# myExperiment Research Objects: Beyond Workflows and Packs

Stian Soiland-Reyes

myGrid, University of Manchester

*BOSC 2013, ISMB, Berlin, 2013-07-18*

# Motivation: Scientific workflows

Coordinated **execution** of *services* and linked *resources*

**Dataflow** between services
- Web services (SOAP, REST)
- Command line tools
- Scripts
- User interactions
- Components (*nested workflows*)

**Method** becomes:
- **Documented** visually
- **Shareable** as single definition
- **Reusable** with new inputs
- **Repurposable** other services
- **Reproducible**?

http://www.myexperiment.org/workflows/3355

http://www.biovel.eu/
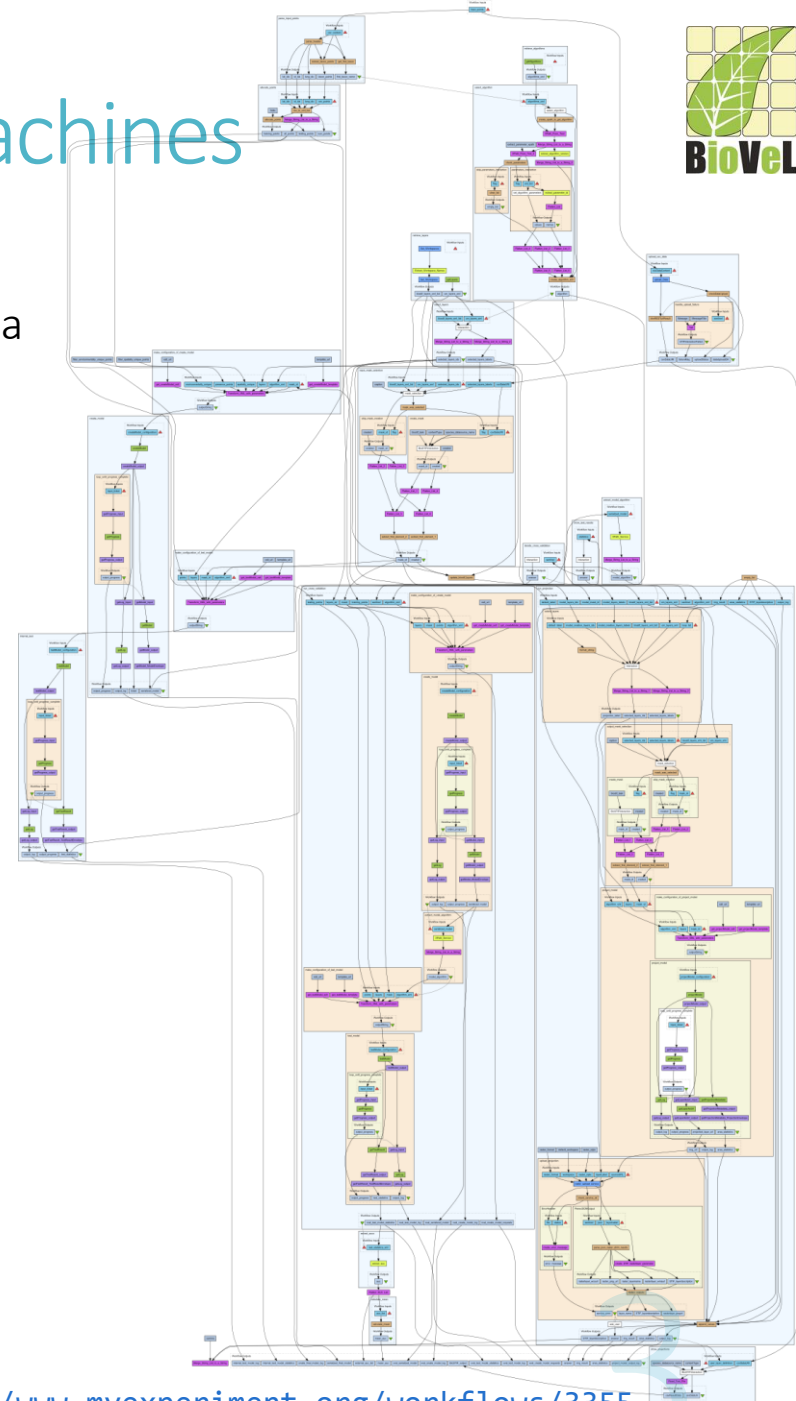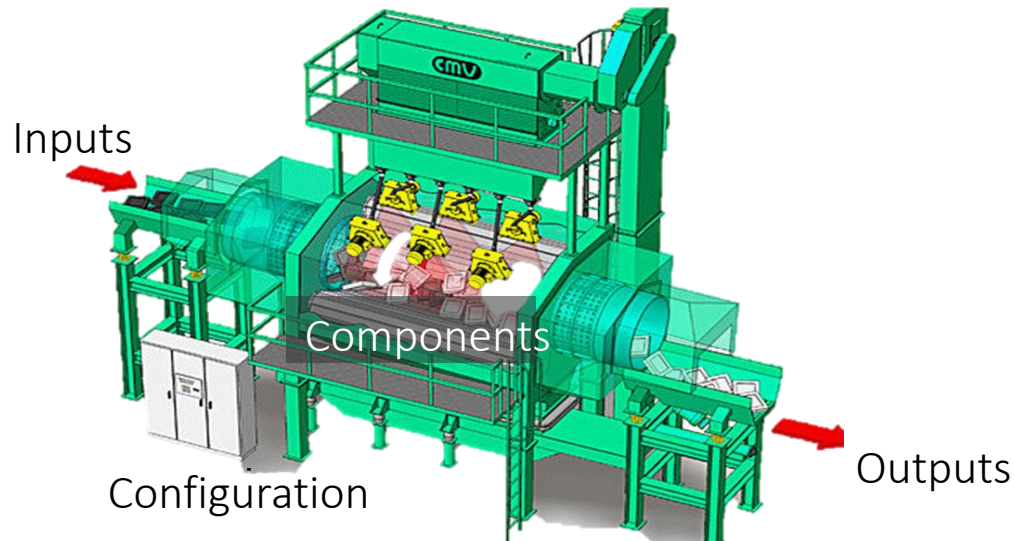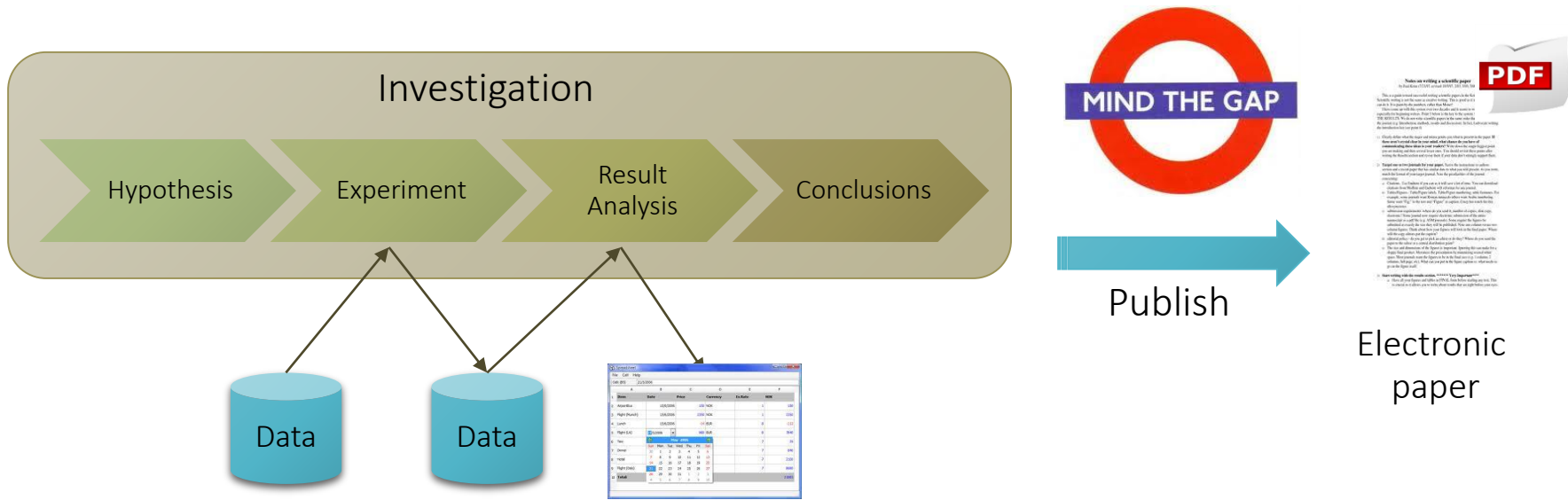
http://www.taverna.org.uk/

