

Streamline your Data Science projects with Ploomber

Eduardo Blancas, JupyterCon 2020

Data Science is all about experimentation

- Three stages for each experiment:
 1. Thinking
 2. Coding
 3. Execution
- **More experiments** = higher chance of **success**

Image: <https://unsplash.com/photos/JeInkKlI2Po>

Twitter: @edublancas / Website: ploomber.io



Poor man's solution

- Code project in a (long) **single file**
- Problem: **Not maintainable**
 1. No clear boundaries among tasks
 2. Leads to unwanted interactions
- We need **reproducible** and **maintainable** projects

Image: <https://www.flickr.com/photos/44674593@N00/7308867666>

Twitter: @edublancas / Website: ploomber.io



Why bother about maintainability?

Without a lot of structure, we forget what we've done in the past, we can't read each others' work, and we can't test whether what we've done is correct.¹

— *Patrick Ball*

Image: <https://unsplash.com/photos/r6mBXuHnxBk>

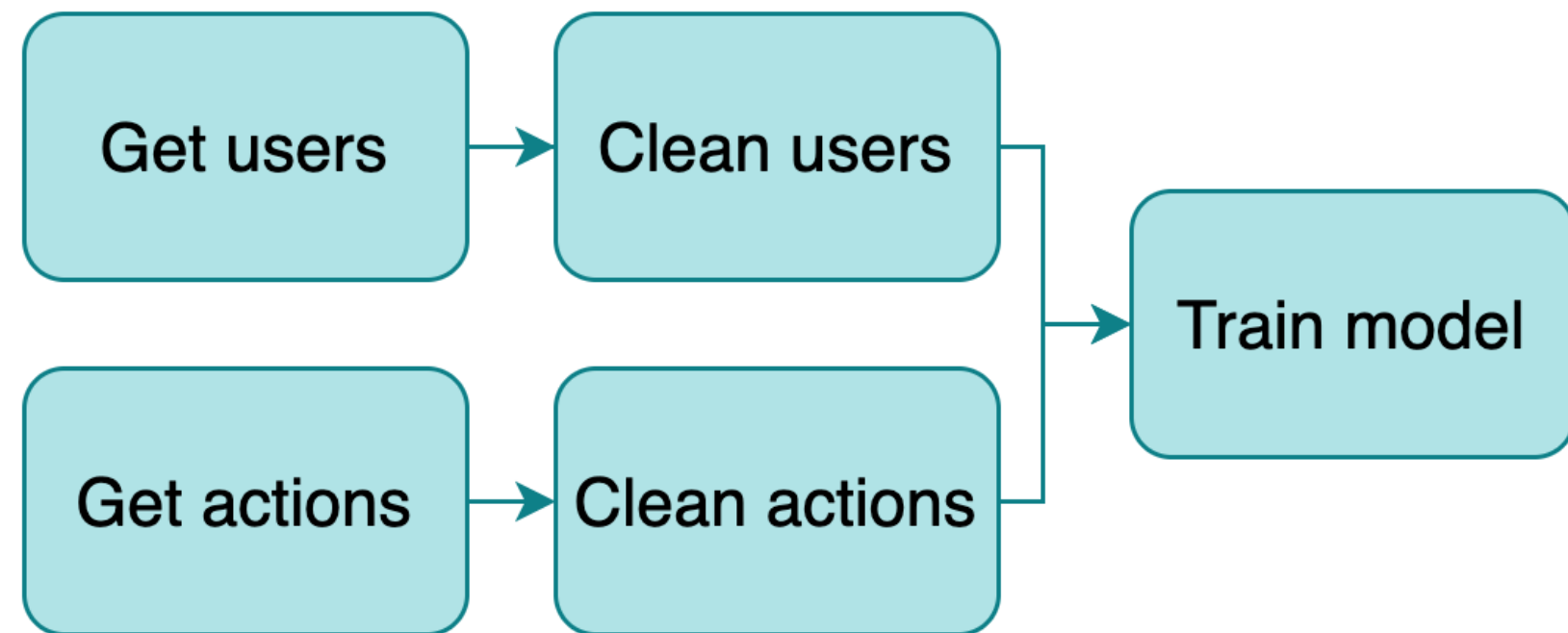
¹ <https://hrdag.org/2016/06/14/the-task-is-a-quantum-of-workflow/>

Twitter: @edublancas / Website: ploomber.io



Build a data pipeline

- **Break down** logic in small tasks
- **Outputs become inputs** of *downstream* tasks
- **New challenge:** Manage structure
 1. Multiple files
 2. Route outputs
 3. Orchestrate execution

