



Towards automated provenance collection for experimental runs of agent-based models

Doug Salt Gary Polhill and Corran Musk Lorenzo
Milazzo Dawn Parker

The James Hutton Institute

September 5, 2023



The James
Hutton
Institute

Structure of the presentation

- Abstract
- Motivation
- The standard approach
- The database
- The scope
- A simple workflow
- Some provenance
- So what?
- But so what?
- Maybe a use?
- What else?



Abstract

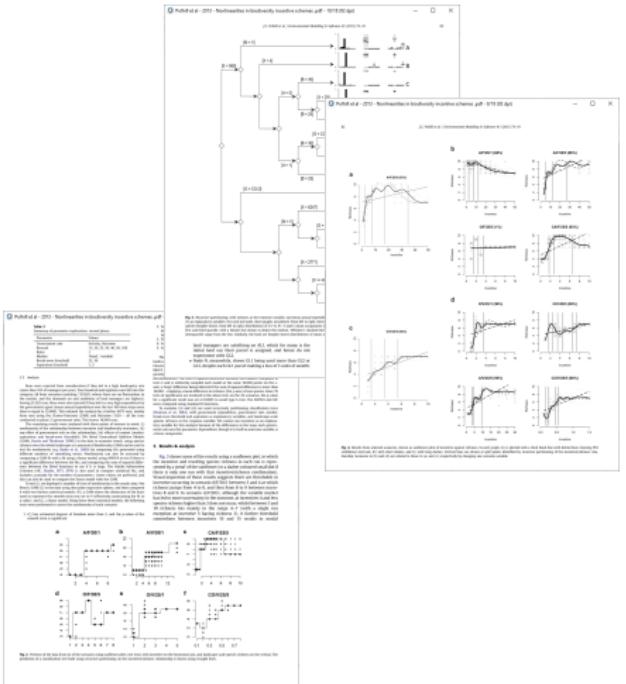
We demonstrate a working framework for the automatic recording of provenance and metadata for primarily agent-based models that could easily be adapted to the other modelling environments. We discuss the need for such a framework, the philosophy behind the design we adopted, the implementation, discuss the results and demonstrate a simple tool for tracing bad data through a provenance graph.