# But workflows are complex machines

- Will it **still work** after a year? 10 years?

- Expanding *components*, we see a workflow involves a series of specific **tools and services** which

  - **Depend** on datasets, software libraries, other tools

  - Are often poorly **described** or **understood**

  - Over time *evolve,* ***change,*** *break* or are *replaced*

- **User interactions** are not reproducible

  - But can be *tracked* and *replayed*



Inputs

Components

Configuration

Outputs

# Electronic Paper Not Enough



Investigation

Hypothesis → Experiment → Result Analysis → Conclusions

Data    Data

MIND THE GAP

Publish

Electronic paper

**Open Research** movement: Openly share the data of your experiments

figshare
credit for all your research

DRYAD

Beyond the PDF²

http://figshare.com/    http://datadryad.org/    http://www.force11.org/beyondthepdf2

4

# RESEARCH OBJECT (RO)



**Research objects** goal: Openly share *everything* about your experiments, including how those things are related

http://www.researchobject.org/

# What is in a research object?

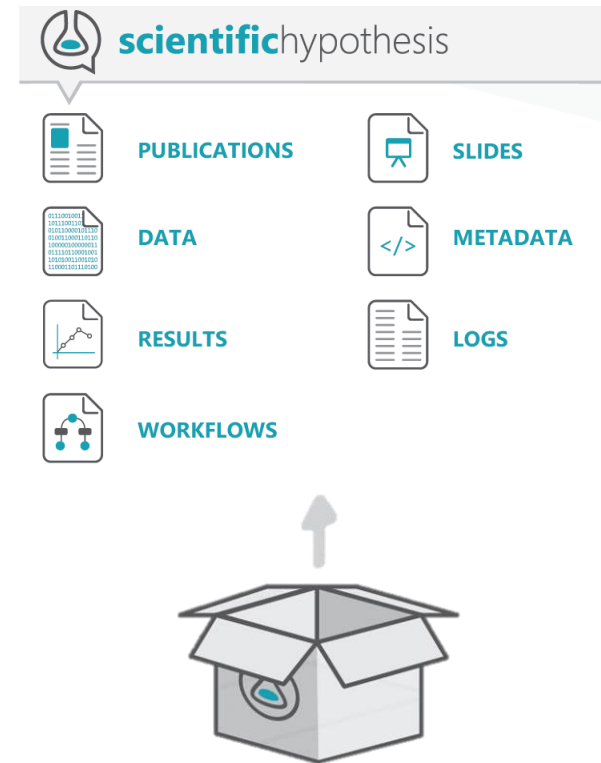A Research Object **bundles** and **relates** digital resources of a scientific experiment or investigation:

**Data** used and results produced in experimental study

**Methods** employed to produce and analyse that data

**Provenance** and settings for the experiments

**People** involved in the investigation

**Annotations** about these resources, that are essential to the understanding and interpretation of the scientific outcomes captured by a research object
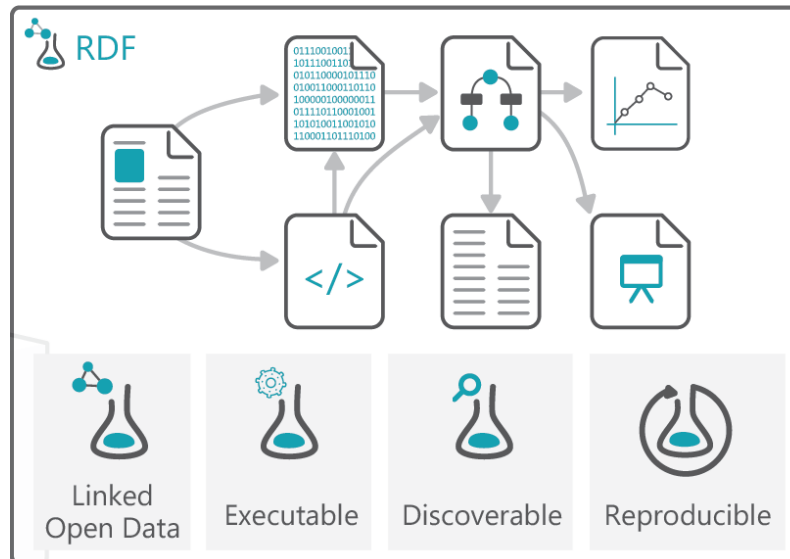
**scientific**hypothesis

PUBLICATIONS    SLIDES

DATA    METADATA

RESULTS    LOGS

WORKFLOWS

# Gathering everything

**Research Objects** (RO) *aggregate* related resources, their *provenance* and *annotations*

Conveys "everything you need to know" about a study/experiment/analysis/dataset/workflow

**Shareable**, **evolvable**, contributable, **citable**

ROs have their own provenance and lifecycles

# Why Research Objects?

i.    To **share** your research materials
   *(RO as a social object)*

ii.    To facilitate **reproducibility** and **reuse** of methods

iii.    To be **recognized** and **cited**
   *(even for constituent resources)*

iv.    To **preserve** results and **prevent decay**
   *(curation of workflow definition;
   using provenance for partial rerun)*

# A Research object

## Hypothesis

- Hypothesis.txt

## Conclusions

- conclusion.pdf

### ⓘ Items (7)

GWAStoPathway_Marco.t2flow (Workflow)

Hypothesis.txt (Hypothesis)

Mining_the_Kegg_path.wfbundle

conclusion.pdf (Conclusions)

datasetmarkers_hgvrs487.txt (Example inputs)

10.1038_ng.507

workflow_sketch_final.jpg (Sketch)

### ⓘ Relationships

datasetmarkers_hgvrs487.txt is selected as input to Mining_the_Kegg_path.wfbundle

### ⓘ Download

📁 **Download Pack Items (ZIP archive)**

**Wf4Ever tools**

- **Browse in portal**
- **Analytics and Quality**
- **Recommendations**

**Prepare for publication:**

- **Checklist**
- **Snapshot**
- **Archive**

**RO status**

live

10

# Quality Assessment of a research object

## GWAS to pathway

This pack is for a workflow that finds KEGG pathways for genes from a GWAS.

✓ Target **Pack387** *fully satisfies* checklist for *ready-to-release*.

✓ Experiment hypothesis is present
✓ Workflow design sketch is present
✓ All workflow definitions are accessible
✓ All web services used by workflows are accessible
✓ Input data is present
✓ Experiment conclusions are present

**Wf4Ever project**

11

# Quality Monitoring

# Annotations in research objects

**Types:** "This document contains an *hypothesis*"

**Relations**: "These datasets are *consumed* by that tool"

**Provenance**: "These results *came from* this workflow run"

**Descriptions**: "*Purpose* of this step is to filter out invalid data"

**Comments**: "This method looks useful, but how do I install it?"

**Examples**: "This is how you could use it"

# What is provenance?

**Attribution**
*who did it?*

**Abstraction levels**
*shallots, sign, photo or flickr page?*

**Activity**
*what happens to it?*

**Date and tool**
*when was it made?*
*using what?*

**Derivation**
*how did it change?*

**Origin**
*where is it from?*

**Aggregation**
*what is it part of?*

**Annotations**
*what do others say about it?*

**Attributes**
*what is it?*

**Licensing**
*can I use it?*

14

# Attribution

Who *collected* this sample? Who helped?