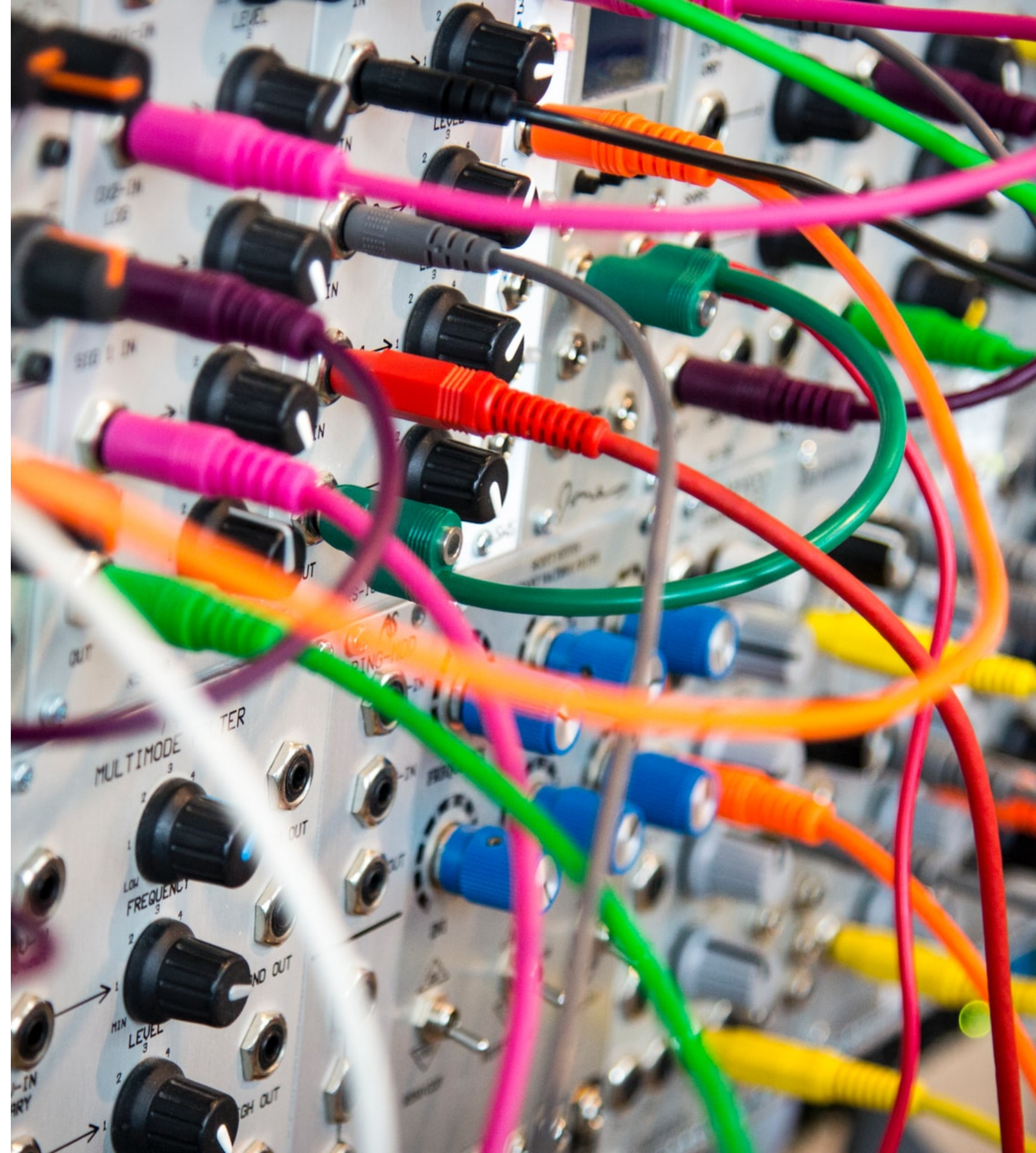


Current approach: Workflow managers

- Frameworks to develop pipelines (Make, Airflow, Luigi, etc)
- Problems:
 1. Learn a new tool
 2. Write pipeline code
 3. Adds unnecessary complexity

Image: <https://unsplash.com/photos/I090uFWoPal>

Twitter: @edublancas / Website: ploomber.io



Ploomber

Adheres to a *convention over configuration* philosophy, to allows us to **write pipelines with a focus on usability**.

