



University of Pittsburgh

# *Tutorial*

## *Workflow Systems: OCCAM*

*Big picture, workflows, examples, demo*

### **Bruce Childers**

Associate Dean for Strategic Initiatives

Professor of Computer Science

School of Computing & Information

[childers@pitt.edu](mailto:childers@pitt.edu)

<http://www.cs.pitt.edu/~childers>



*If a tree falls in a forest and no one is around to hear it, does it make a sound?*

*If you create a model and use it in an experiment, and nobody can access it (& repeat!), does it exist?*

**nature**

NATURE | NEWS FEATURE

1,500 scientists lift the lid on reproducibility

Survey sheds light on the 'crisis' rocking research.

70% failed replication trial (50% of their own!); 45% due to data/code, 80% contributed



**SCIENTIFIC AMERICAN**

Why Should Scientific Results Be Reproducible?

Climate Science Can Be More Transparent, Researchers Say

open data

Virginia Gewin

Nature 529, 117–119 (07 January 2016)

COMMUNICATIONS OF THE ACM

HOME | CURRENT ISSUE | NEWS | BLOGS | OPINION | RESEARCH

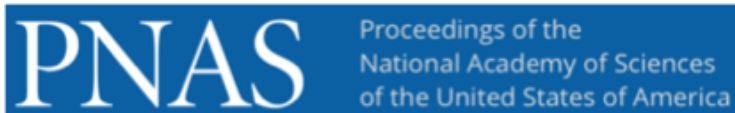
Home / Magazine Archive / March 2016 (Vol. 59, No. 3) / Repeatability in Computer Systems Research / Full Text

CONTRIBUTED ARTICLES

Repeatability in Computer Systems Research



Weak repeatability for 32.3% of 601 papers in computer systems research



Cluster failure: Why fMRI inferences for spatial extent have inflated false-positive rates

Anders Eklund, Thomas E. Nichols, and Hans Knutsson



Why Most Pub

John P. A. Ioannidis

# ~~Model~~ Path to the Emerald City



## I. Community and culture

- Galvanize, **incentivize**, & **educate**
- Review, mandates, **funding**, expectations
- Governance, policy & procedures
- Work with **existing practices to push boundaries**

## II. Technology

- Ease effort to **accelerate research**
- **Leverage** for higher quality, new discovery
- **Path of least resistance** is “right one”
- Work with **existing methods to push boundaries**

# Technology: Workflow Systems

## A place to do experiments

- Computational experiments (model + data + compute)
- Setup your own, or use one that is available

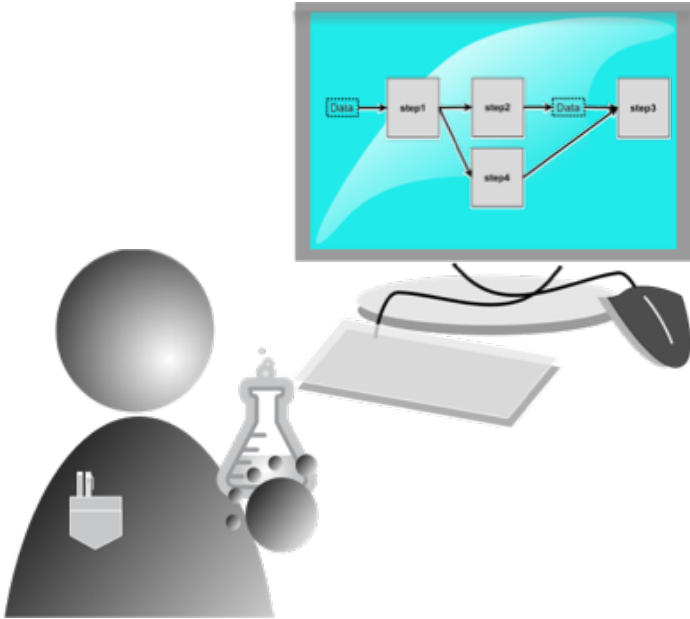
## Accelerate your research

- Focus on research rather than the infrastructure
- Leverage shared data, models, & other resources
- Capture & document experiment (FAIR)

# Workflow System

# Workflow System

## (1) Capture



### **Front-end User Interaction**

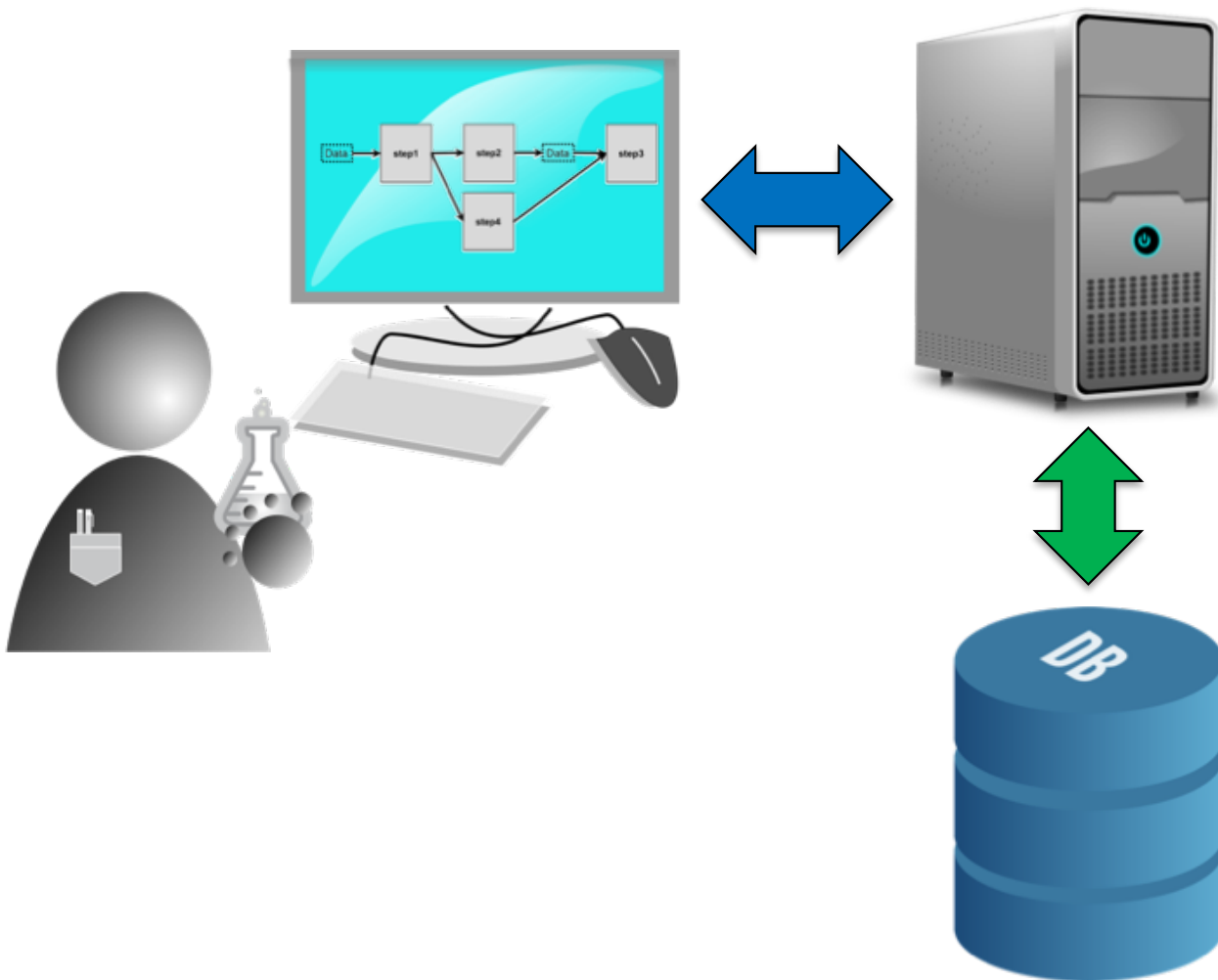
- Create & edit workflow
- View & manipulate results
- User workspace



# Workflow System

(1) Capture

(2) Engine

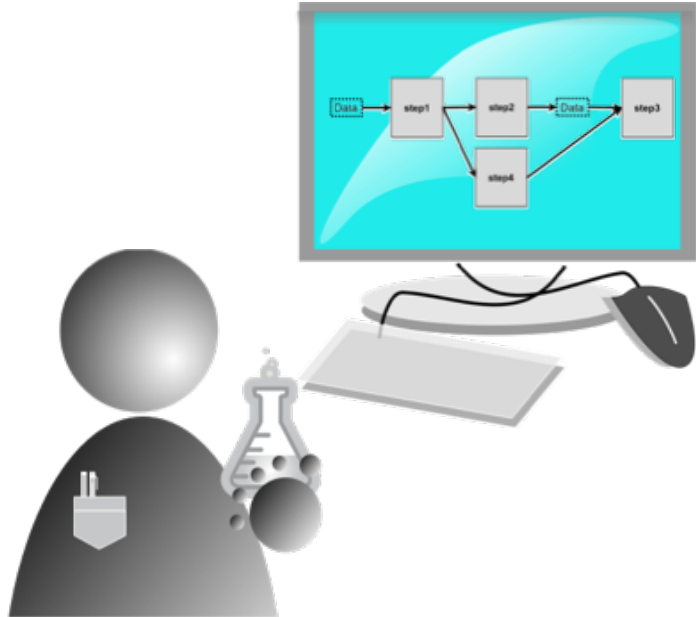


## Middle-End Management

- User interface
- Generate & dispatch
- Repository (database)

# Workflow System

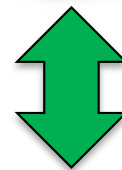
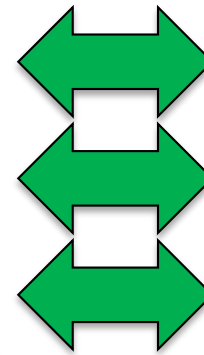
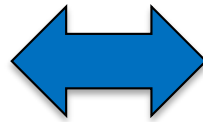
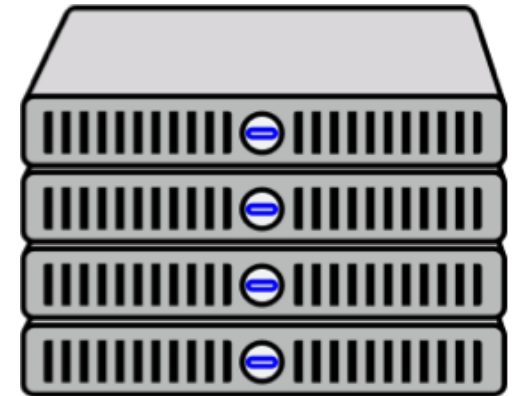
(1) Capture