Which lab performed the *sequencing*?
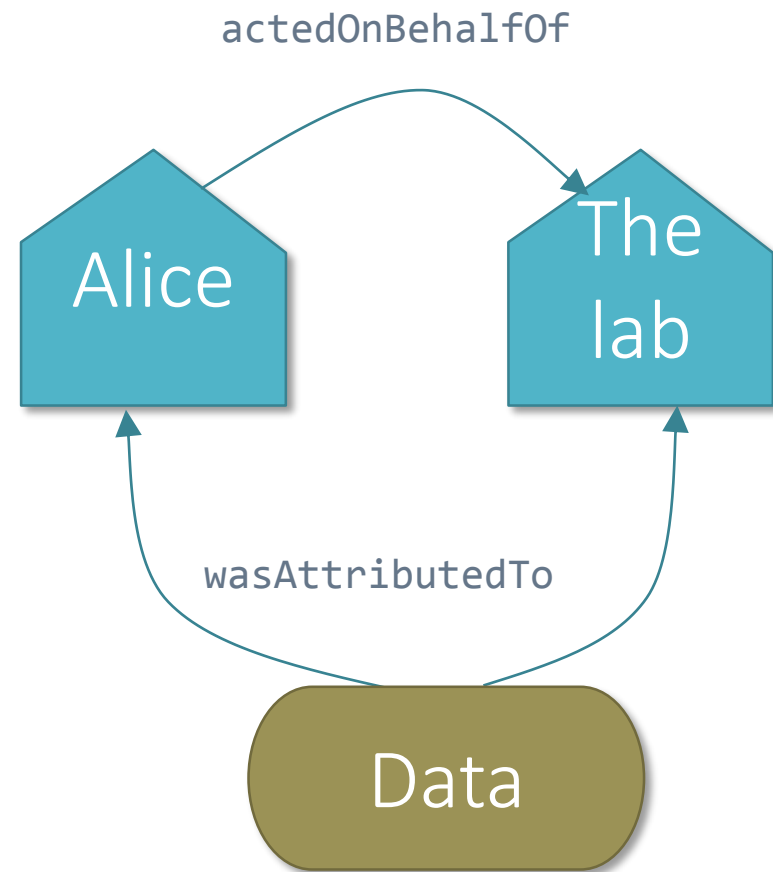
Who did the *data analysis*?

Who wrote the *analysis* workflow?

Who made the *data set* used by analysis?

Who *curated* the results?

*Why do I need this?*

i.      To be **recognized** for my work

ii.     Who should I give **credits** to?

iii.    Who should I **complain** to?

iv.     Can I **trust** them?

v.      Who should I make **friends** with?

actedOnBehalfOf

Alice

The lab

wasAttributedTo

Data

# Derivation

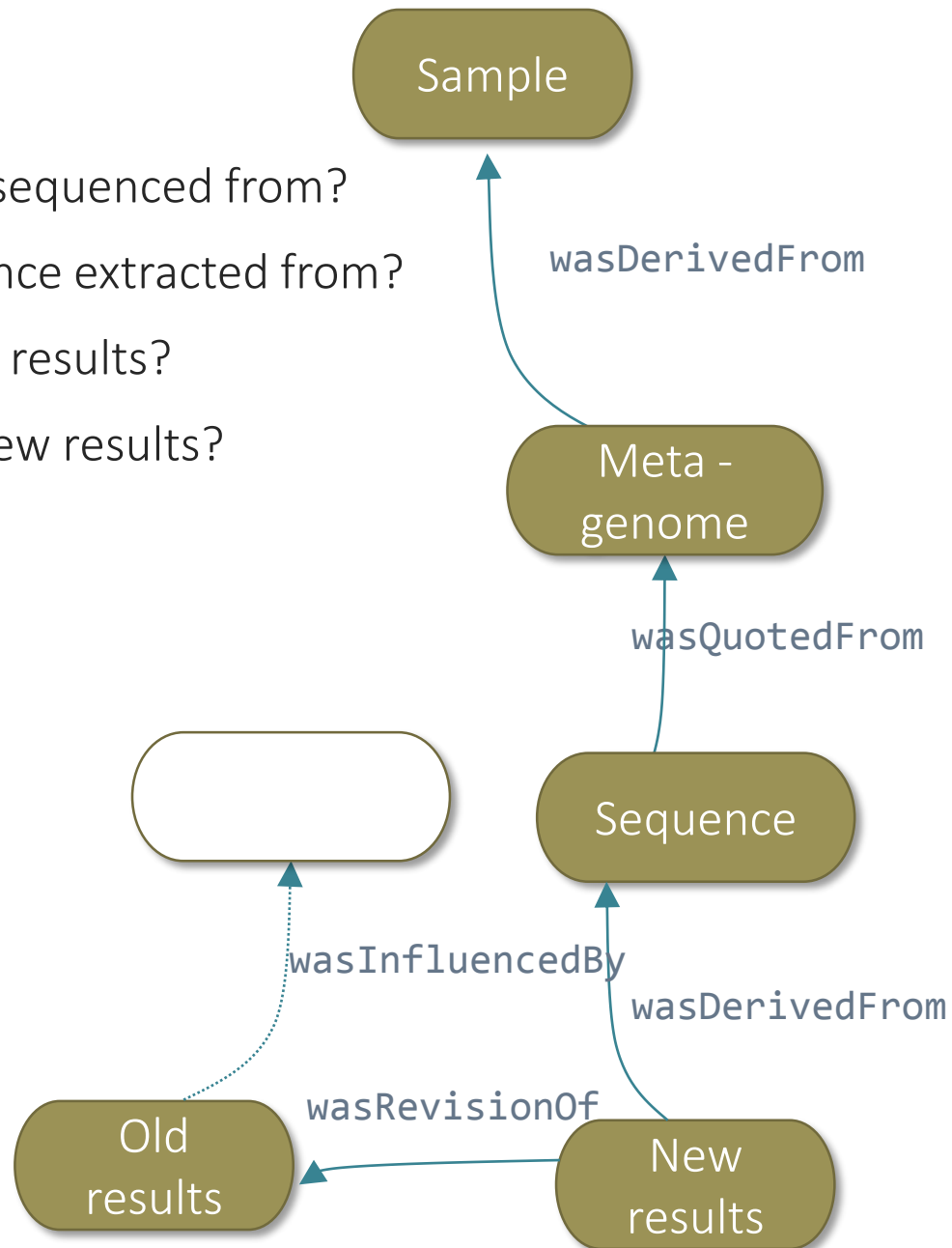Which sample was this metagenome sequenced from?