Three conventions:

1. Each task is a script
2. Scripts declare dependencies via an `upstream` variable
3. And outputs via a `product` variable

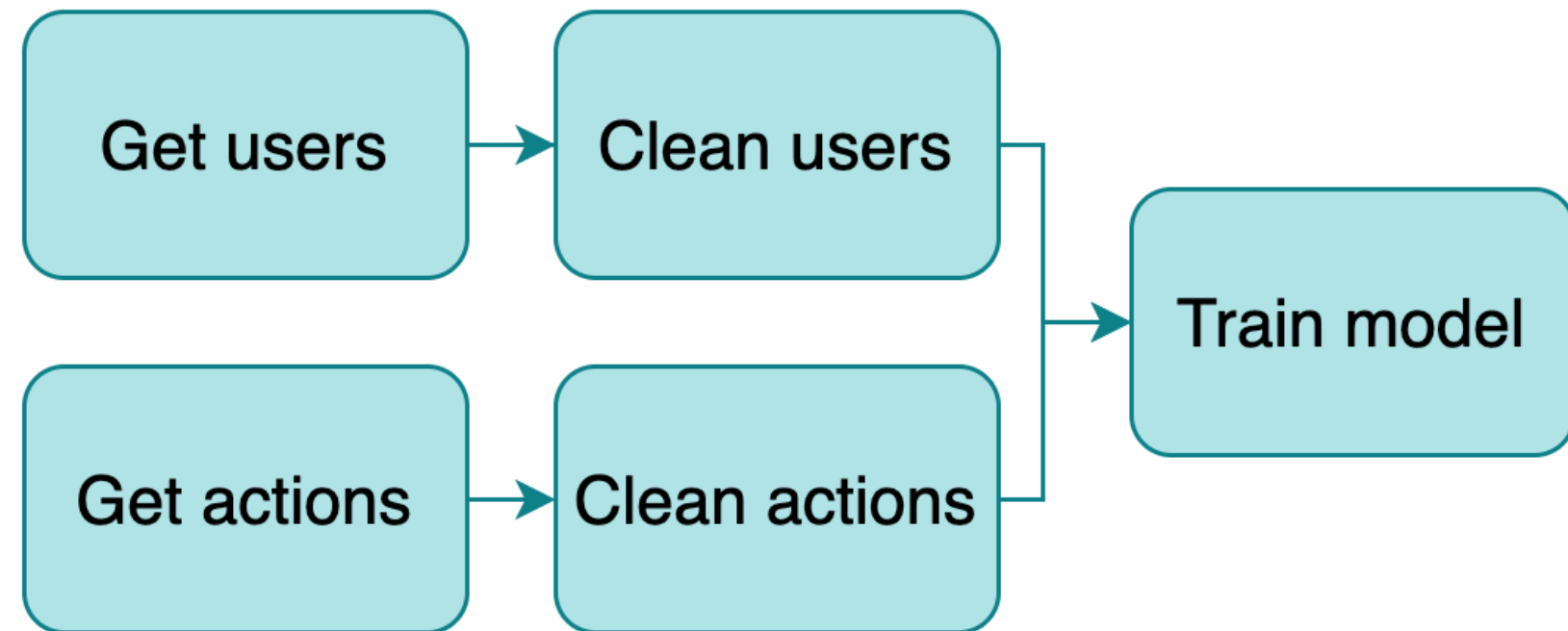
Python example

train-model.py:

```
import pandas as pd

# + tags=["parameters"]
upstream = ['clean-users', 'clean-actions']
product = {'nb': 'output/model-evaluation.ipynb',
           'model': 'output/model.pickle'}

# +
users = pd.read_parquet(
    upstream['clean-users']['data'])
actions = pd.read_parquet(
    upstream['clean-actions']['data'])
# code continues...
```



Cell injection

```
import pandas as pd
```

```
# + tags=["parameters"]
upstream = ['clean-users', 'clean-actions']
product = {'nb': 'output/model-evaluation.ipynb',
           'model': 'output/model.pickle'}
```

```
# +
users = pd.read_parquet(
    upstream['clean-users']['data'])
actions = pd.read_parquet(
    upstream['clean-actions']['data'])
# code continues...
```

```
1 import pandas as pd
```