(2) Engine



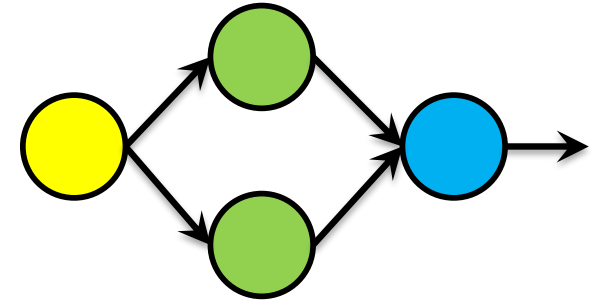
(3) Workers



## Back-End Compute

- Model/simulation workers
- Compute servers
- Local, private HPC, cloud

# Workflow



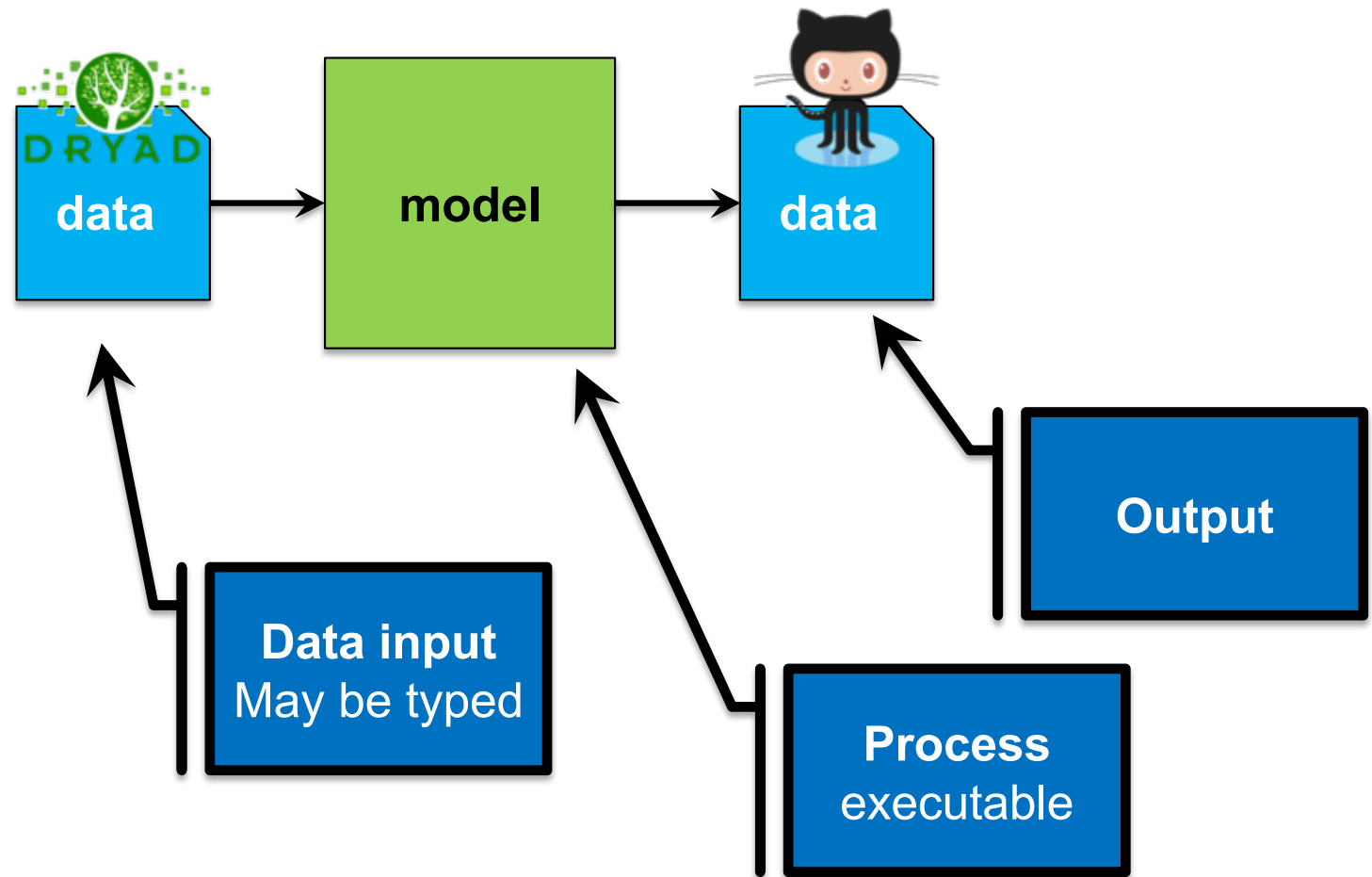
Sequence of computational steps

- **Directed acyclic graph** (DAG)
- **Nodes:** Process (executes), data (input/output)
- **Edges:** data flow of one step to the next
- **Operators:** of the language (e.g., transformation)

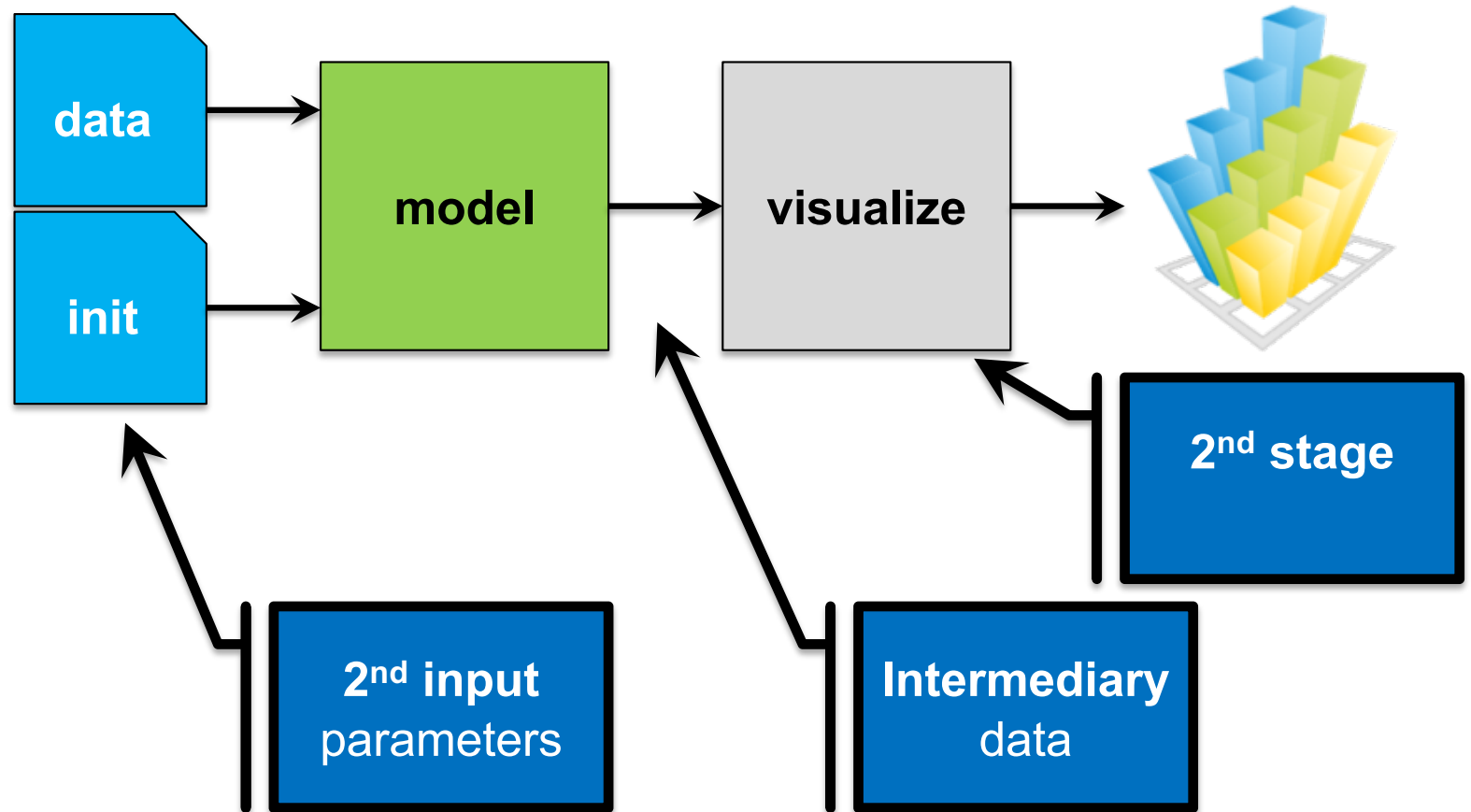
By itself, ***workflow is a specification***

- Realized: Program, script, *visual language*
- Represents experiment structure, data, steps