Motivation



The James
Hutton
Institute



The standard approach

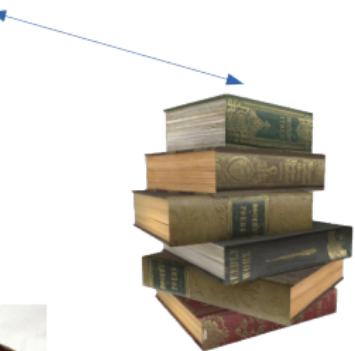
- Scripts
- Storage archives
- Diaries



How we did it



The James
Hutton
Institute



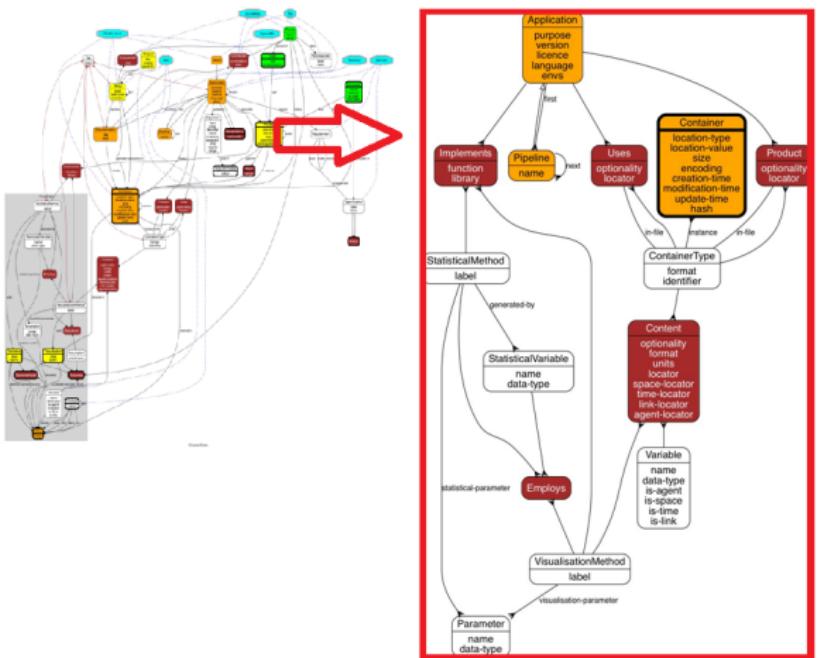
Python library



The database



The James
Hutton
Institute



Social Simulation REplication Interface





The scope

- **Analysis** - fine grain provenance pertaining to statistical and visualisation outputs.
- **Finegrain** - a provenance diagram down to the level of variables.
- **Folksonomy** - a diagram showing annotations against the database, produced and categorised at the discretion of the user doing the annotation
- **Project** - management metadata. Largest granularity of metadata supported
- **Provenance** - provenance diagram at the level of file and parameter
- **Services** - service provided and requirement description
- **Workflow** - the actual workflow



What was needed

Primitive	Type	Purpose
SSREPI_require_minimum	M	Lower bound on software/hardware required
SSREPI_require_exact	M	Exact bound on software/hardware required
SSREPI_application	P & M	Specifies some executable
SSREPI_me	P & M	Determines executable or reference to executable being run.
SSREPI_argument	P	An argument type to an executable
SSREPI_output	P	An output type from an executable
SSREPI_input	P	An input type for an executable
SSREPI_person	M	Provide metadata for a particular actor within this system
SSREPI_project	M	Specifies a project which contains all studies
SSREPI_study	M	A set of experiments makes up a single study
SSREPI_set	M	Sets the default licence and other metadata
SSREPI_involvement	M	Links personnel to a study
SSREPI_paper	M	A paper associated with this study
SSREPI_make_tag	M	Used for building a folksonomy
SSREPI_tag	M	Used to tag any entity with a folksonomy tag
SSREPI_contributor	M	A person with some kind of relation to an executable or script.
SSREPI_statistical_method	M	Record a statistical method
SSREPI_visualisation	M	Record a method to create an image to depict one or more values.
SSREPI_statistics	M	Record activities that populate the values of statistical variables.
SSREPI_visualisation_method	M	Methods for generating visualisations.
SSREPI_implements	M	Links a statistical or visualisation method to an application
SSREPI_parameter	M	Record the name taken by a statistical or visualisation method.
SSREPI_statistical_variable	M	A name for (one of) the result(s) of a statistical method.
SSREPI_visualisation_variable	P & M	Declares a named variable of interest
SSREPI_variable	M	Names a variable of interest
SSREPI_statistical_variable_value	P & M	Sets an actual value for a named statistical variable
SSREPI_value	P	Sets a value.
SSREPI_content	M	Links a kind of output/input/argument to a variable
SSREPI_person_makes_assumption	M	Links a person to an assumption