Which meta-genomes was this sequence extracted from?

Which sequence was the basis for the results?

What is the previous revision of the new results?

*Why do I need this?*

i.    To **verify** consistency (did I use the correct sequence?)

ii.   To find the latest **revision**

iii.  To **backtrack** where a diversion appeared after a change

iv.   To **credit** work I depend on

v.    **Auditing** and defence for peer review

Sample

wasDerivedFrom

Meta - genome

wasQuotedFrom

Sequence

wasInfluencedBy

wasDerivedFrom

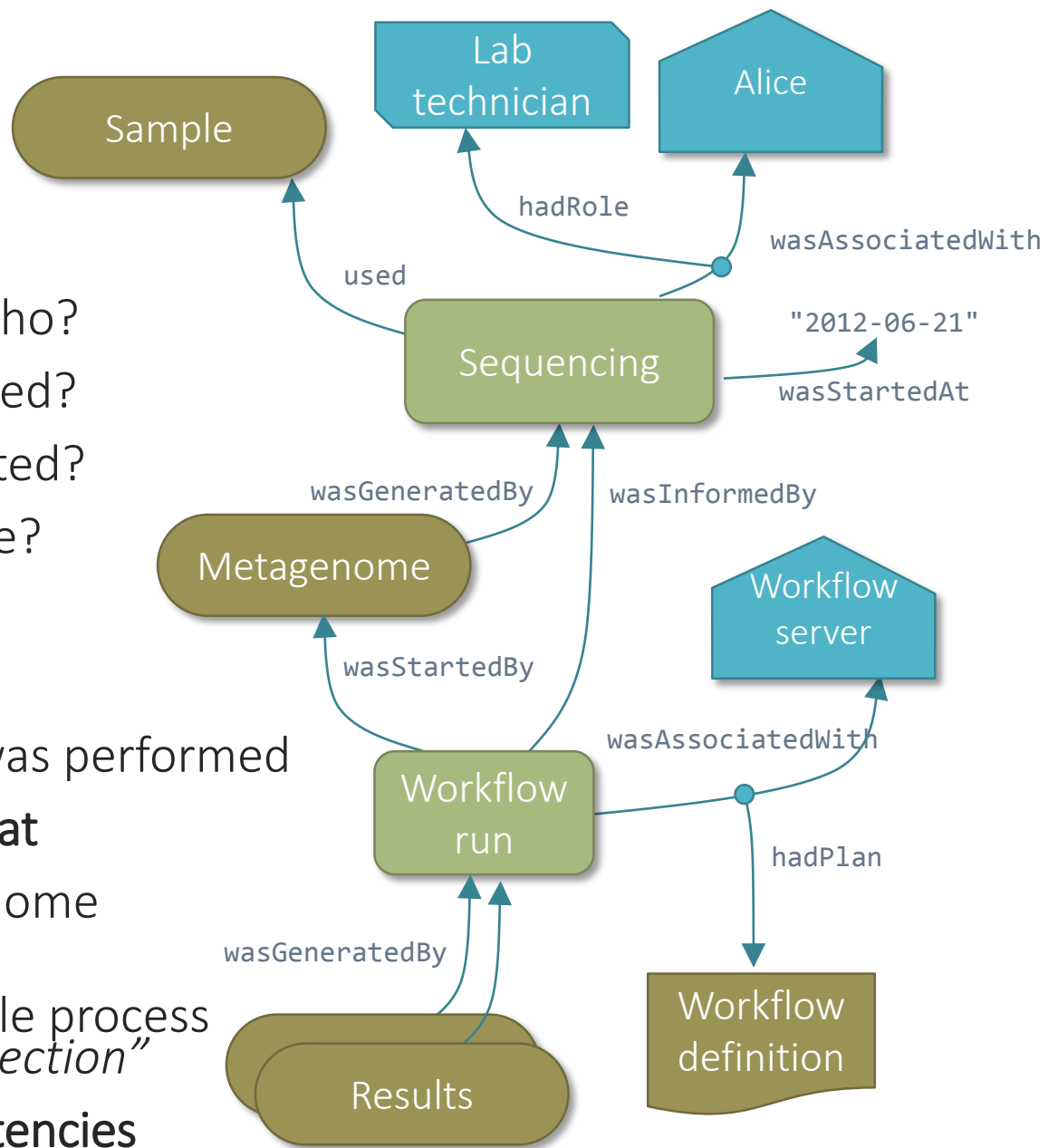Old results

wasRevisionOf

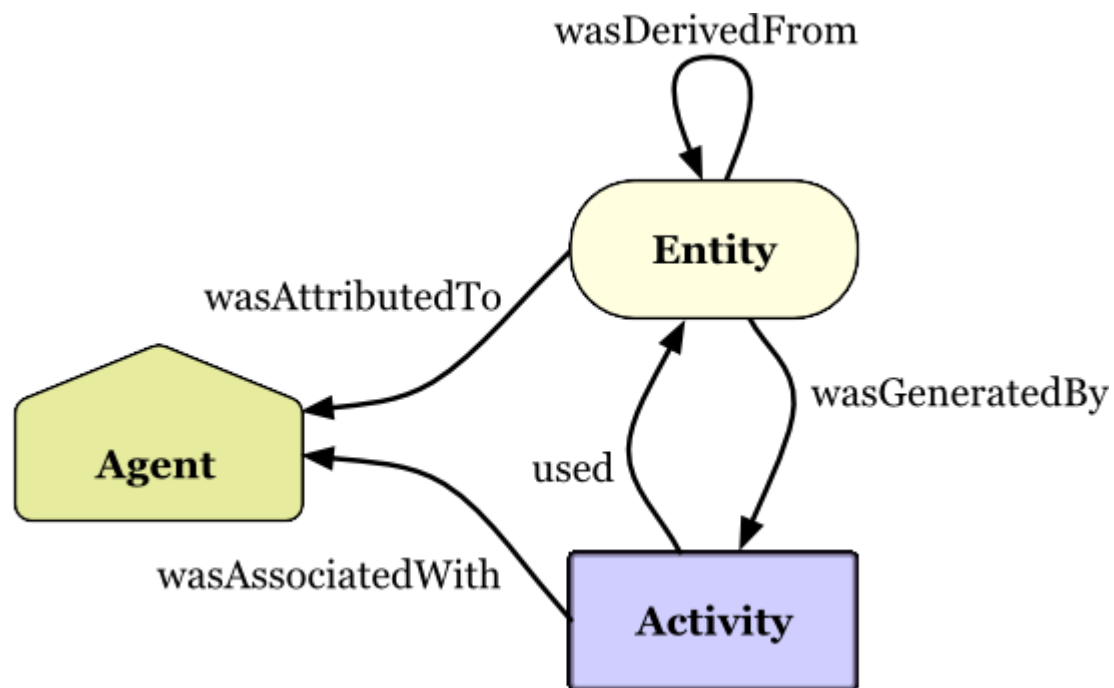New results

# Activities

What happened? When? Who?

What was used and generated?

Why was this workflow started?

Which workflow ran? Where?

*Why do I need this?*