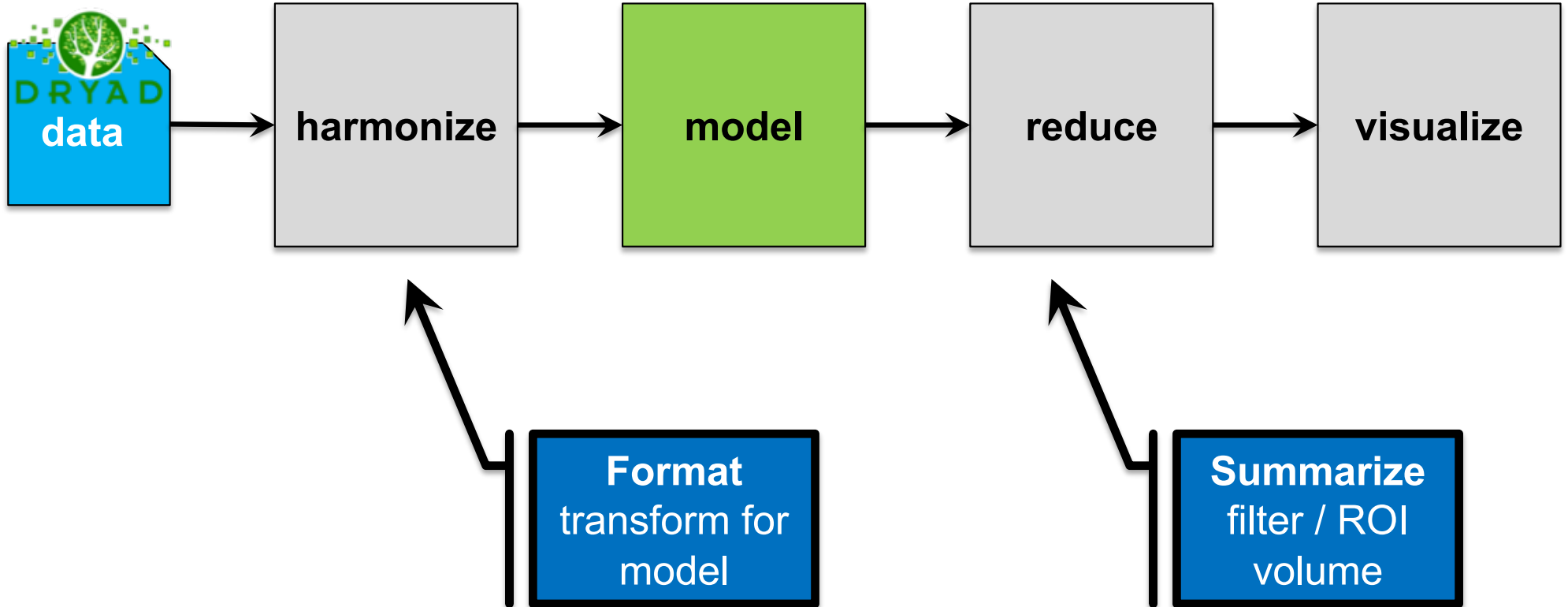
# Workflow Pattern: Process



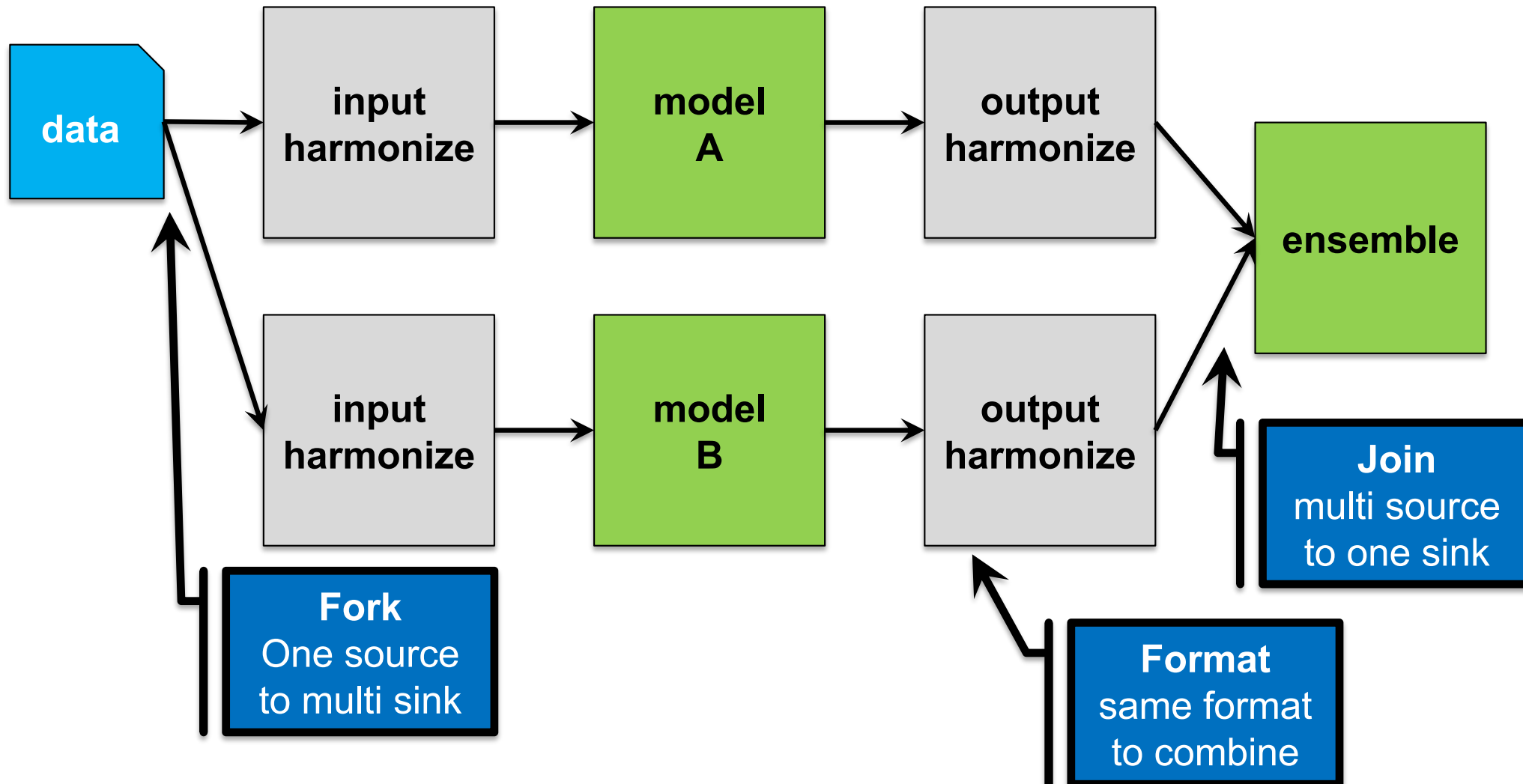
# Workflow Pattern: Pipeline



# Workflow Pattern: Pipeline



# Workflow Pattern: Fork & Join



# Capture: How to Specify

- Text

Specification  
print.cwl

```
cwlVersion: cwl:draft-3
class: CommandLineTool
baseCommand: echo
inputs:
  - id: message
    type:string
    inputBinding:
      position: 1
outputs: []
```

Data Input  
print.cwl

```
message: Hello World Modelers!
```

Execution  
Output

```
$cwl-runner print.cwl msg.yml
Hello World Modelers!
```

*Credit: CWL 3.0 hello world example*



# Capture: How to Specify

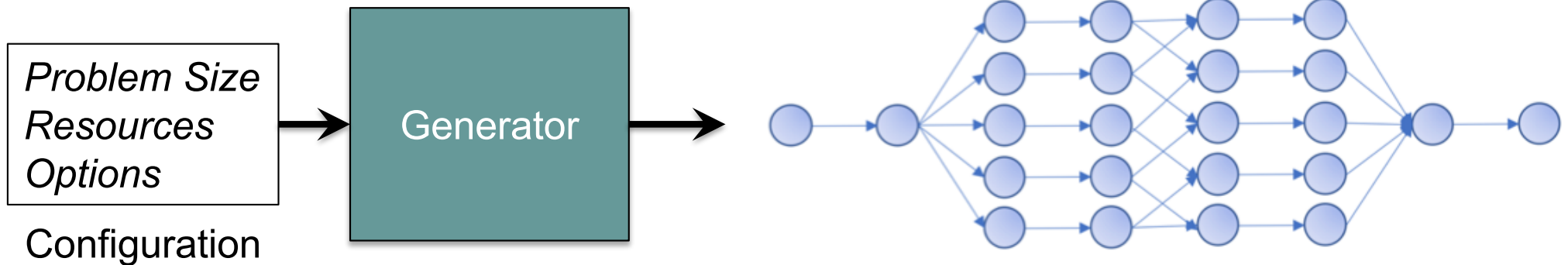
- Text
- Draw (visual)

The screenshot shows a web browser window titled 'exp100 [experiment] - OCCAM'. The address bar shows the URL: <https://occam.cs.pitt.edu/worksets/a9e6c084-b5bd-11e7-9e8d-000af7451cc2/experiments/b03cfb2e-b5bd-11e7-9cc1-00...>. The page header includes 'exp4 >', 'EXPERIMENT exp100', and utility buttons for 'bookmark', 'fork', and 'history'. Below the header is an 'Edit Description' link. A navigation bar contains 'Workflow', 'Run', 'Details', 'Files', and 'Output'. The main content area is split into two panels. The left panel, titled 'Select Object', contains a 'Type' section with a 'filter by type' input, a 'select object' button, and an 'attach' button. Below this is a 'Recently Used Objects' section showing 'None'. The right panel displays a workflow diagram on a grid. The workflow consists of several nodes: 'DSMS Aqslas2.0', 'trace-generator TracerPin', 'trace generated', 'simulator HMMSim', and 'application/json generated'. A new node, represented by a cube icon, is being added to the workflow and is highlighted with a green border. An 'Attach' button is visible below this new node.

**Editor**  
Drawing Workflow

# Capture: How to Specify

- Text
- Draw (visual)
- Generate



Example generators:

- ❑ Montage: Image stitching for astronomy
- ❑ CyberShake: Earthquake modeling for So. California
- ❑ LIGO: Analyze gravitational waveforms

# Capture: How to Specify

- Text
- Draw (visual)
- Generate
- Program

```
input_file <- "data/data.csv"
output_file <- "data/results.csv"

# read input
input_data <- read.csv(input_file)
# get number of samples in data
sample_number <- nrow(input_data)
# generate results
results <- some_other_function(input_file, sample_number)
# write results
write.table(results, results_file)
```

# Workflow Execution

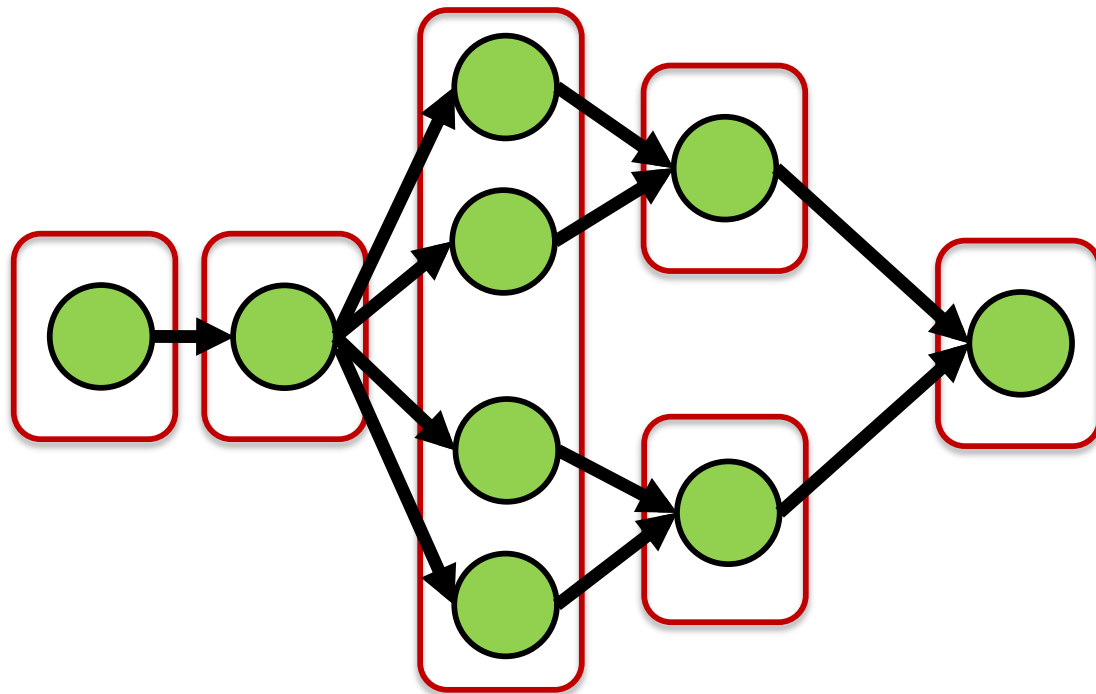
Need to run (i.e., “conduct experiment”)