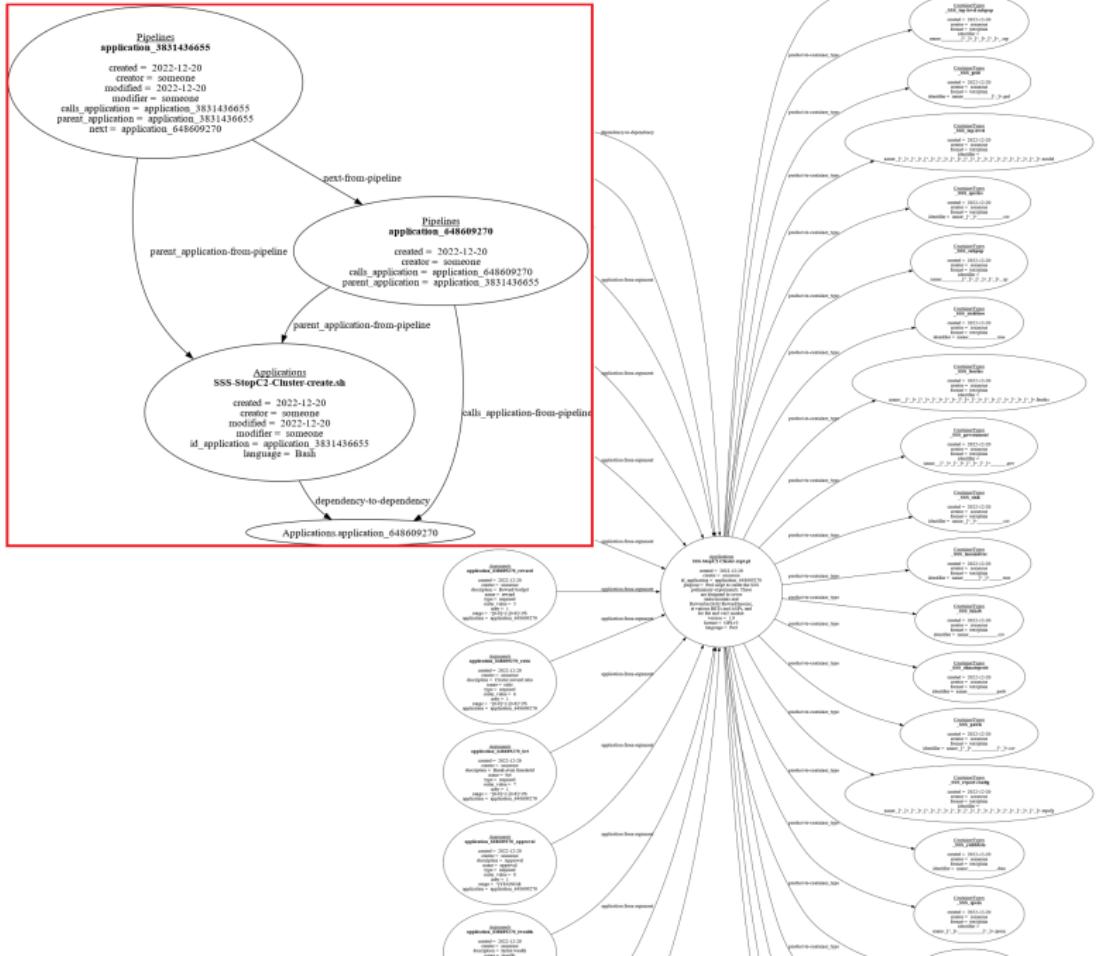


What was needed? (cont.)

We encoded the previous actions in Python and then then called these from Bash

- `create_database.py` - creates a database idempotently.
- `exists.py` - checks if a particular row in a table exists.
- `get_value.py` - gets any specified single value from a table given the primary key.
- `get_values.py` - gets one or more rows given the search values.
- `next_study.py` - gets the next available and unique study number.
- `update.py` - idempotently updates a particular row in a table.

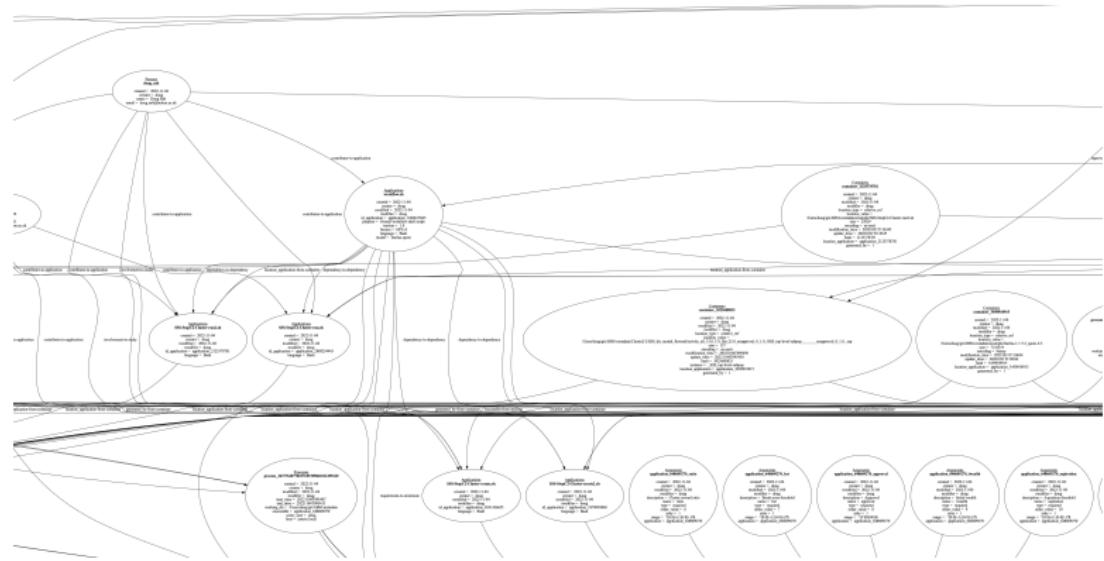




Some provenance!!!



