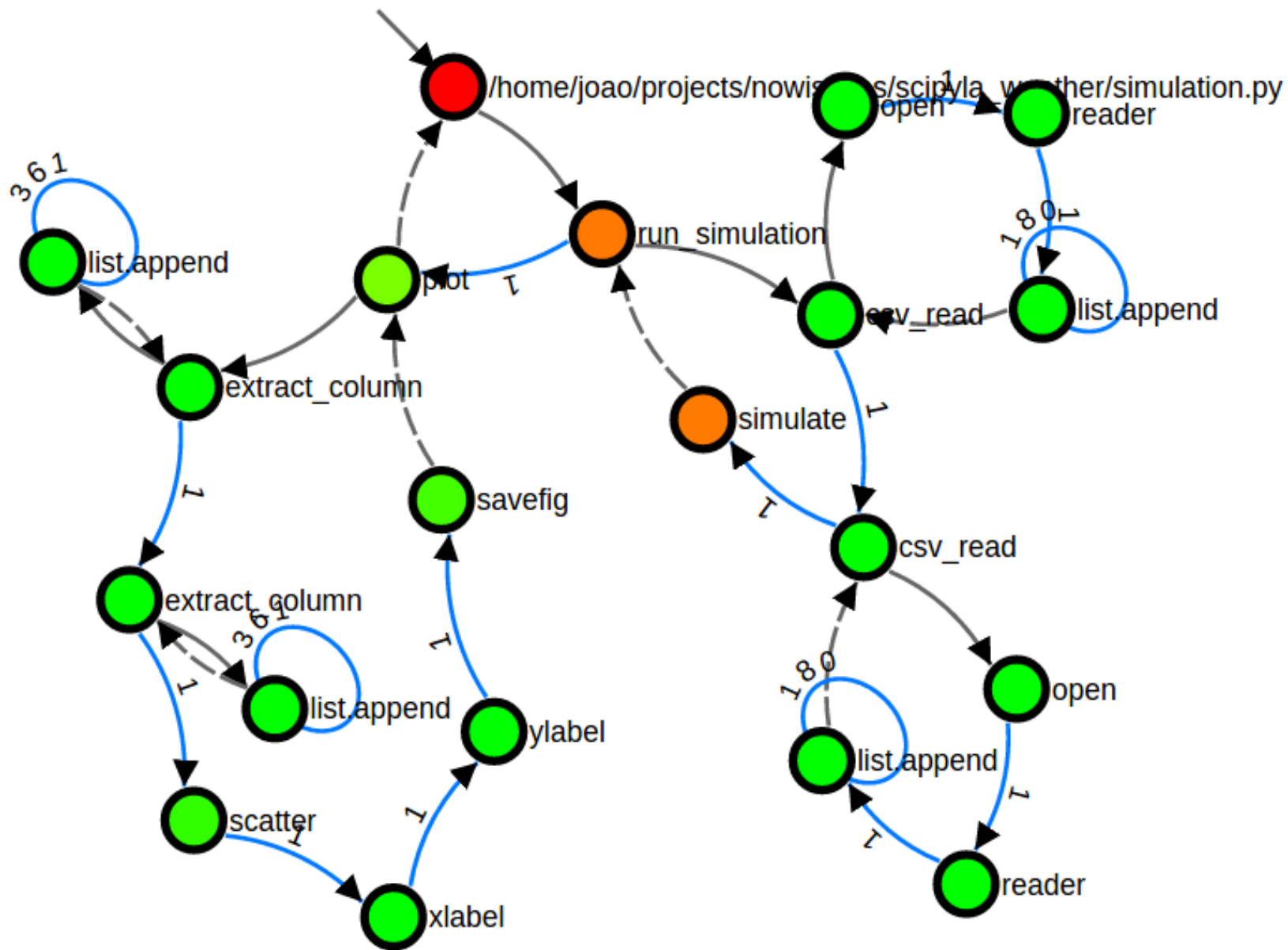
# Export

- Trial to Prolog: **$ now export –r 1 > trial1.pl**

- Trial to IPython Notebook: **$ now export –i 1**

- History to Notebook: **$ now export –i history**
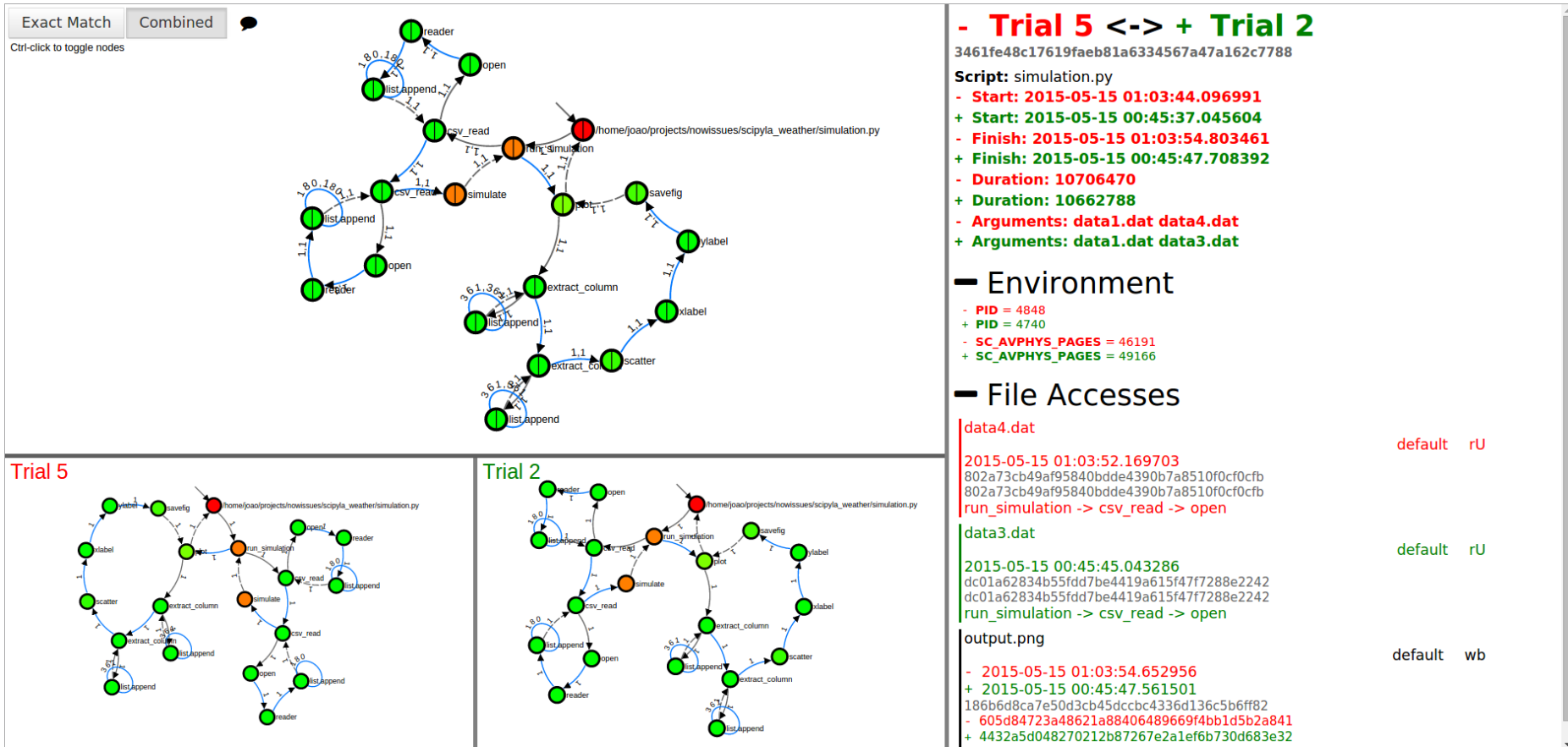
- Diff to Notebook: **$ now export –i diff:1:2**

# $ now vis

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

noWorkflow: Capturing, Analyzing, and Managing
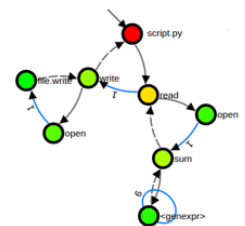Provenance from Python Scripts