- Workflow (few steps to 100Ks of steps)
- Actual input and configuration values
- Resources (e.g., machines)

***Workflow engine*** handles the “magic”

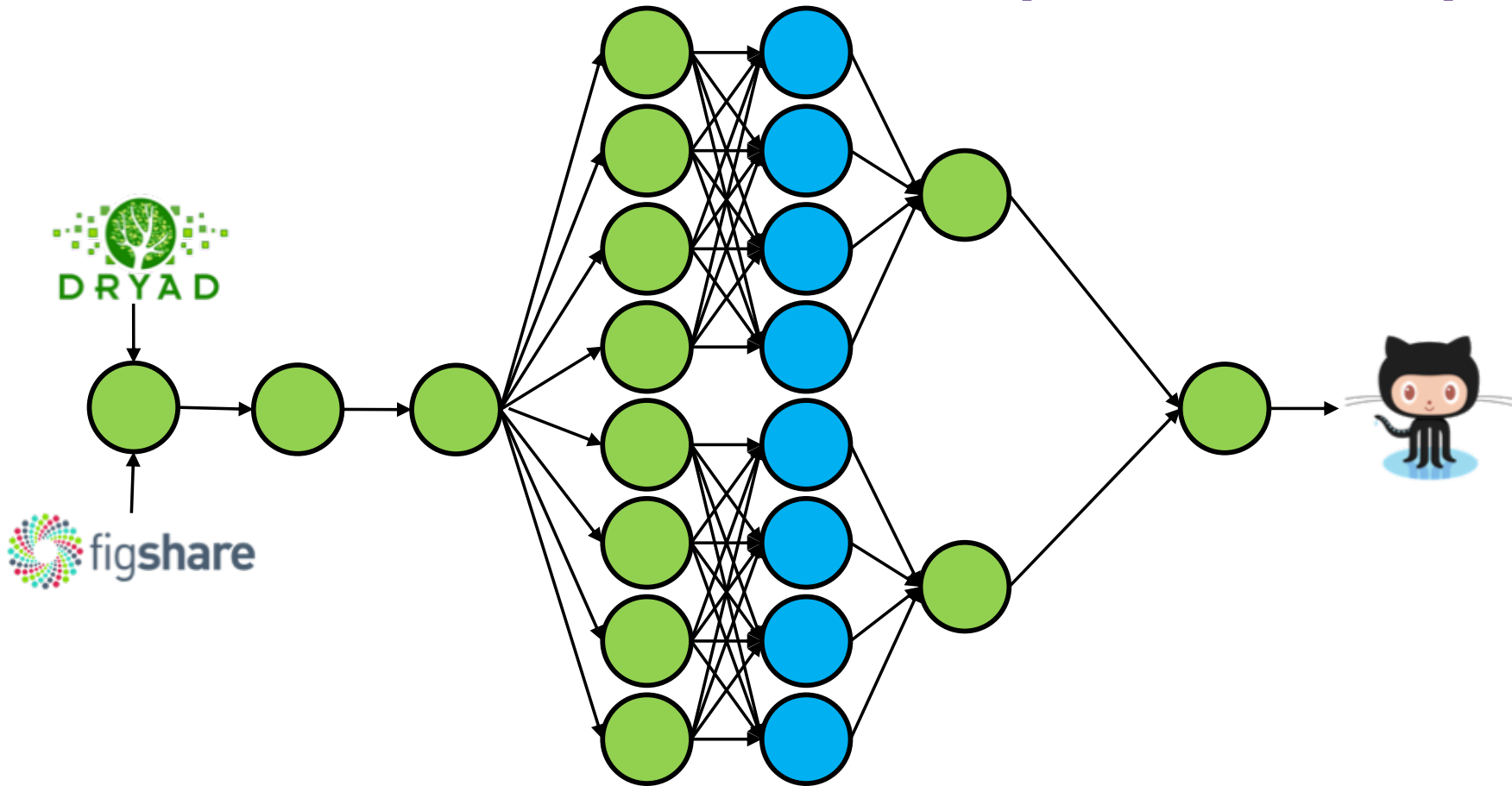
- Compute scheduling: dispatch processes to resources
- Data orchestration: input, intermediary, result
- Management: Monitoring, logging, error recovery
- Metadata for experiment, workflow preserved

# Workflow Execution (Simple)

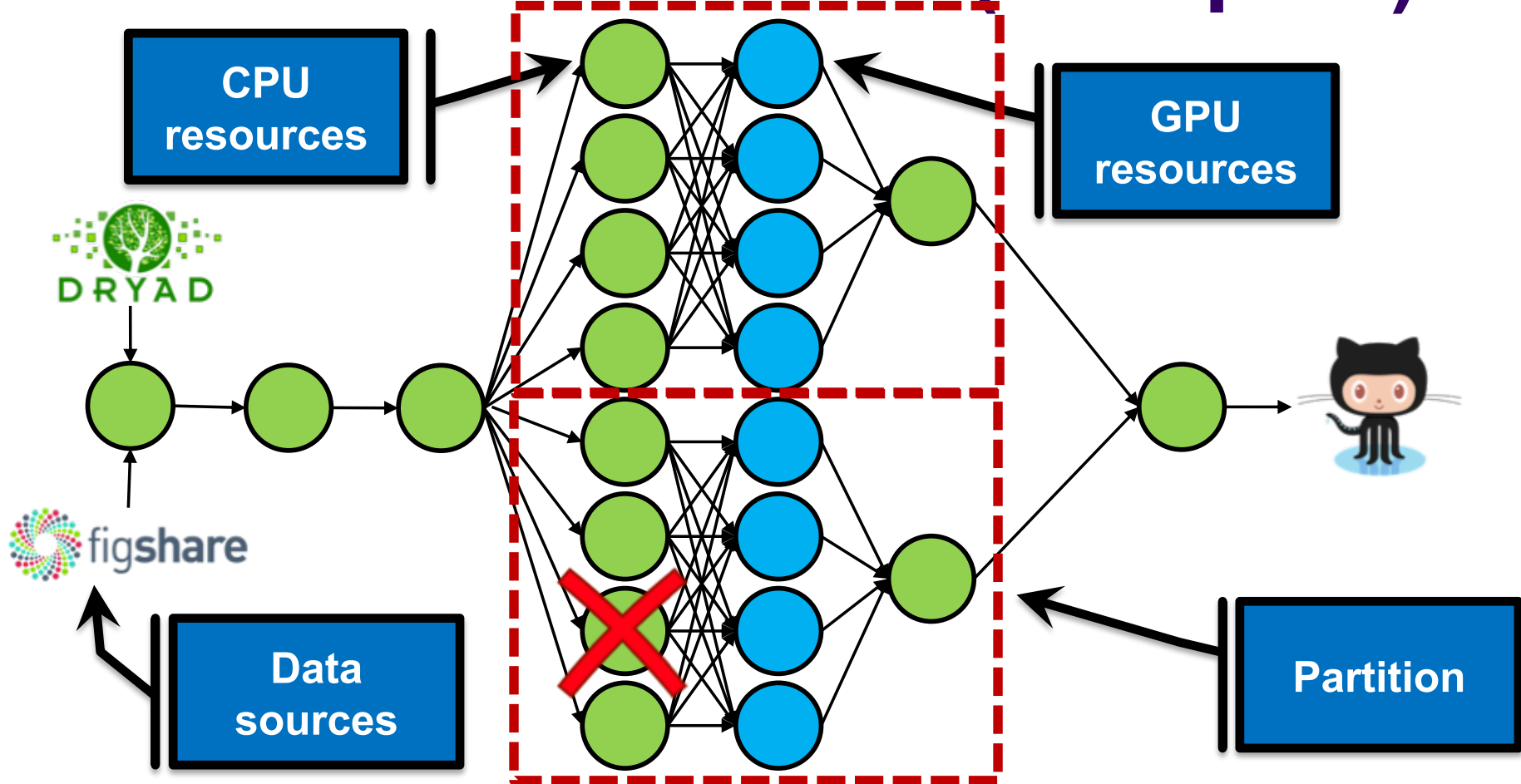


Simple scheduling on non-distributed resources  
Order imposed by workflow  
Inputs and resource available to execute step

# Workflow Execution (Complex)



# Workflow Execution (Complex)



Data sources and sinks

Grouping computational steps & allocate resources

Failure and recovery

# Some Factors to Consider

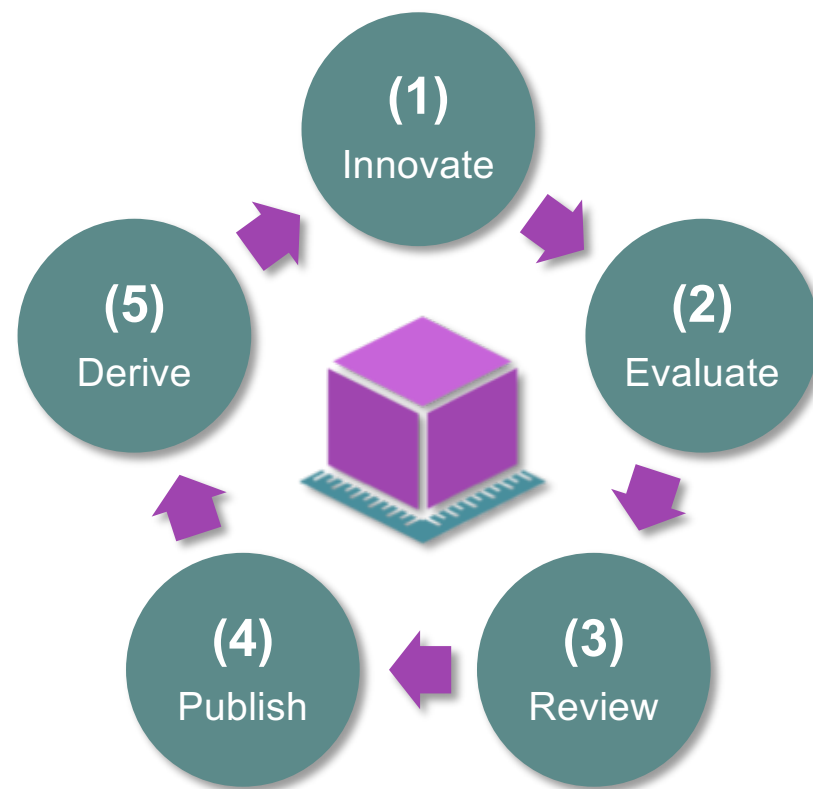
- Science domain
- Interactivity (vs. generated)
- Scalability
- Resources (e.g., laptop vs. distributed cloud)
- Collaboration (e.g., workspaces, identifiers)
- Metadata & preservation



| System    | Application         | Notable features                       |
|-----------|---------------------|--|
| Galaxy    | Genomics            | Commons, HPC, repeatability            |
| Pegasus   | General science     | Data-centric, scalability, community   |
| Kepler    | Astronomy, others   | Derived embedded systems modeling      |
| Traverna  | Biomedicine, others | Workbench, provenance                  |
| VisTrails | General science     | Provenance, visualization, exploration |
| CK        | Machine learning    | Packaging, comparison, auto-tuning     |
| ReproZip  | Data science        | Record & replay for repeatability      |
| Pachyderm | General science     | Repeatability, provenance, commercial  |
| CodeOcean | General science     | Publication, editing & run, commercial |