The James
Hutton
Institute



So what?

It now takes me the following lines of code to reproduce an experiment

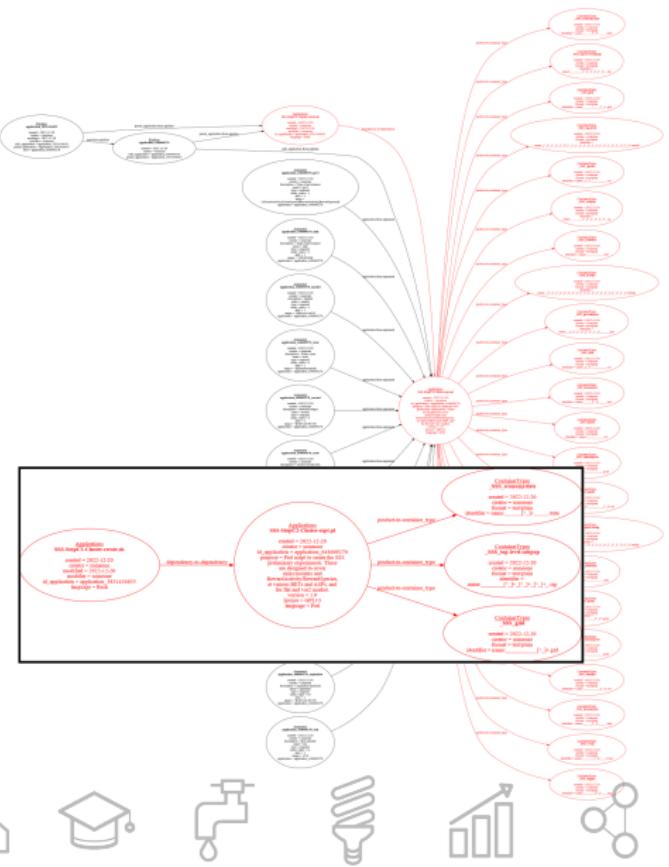
1. Create new directory
2. Enter new directory
3. `mkdir bin`
4. `cp /mnt/storage/doug/git/ABM-metadata/example/*.sh bin`
5. `cp /mnt/storage/doug/git/ABM-metadata/example/*.pl bin`
6. `cp /mnt/storage/doug/git/ABM-metadata/example/*.R bin`
7. `ln -s /mnt/storage/doug/git/ABM-metadata/cfg`
8. `ln -s /mnt/storage/doug/git/ABM-metadata/lib`
9. `ln -s /mnt/storage/doug/git/ABM-metadata/doc`
10. Amend `bin/path.sh` to `bin/path.directory_name.sh`
11. `source bin/path.directory_name.sh`
12. `psql`
13. `create database third`
14. `\q`
15. profit!



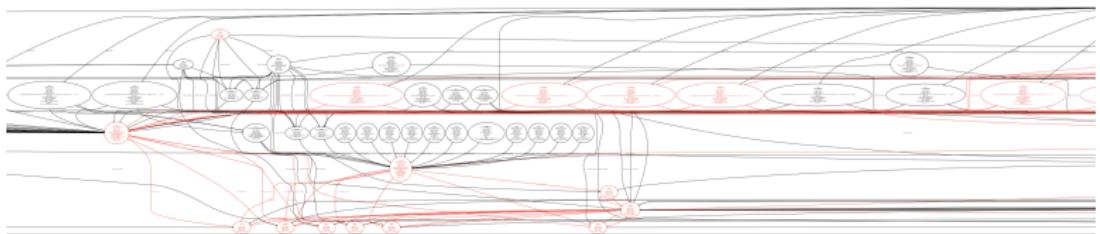
But so what?



The James
Hutton
Institute



Finally a use?



What else?

Thank you very much