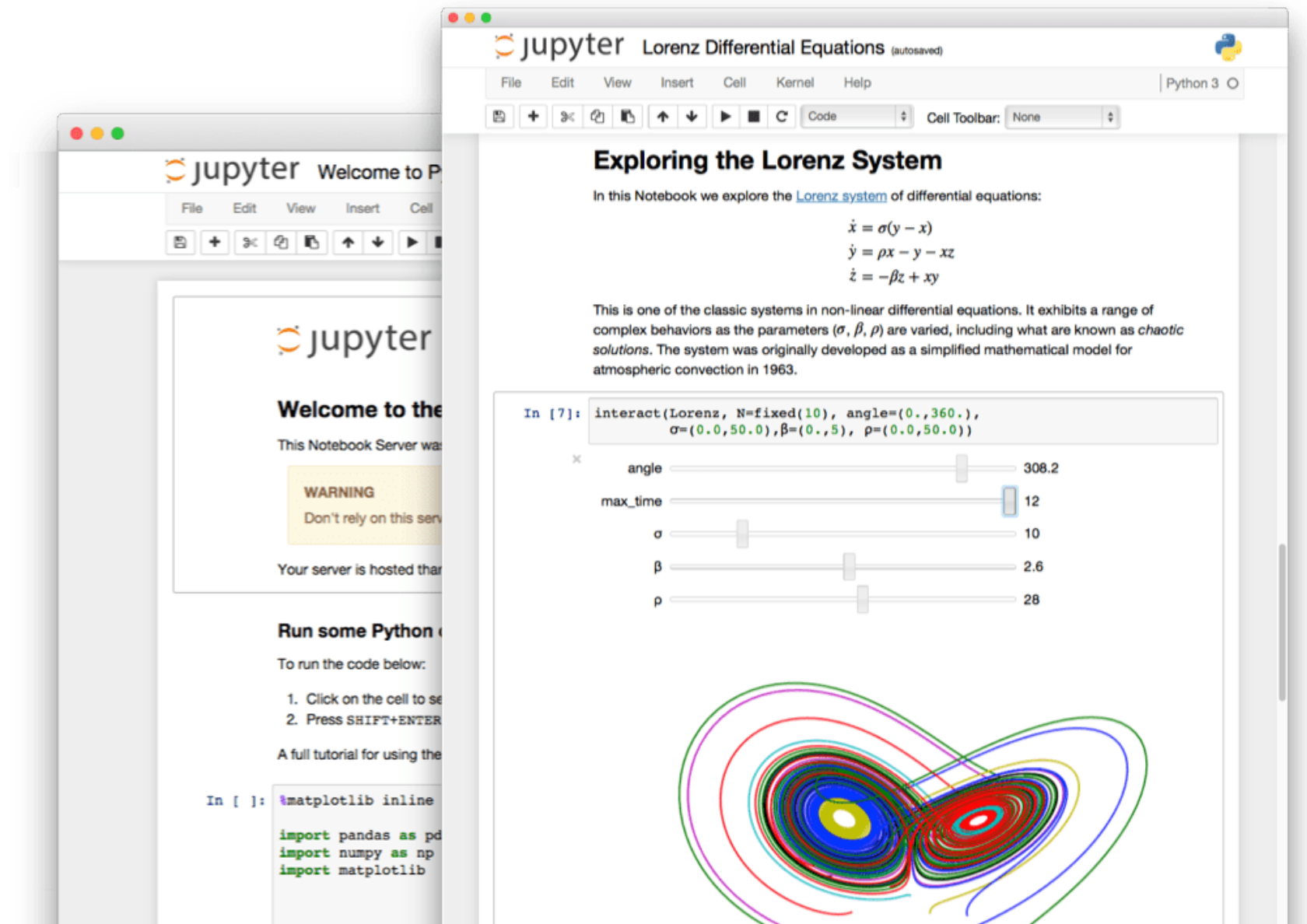
```
1 upstream = ['clean-users', 'clean-actions']
2 product = {
3     'nb': 'output/model-evaluation.ipynb',
4     'model': 'output/model.pickle'
5 }
```

```
1 # Injected cell
2 upstream = {
3     'clean-users': {
4         'nb': 'output/clean-users.ipynb',
5         'data': 'output/clean-users.parquet'
6     },
7     'clean-actions': {
8         'nb': 'output/clean-actions.ipynb',
9         'data': 'output/clean-actions.parquet'
10    }
11 }
```

```
1 users = pd.read_parquet(
2     upstream['clean-users']['data'])
3 actions = pd.read_parquet(
4     upstream['clean-actions']['data'])
```