# OCCAM: *Experiment System*

- Workflows
- Research life cycle
- Contributory & collaborate
- Preservation, provenance
- FAIR metadata



Dropbox

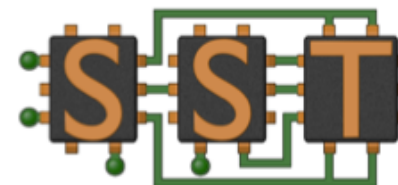


python™



git

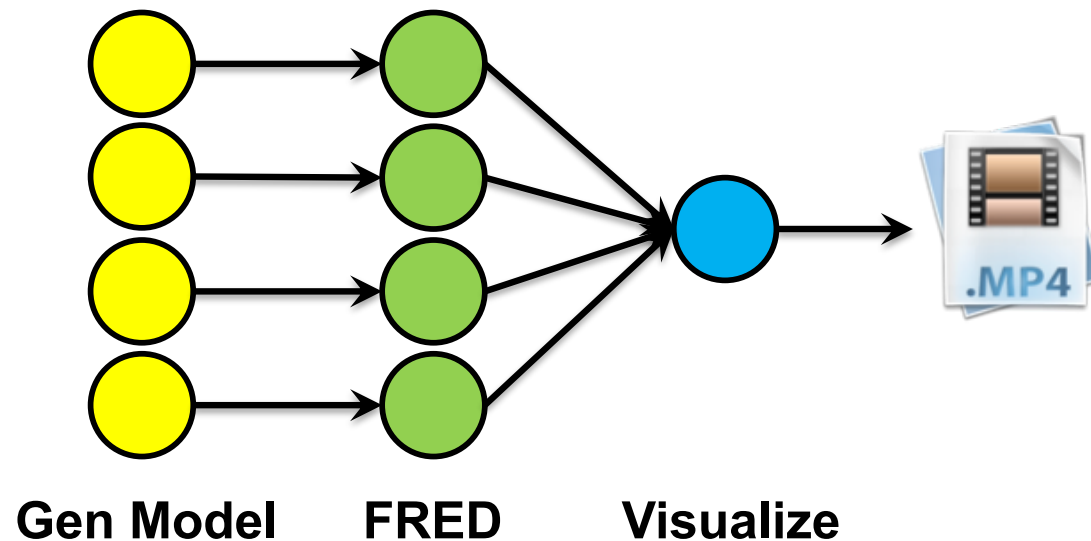
GitHub

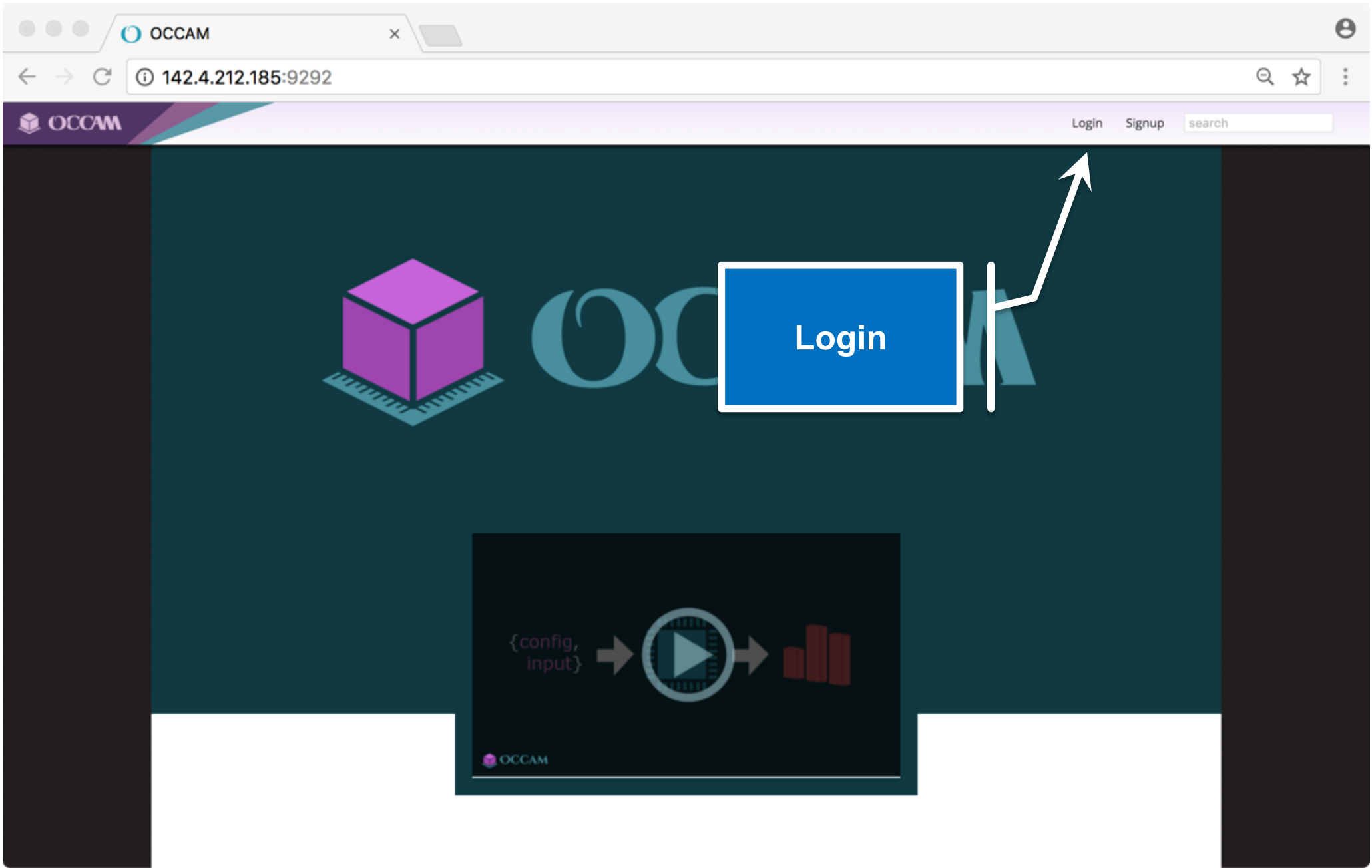


# Demo! #1

- Start of research lifecycle

- Creating, running & sharing an experiment
- Influenza model with FRED
- Visualize results



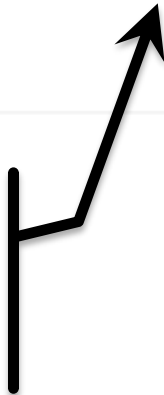




Add a New Object


test1 0 collaborators

**Workspace Experiments**



test1 [workset] - OCCAM

Not Secure | 142.4.212.185:9292/94f48412-5bdb-11e8-a324-4c72b9434eda/19c227b8e4f5ea7b1791be661389833774e5...

 WORKSET  
test1  
Childers

bookmark  
fork

Contents Details Files Output History Access

Authors

+ username add

Collaborators

+ username add

Directory ?

- fred1
- fred2

**New Experiment**

Add a New Object

**Prior Experiments**

 experiment Fred3 add

fred3 [experiment] - OCCAM x

142.4.212.185:9292/94f48412-5bdb-11e8-a324-4c72b9434eda/cdc4e1cf2fc20e6ab0f6b47831d7e087517de7de/2?link=8

OCCAM Childers search

test1

EXPERIMENT  
fred3  
Childers

bookmark  
fork

Workflow Contents Details Run Files Output History Access

Workflow Canvas

Select Object

Type

simulator

Name

FRED

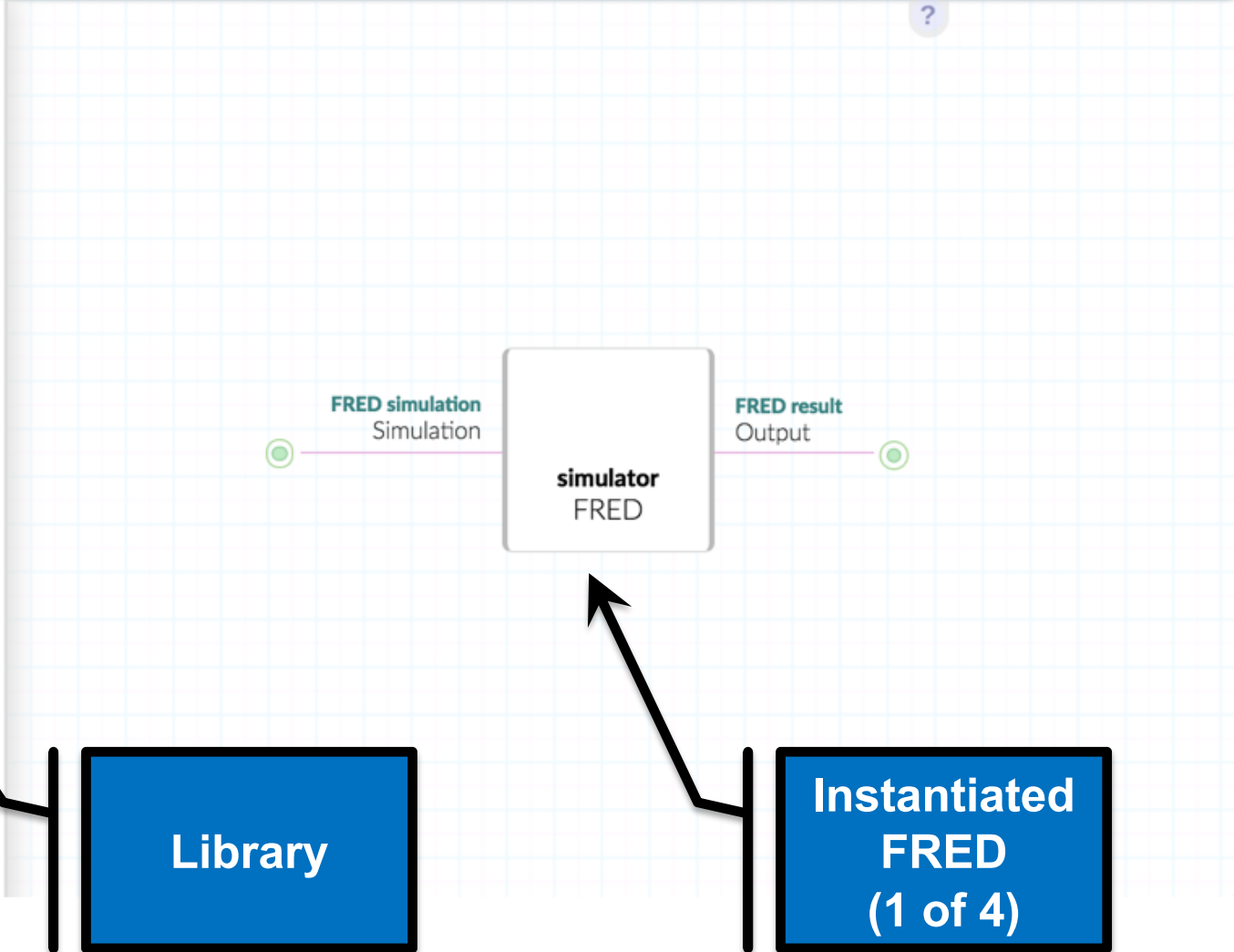
FRED simulation  
Simulation



FRED result  
Output

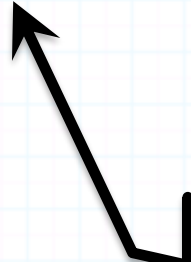
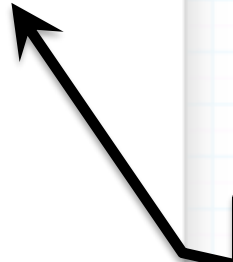
Workflow Templates  
Recently Used Objects