i.   To see which **analysis** was performed

ii.  To find out **who** did **what**

iii. What was the metagenome **used** for?

iv.  To **understand** the whole process *"make me a Methods section"*

v.   To track down **inconsistencies**

Sample

Lab technician

Alice

hadRole

wasAssociatedWith

used

"2012-06-21"

Sequencing

wasStartedAt

wasGeneratedBy

wasInformedBy

Metagenome

Workflow server

wasStartedBy

wasAssociatedWith

Workflow run

hadPlan

wasGeneratedBy

Results

Workflow definition

17

# Core PROV model

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

# Provenance of what?

Who **made** the (content of) research object? Who **maintains** it?

Who **wrote** this document? Who **uploaded** it?

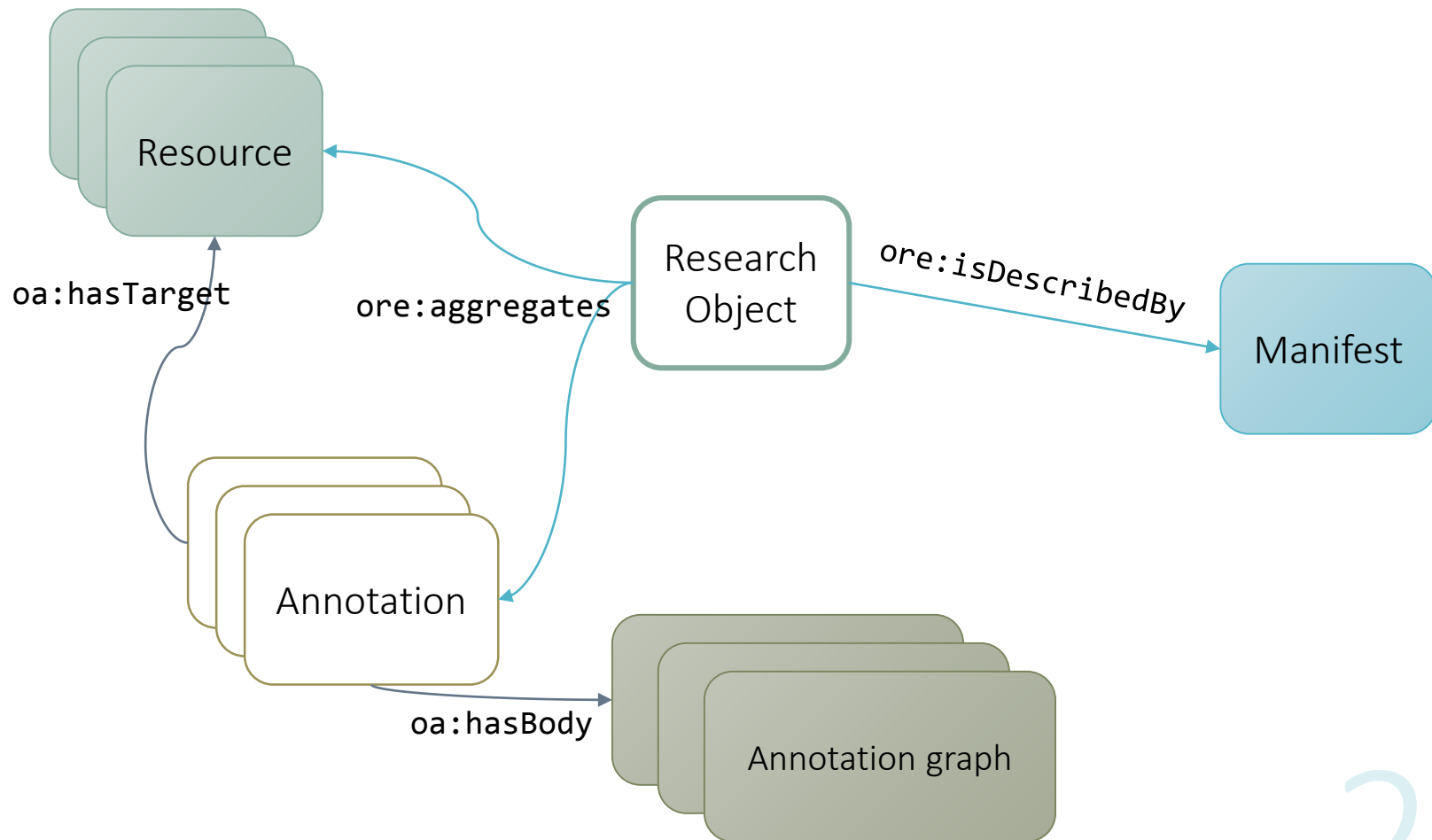Which CSV was this Excel file **imported from**?

