# FACILITATING REPRODUCIBILITY AFTER THE FACT

**Fernando Chirigati**

ViDA – Visualization and Data Analysis Lab
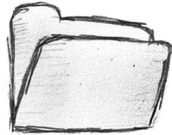
NYU Polytechnic School of Engineering
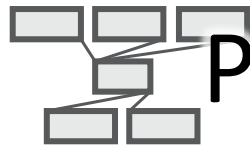
**NEW YORK UNIVERSITY**

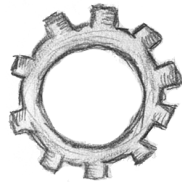Reproducibility may require a lot of computational effort.
Why?

# Too many dependencies!

 DATA

 PROVENANCEOW

 ENVIRONMENT

# Too much to do, too little time!

*"authors have complained that the process **requires too much work for the benefit derived**"*

Bonnet et al., SIGMOD Record 2011

*"**Insufficient time** is the main reason why scientists do not make their data and experiment available and reproducible."*

Carol Tenopir, Beyond the PDF 2 Conference

*"**77%** claim that they do not have **time to document and clean up the code**."*

Victoria Stodden, Survey of the Machine Learning Community – NIPS 2010

*"It would require **huge amount of effort** to make our code work with the latest versions of these tools."*

Collberg et al., Repeatability and Benefaction in Computer Systems Research, University of Arizona TR 14-04

# Planning for Reproducibility

Scientific Workflow Systems (VisTrails, Taverna, Kepler, …)

Virtual Machines and Containers (VirtualBox, Vagrant, Docker, …)

Configuration Management Tools (Chef, Puppet, …)

… and many others !

But what about ***reproducibility after the fact***?

Again, time-consuming and error-prone!

# noWORKFLOW
## CAPTURING AND ANALYZING
## PROVENANCE OF SCRIPTS

**Joint work with:**   João Felipe Pimentel **(UFF)**
Leonardo Murta **(UFF)**
Vanessa Braganholo **(UFF)**
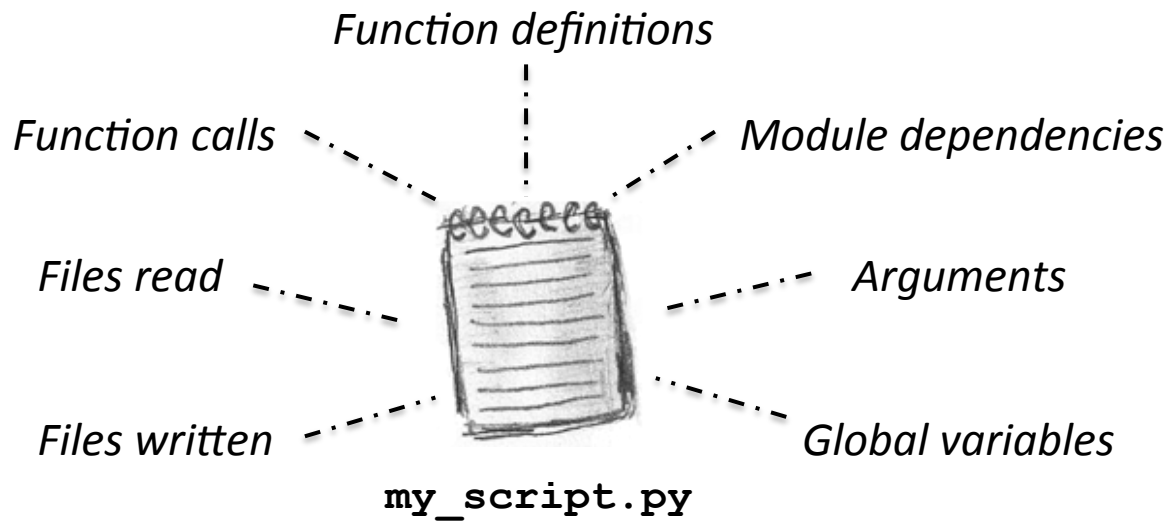David Koop **(UMass-Dartmouth)**
Juliana Freire **(NYU)**

**NEW YORK UNIVERSITY**

Instituto de
**Computação**

# noWorkflow

Transparently captures the *provenance* of a Python script and its various executions (trials)



Function definitions

Function calls

Module dependencies

Files read

Arguments

Files written

Global variables

**my_script.py**

*Non-intrusive*: no need for user-defined annotations, instrumented environment, or other requirements

Instead of running

```
$ python my_script.py
```

users run

```
$ now run my_script.py
```

That's it.