None

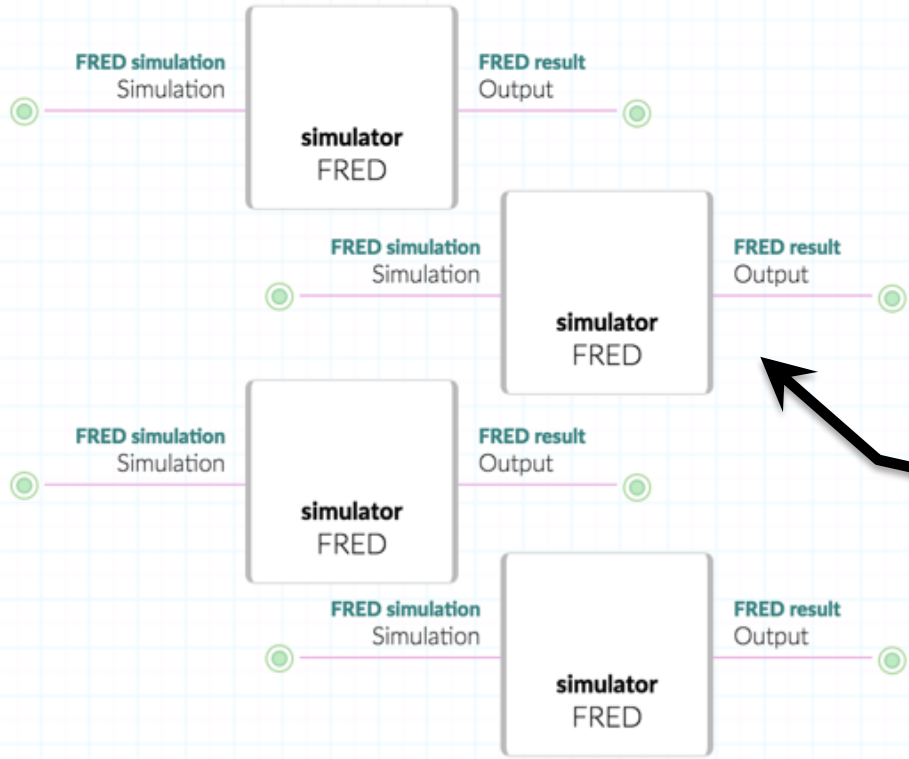


**Library**

**Instantiated  
FRED  
(1 of 4)**







**Instantiated  
FREDs  
(4 of 4)**

Select Object

Type

simulation

Name

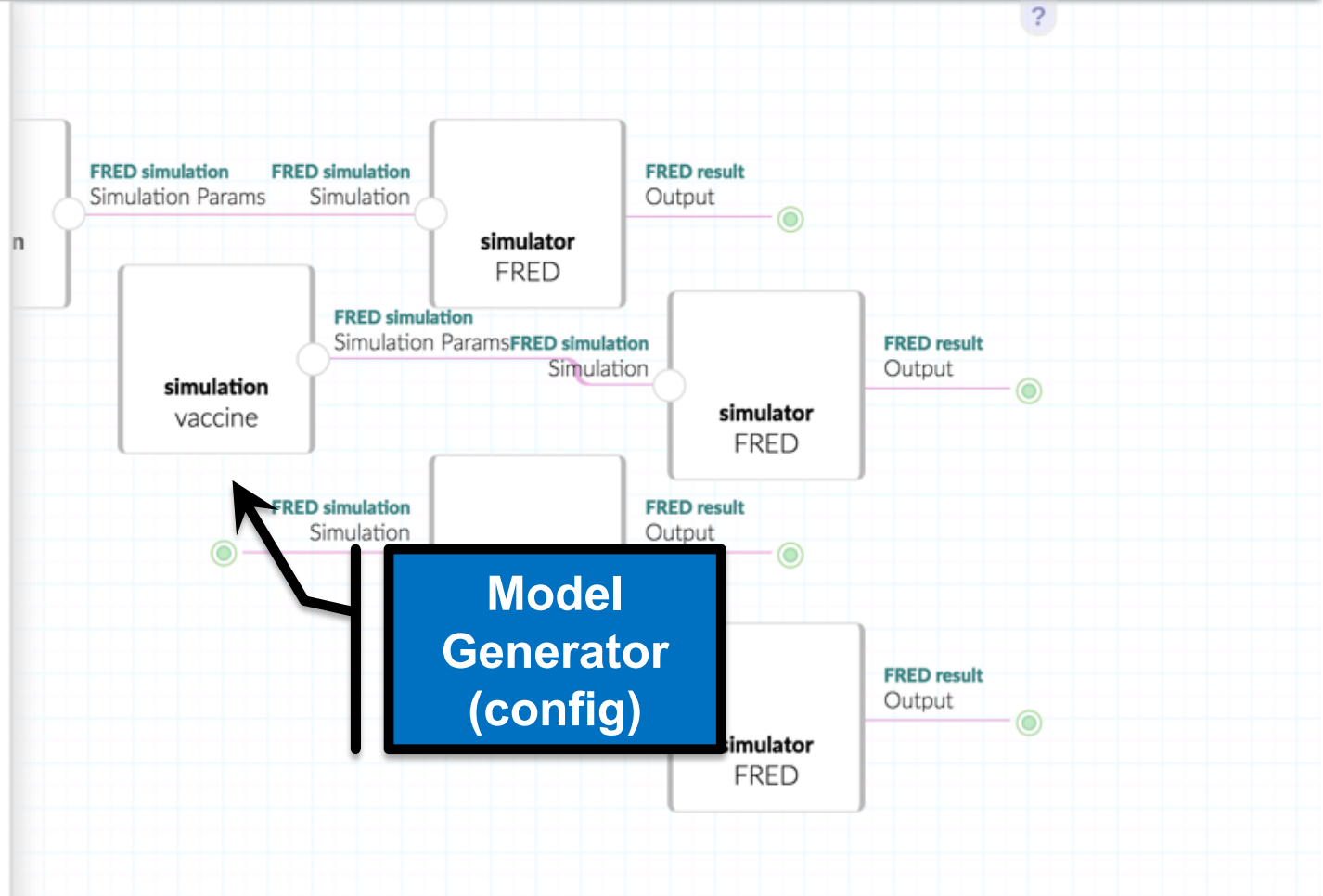
vaccine



FRED simulation Simulation Params

Workflow Templates  
Recently Used Objects

None



**Model Generator (config)**



FRED simulation Simulation Params FRED simulation Simulation



FRED result Output

Visualization (video gen)



FRED simulation Simulation Params FRED simulation Simulation



FRED result Output



FRED simulation Simulation Params FRED simulation Simulation



FRED result Output

FRED result North-East FRED result North-West FRED result South-East FRED result South-West



video Stitched Vic

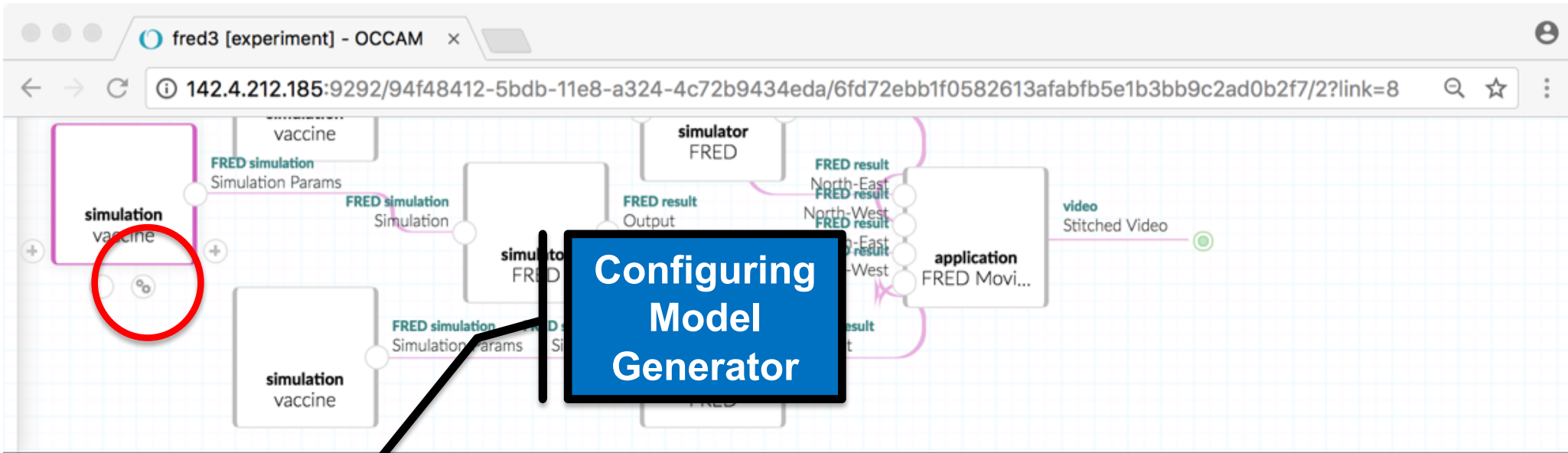


FRED simulation Simulation Params FRED simulation Simulation



FRED result Output

FRED result Output



Configuration

FRED | FRED | FRED | FRED | vaccine | vaccine | **vaccine** | vaccine | FRED

Simulation Vaccine **Influenza** Vaccine Events

Influenza model

**Household Confinement**

▼ Influenza model

▼ State - Infected and Symptomatic

Visualize (+)