Tasks as notebooks

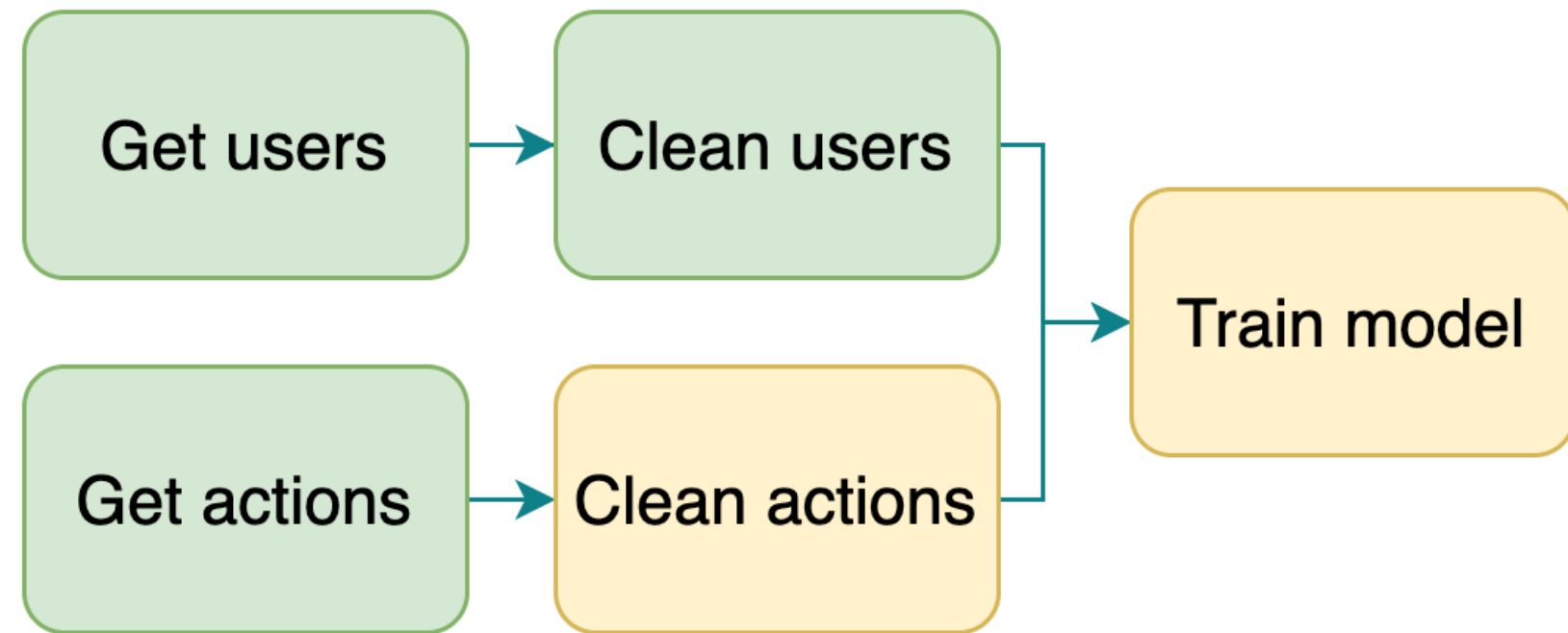
- Tasks are converted to notebooks (.ipynb) before execution²
- No need to write code to save tables or charts
- Generate logs for each execution



² Thanks to papermill! <https://github.com/nteract/papermill>

Build process

1. **Extract** upstream and product
2. **Determine execution order** using extracted upstream dependencies
3. **Inject code.** Replace the list of dependency names with a map of names to products
4. **Execute** the pipeline (skips up-to-date tasks)



Sample workflow

1. Input

```
get-users.py  
clean-users.py  
get-actions.py  
clean-actions.py  
train-model.py
```

2. Execute

```
ploomber build --entry-point .
```

3. Result

```
output/  
    get-users.ipynb  
    users.parquet  
  
    clean-users.ipynb  
    clean-users.parquet  
  
    get-actions.ipynb  
    actions.parquet  
  
    clean-actions.ipynb  
    clean-actions.parquet  
  
    model-evaluation.ipynb  
    model.pickle
```


Integration with Jupyter

- Scripts open as notebooks³
- (Temporary) cell injection
- Ensures no hidden state by running `ploomber build`



³ Thanks to jupyter! <https://github.com/mwouts/jupytex>

Demo

Main features

- Python, R and SQL⁴ supported⁵
- Parametrized pipelines with auto-generated CLI
- Testing (i.e. run data tests upon task execution)
- Debugger integration (pdb and ipdb)

⁴ Any database with a DBAPI 2.0 client or compatible with SQLAlchemy is supported

⁵ Languages that have Jupyter kernels are supported with a few limitations

Resources

- Install: `pip install ploomber`
- For updates/questions/feedback, follow me on Twitter: [@edublancas](https://twitter.com/edublancas)
- Code: github.com/ploomber/ploomber
- Examples: github.com/ploomber/projects
- Website: ploomber.io
- This presentation: blancas.io/talks/jupytercon-20.pdf

Thanks for watching!