# Provenance Analysis
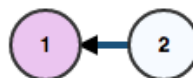
Script Evolution

Diff Analysis

Querying

Interactive Visualization
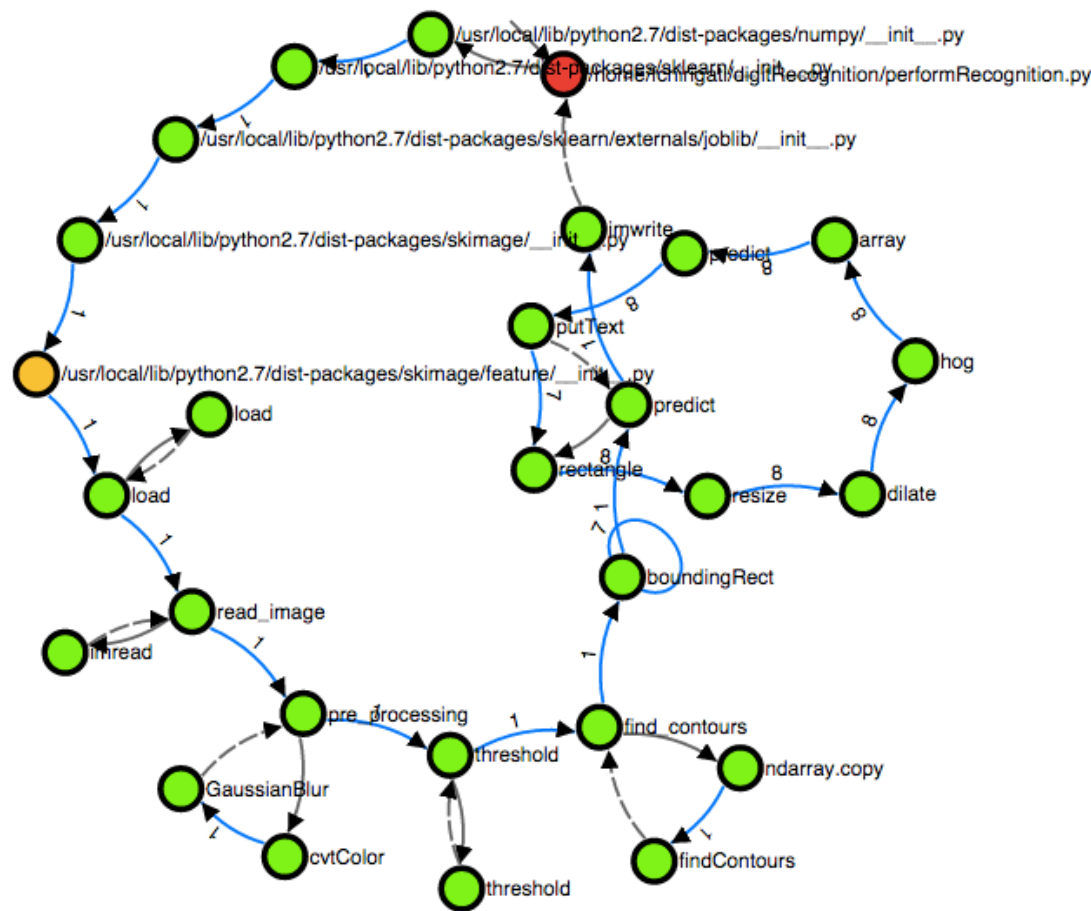
Support for Analysis in IPython Notebooks

# Script Evolution

# Diff Analysis

# Provenance Analysis

In [1]:
```
%load_ext noworkflow
%now_set_default graph_width=430 graph_height=150
nip = %now_ip
```

In [2]:
```
dry = 0
trial = %now_run --name ipython_script script.py $dry
trial
```

Out[2]:

Trial 6. Ctrl-click to toggle nodes

# ReproZip
## Packing Experiments for Reproducibility

**Joint work with:**    Rémi Rampin
Dennis Shasha
Juliana Freire

NEW YORK UNIVERSITY

# ReproZip is a packaging tool

PACKING STEP

UNPACKING STEP

*From reputablemoving.com*

*From wykop.pl*

# Packing Experiments

Computational Environment **E** (Linux)

Experiment → Executing → **reprozip** → Capturing Provenance (*ptrace + SQLite*) → Experiment Provenance

**Experiment Provenance**

**Data**
> Input files, output files, parameters, …

**Workflow**
> Executable programs and steps

**Environment**
> Environment variables, dependencies, …

Reproducible Package ← Configuring / Packing ← Configuration File ← Creating Configuration

# Unpacking Experiments

Computational Environment **E'** (potentially different than **E**)



**directory** — unpacks and reproduces from a single directory (Linux)

**chroot** — unpacks in a single directory and builds a full system environment (Linux)

**vagrant** — unpacks in a virtual machine using Vagrant (Linux, Mac OS X, Windows)

**docker** — unpacks in a Docker container (Linux, Mac OS X, Windows)

*Reproducible Package* → Extracting → **reprounzip**

# Unpacking Experiments

Running an unpacker:

Setting up: *setup*
Replicating results / changing command line parameters: *run*
Changing input files: *upload*
Getting output files: *download*

Natively installing required software dependencies:

*installpkgs*

# News!

**ReproZip** …

… has been adopted in the Bonneau Lab (NYU)

http://bonneaulab.bio.nyu.edu/

… will be used by the ACM SIGMOD 2015 Reproducibility Review

http://db-reproducibility.seas.harvard.edu/

… will be used by the Information Systems journal      http://www.journals.elsevier.com/information-systems/

# Try!

noWorkflow Website: *https://github.com/gems-uff/noworkflow*

L. Murta, V. Braganholo, F. Chirigati, D. Koop, and J. Freire: *noWorkflow: Capturing and Analyzing Provenance of Scripts*. In Provenance and Annotation of Data and Processes, vol. 8628, Lecture Notes in Computer Science (LNCS), pp. 71-83, Springer International Publishing, 2015

ReproZip Website: : *http://vida-nyu.github.io/reprozip/*

F. Chirigati, D. Shasha, and J. Freire: *Packing Experiments for Sharing and Publication*. In Proceedings of the 2013 International Conference on Management of Data (SIGMOD), pp. 977-980, 2013

F. Chirigati, D. Shasha, and J. Freire: *ReproZip: Using Provenance to Support Computational Reproducibility*. In Proceedings of the 5th USENIX conference on Theory and Practice of Provenance (TaPP), 2013

Send your feedback and interesting use cases!