Who wrote this **description**? When? How did we get it?

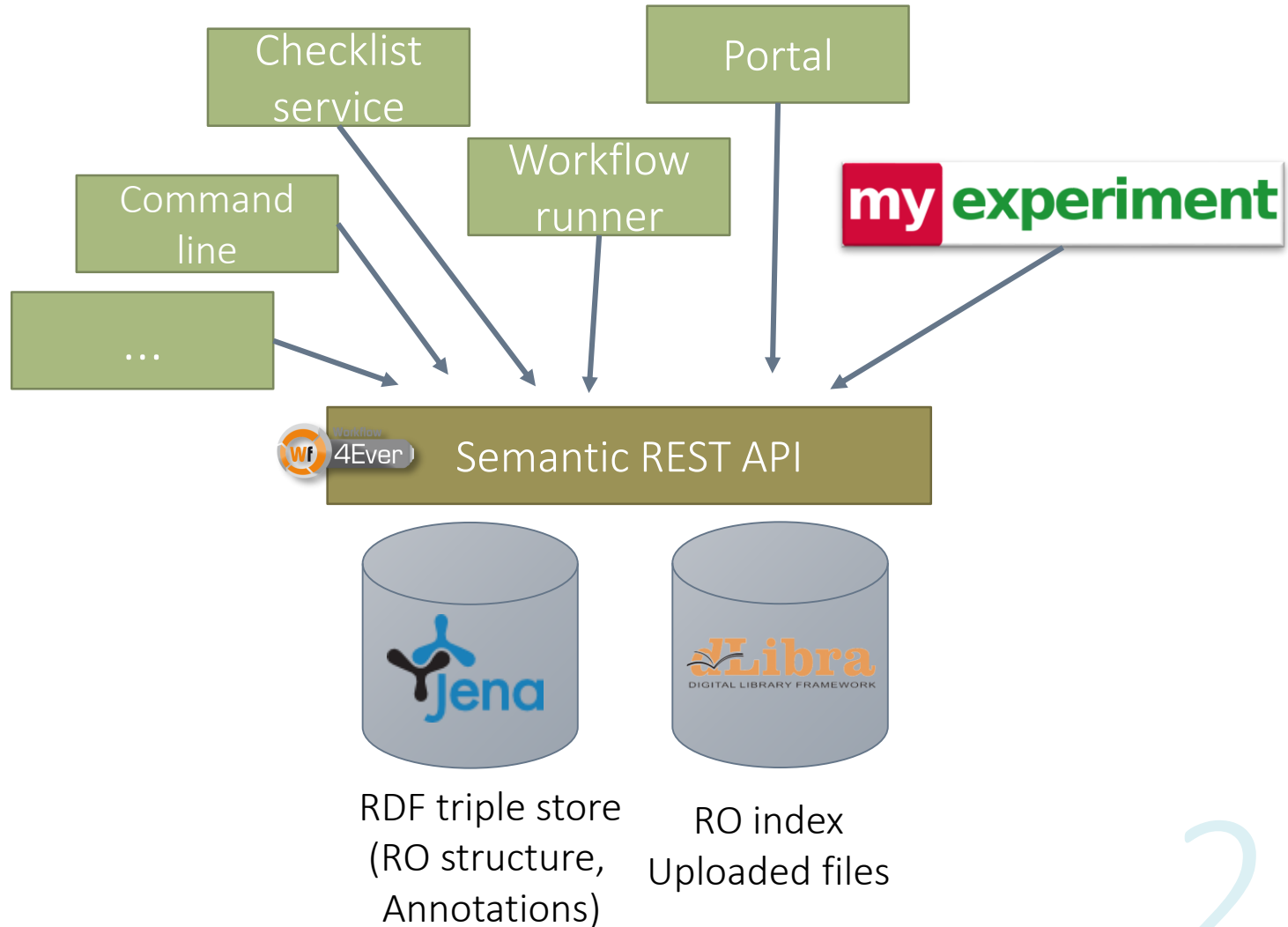What is the **state** of this RO? (Live or Published?)

What did the research object look like before? (**Revisions**) – are there newer versions?

Which research objects are **derived** from this RO?

19

# Research object model at a glance

# Wf4Ever architecture

Checklist service

Portal

Command line

Workflow runner

**my** experiment

...