Probability of household confinement (+)

0.3

fred3 [experiment] - OCCAM x

142.4.212.185:9292/94f48412-5bdb-11e8-a324-4c72b9434eda/6fd72ebb1f0582613afabfb5e1b3bb9c2ad0b2f7/2?link=8

Simulation Vaccine Influenza **Vaccine Events**

Vaccines become available

▼ Vaccines become available

|   |                              |                                  |
|---|------------------------------|----------------------------------|
| Start of event [+]                                  | <b>Candidates Vaccinated</b> | <input type="text" value="10"/>  |
| End of event [+]                                    |                              | <input type="text" value="10"/>  |
| Maximum number of events [+]                        |                              | <input type="text" value="0"/>   |
| Percentage of candidates vaccinated [+]             |                              | <input type="text" value="0.2"/> |
| Center latitude of circle of exposure [+]           |                              | <input type="text" value="0"/>   |
| Center longitude of circle of exposure [+]          |                              | <input type="text" value="0"/>   |
| Radius of exposure (km) [+]                         |                              | <input type="text" value="0"/>   |
| Fips code of affected region (zero - disabled). [+] |                              | <input type="text" value="0"/>   |
| Minimum age [+]                                     |                              | <input type="text" value="0"/>   |
| Maximum age [+]                                     |                              | <input type="text" value="999"/> |



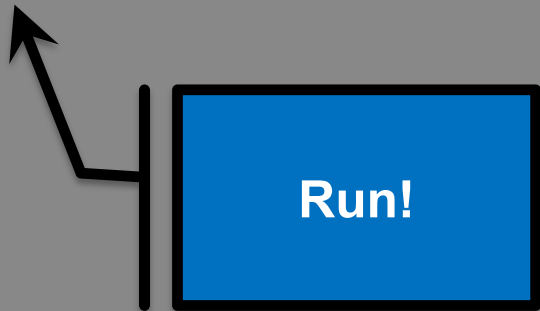
# fred3

Childers

bookmark  
fork

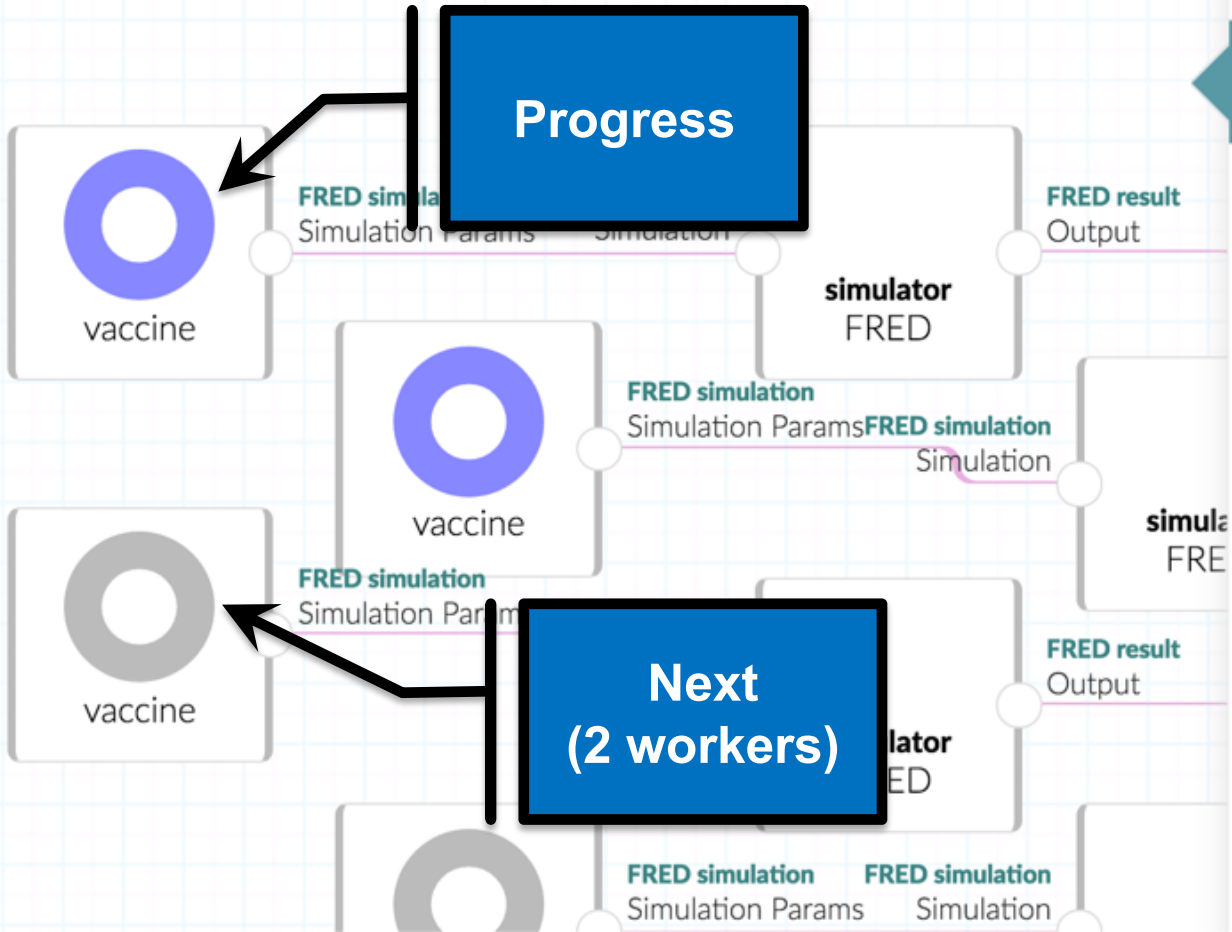
Workflow Contents Details **Run** Files Output History Access