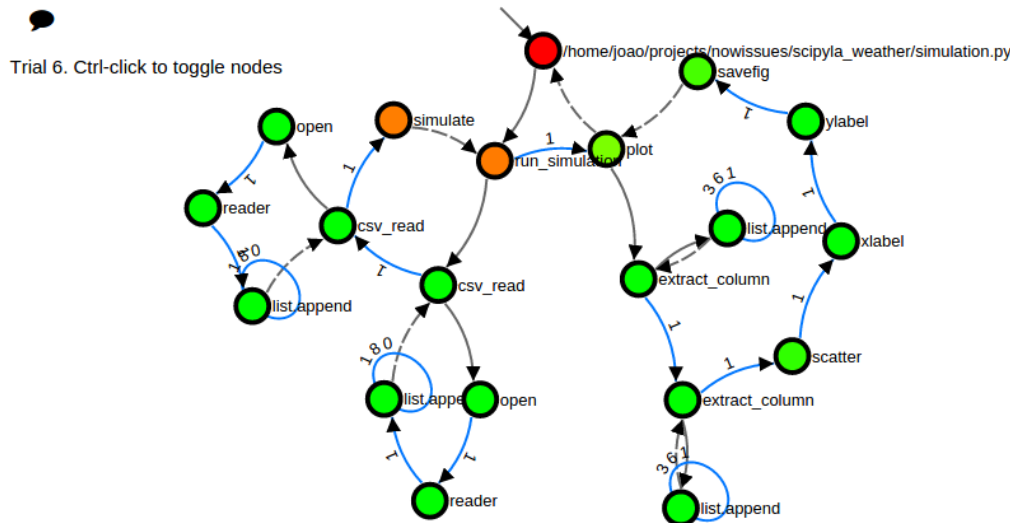# $ now vis

# IPython meets noWorkflow

# $ ipython notebook

```
In [1]: %load_ext noworkflow
        %now_set_default graph_width=800 graph_height=300
        import noworkflow.now.ipython as nip
```

```
In [2]: data1, data2 = 'data1.dat', 'data2.dat'
        trial = %now_run --name ipython_script simulation.py $data1 $data2
        trial
```

Out[2]:



Trial 6. Ctrl-click to toggle nodes

```
In [3]: trial.environment()['PWD']
```

Out[3]: u'/home/joao/projects/nowissues/scipyla_weather'

# Collection
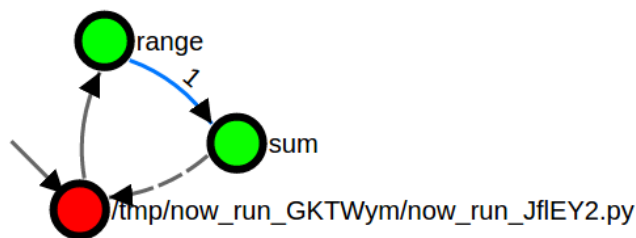
```
In [5]: size = 5

In [6]: %%now_run --name ipython_script --out=out_var $size
        import sys
        l = range(int(sys.argv[1]))
        c = sum(l)
        print(c)

Out[6]:
```

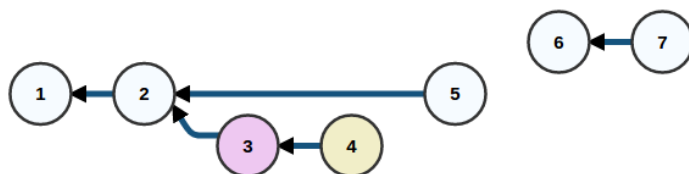Trial 7. Ctrl-click to toggle nodes



```
In [7]: out_var
Out[7]: '10\n'
```

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# Analysis

noWorkflow: Capturing, Analyzing, and Managing
Provenance from Python Scripts

# Queries

```
In [14]: trial.id
Out[14]: 6
```

```
In [15]: %%now_prolog --result result {trial.id}
         duration({trial.id}, simulate, X)
```

```
In [16]: for match in result:
             print(match['X'])

         2.00223684311
```

```
In [17]: %%now_sql
         SELECT DISTINCT script FROM trial
Out[17]:
```

| script |
| --- |
| simulation.py |
| ipython_script |

# Conclusion

- noWorkflow allows users to capture and analyze provenance from Python Scripts

- It is easy to install and use it
  - **$ pip install noworkflow[all]**

- Open source. Please, submit issues at
  - https://github.com/gems-uff/noworkflow

# noWorkflow: Capturing, Analyzing, and Managing Provenance from Python Scripts

joaofelipenp@gmail.com
https://github.com/gems-uff/noworkflow

**João Felipe Nicolaci Pimentel (UFF),** Juliana Freire (NYU), Vanessa Braganholo (UFF), Leonardo Murta (UFF)