Queue



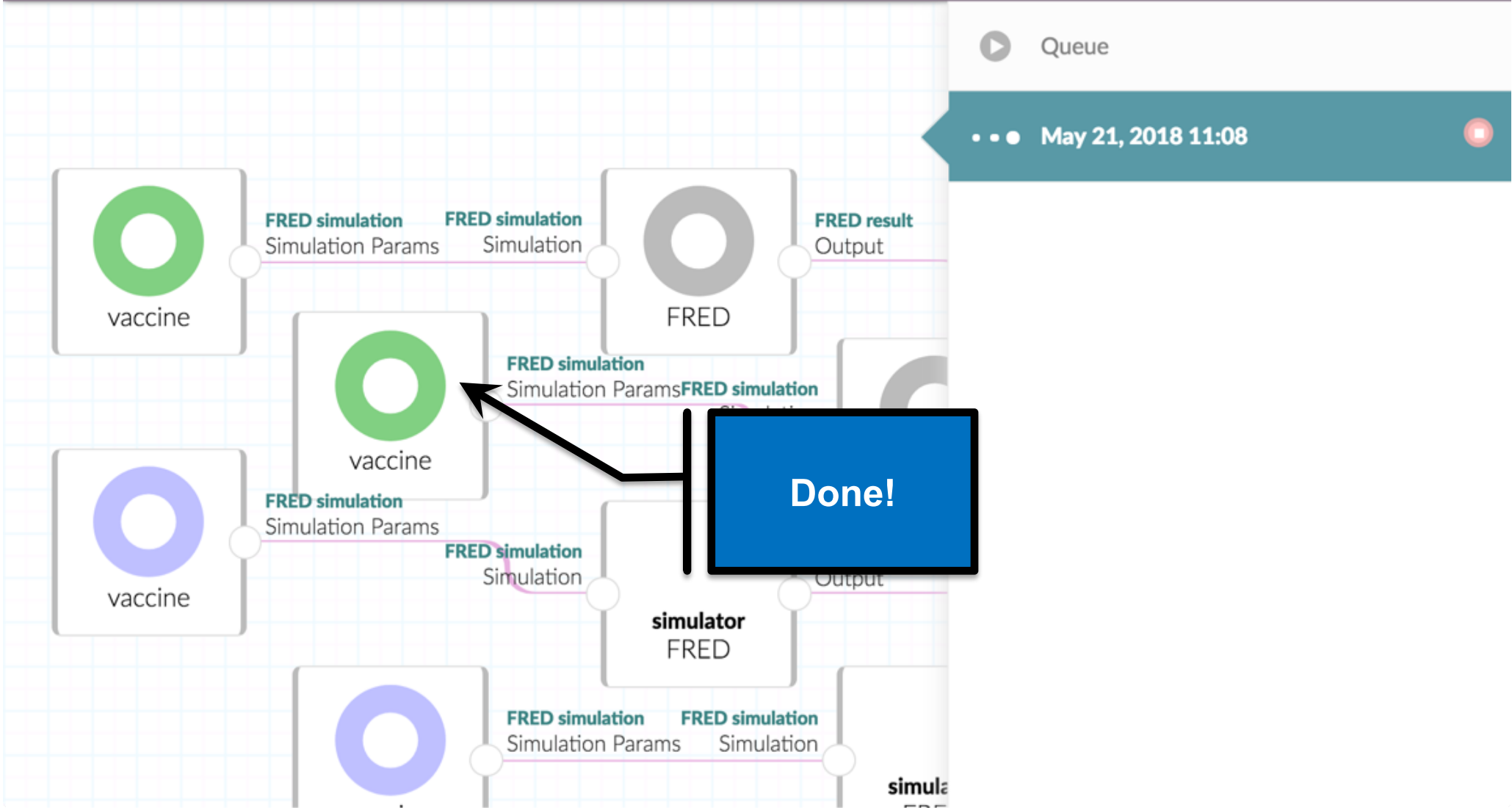
Dispatch To

Run

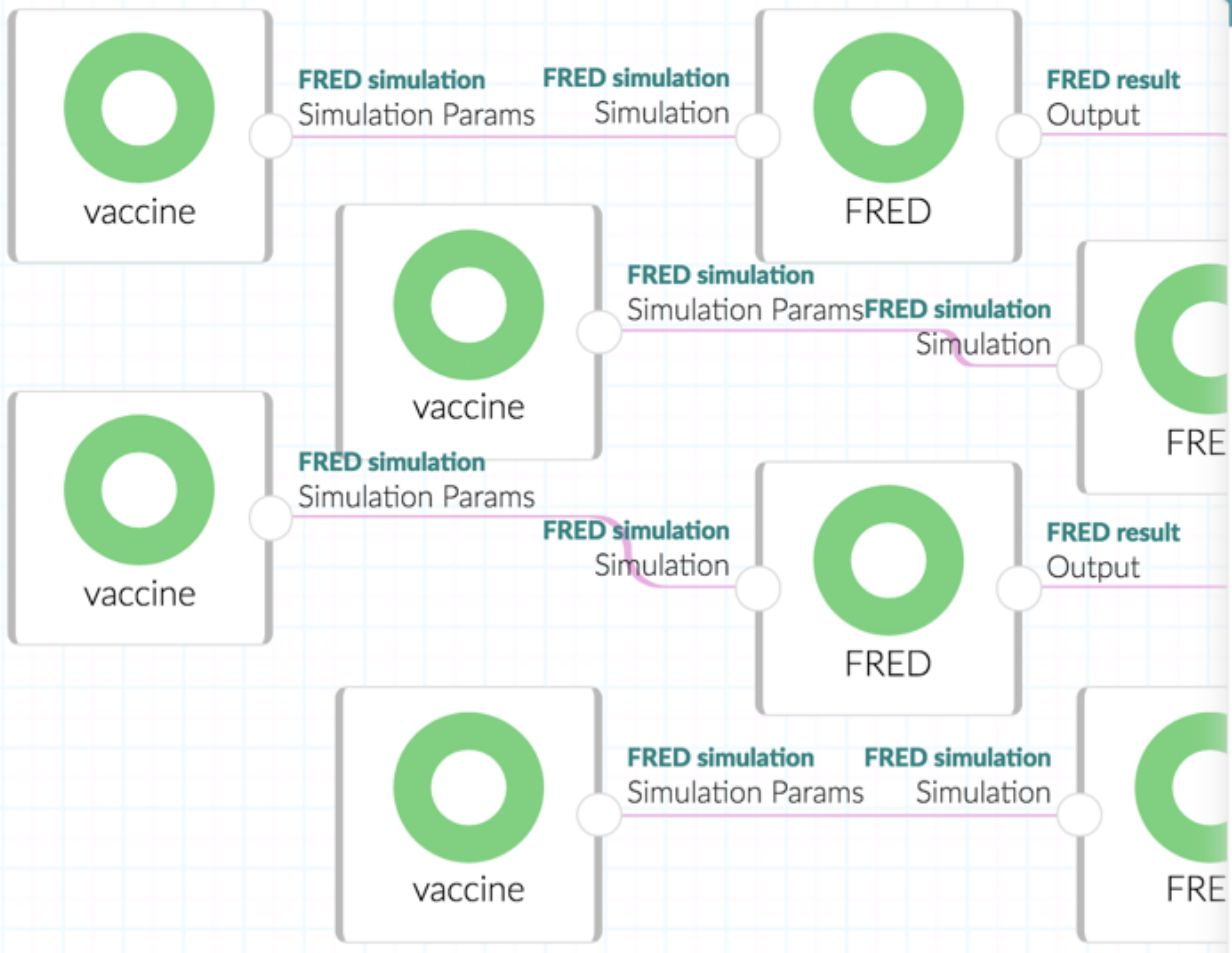


Queue

May 21, 2018 11:08







Queue

✓ May 21, 2018 11:08

Completed!

fred3 [experiment] - OCCAM x

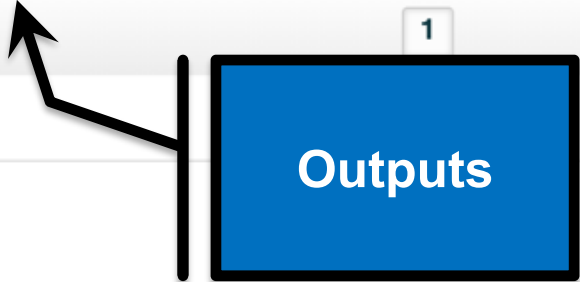
142.4.212.185:9292/94f48412-5bdb-11e8-a324-4c72b9434eda/6fd72ebb1f0582613afabfb5e1b3bb9c2ad0b2f7/2/output?link...

EXPERIMENT  
**fred3**  
Childers

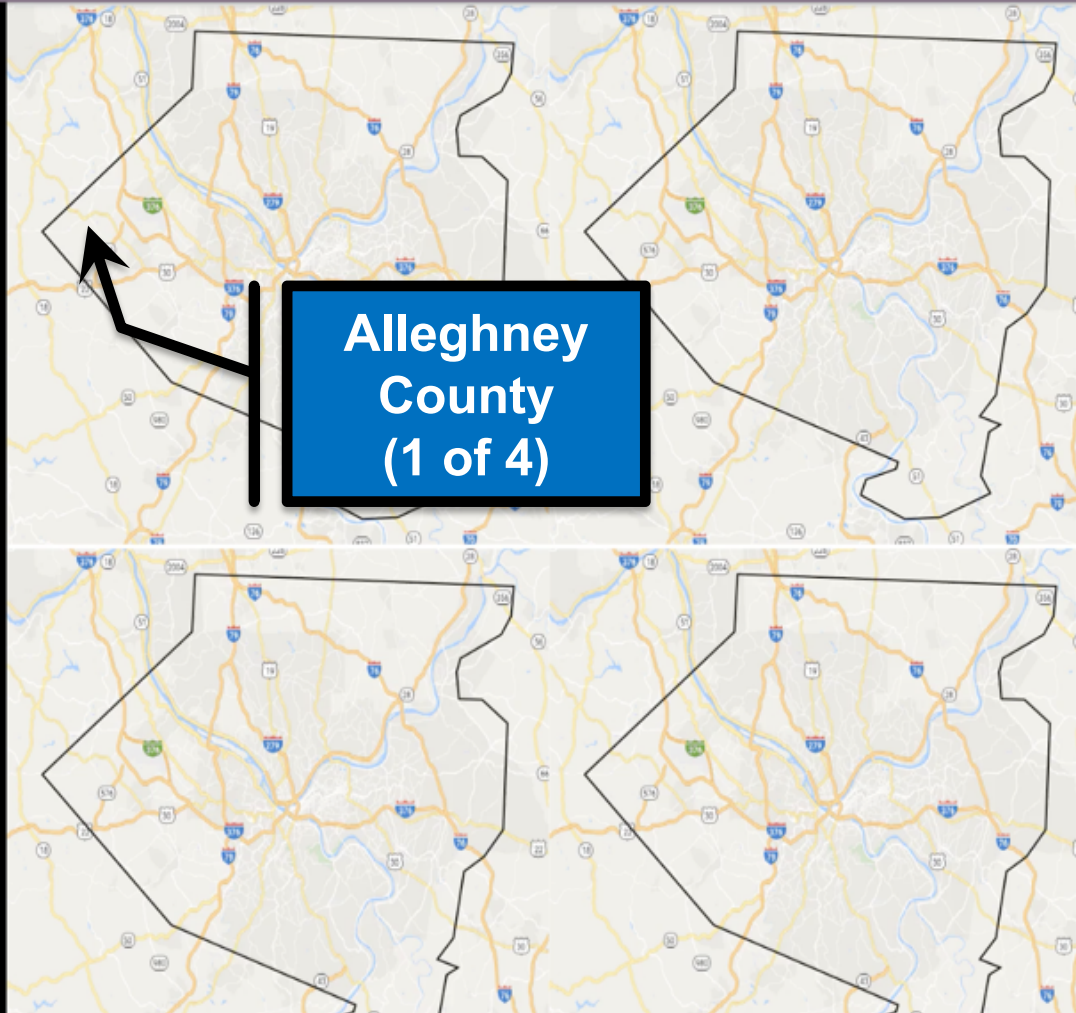
bookmark  
fork

Workflow Contents Details Run Files **Output** History Access

List View | Preview



- FRED simulation Simulation Params
- FRED simulation Simulation Params
- FRED simulation Simulation Params
- FRED simulation Simulation Params
- FRED result Output
- FRED result Output
- FRED result Output
- video Stitched Video

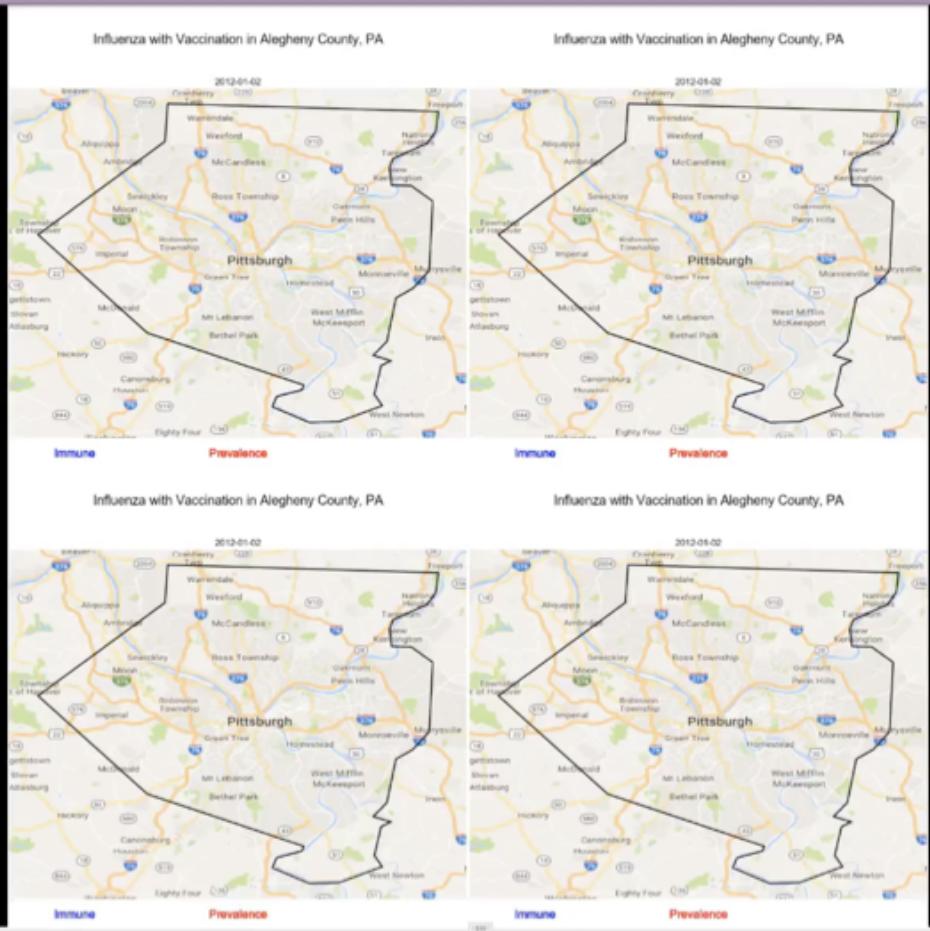




# VIDEO Influenza with Vaccination in Allegheny County, PA

bookmark  
fork

View Details Files Output History Access Admin



Using video.js

Open Options

# Summary

## Workflow systems!

- Accelerate research progress
- Leverage models, data, experiments
- Collaboration for modeling

Part of solution to “reproducibility crisis”

Collaboration, sharing, metadata as well

Community and culture are equally important