Wf 4Ever Semantic REST API

Jena

dLibra
DIGITAL LIBRARY FRAMEWORK

RDF triple store
(RO structure,
Annotations)

RO index
Uploaded files

21

# Saving a research object: RO bundle

Single, **transferrable** research object

    Self-contained **snapshot**

    Which files in ZIP, which are URIs? (Up to user/application)

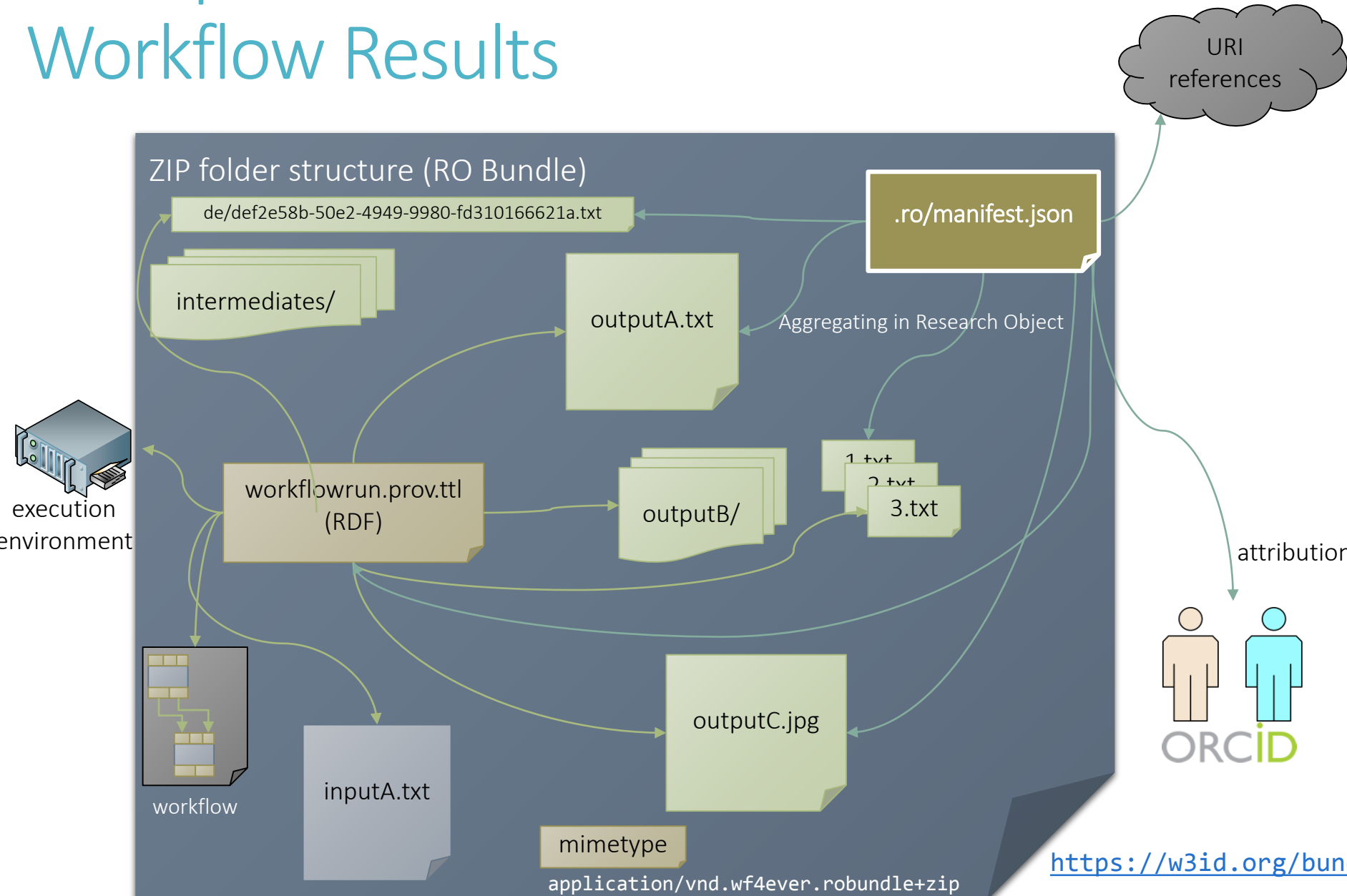Regular **ZIP file**, explored and unpacked with standard tools

**JSON manifest** is programmatically accessible without RDF understanding

Works **offline** and in desktop applications – no REST API access required

Basis for RO-enabled **file formats**, e.g. Taverna run bundle

**Exchanged** with myExperiment and RO tools

# Example RO bundle: Workflow Results



ZIP folder structure (RO Bundle)

de/def2e58b-50e2-4949-9980-fd310166621a.txt

.ro/manifest.json

intermediates/

outputA.txt

Aggregating in Research Object

execution environment

workflowrun.prov.ttl (RDF)

outputB/

1.txt
2.txt
3.txt

attribution

workflow

inputA.txt

outputC.jpg

ORCID

mimetype

application/vnd.wf4ever.robundle+zip

URI references

https://w3id.org/bund

# W3C community group for RO

# Summary

**Provide provenance records**: Use (and extend) PROV.

**Share your methods** (tools, workflows), not just your data

**Make research objects**: Collect resources and relate them

**Pick your RO integration level** to adapt**:**
1. *Conceptual model*
2. *RO Core ontology (aggregation and annotation)*
3. *Workflow description and provenance*
4. *Wf4Ever architecture for APIs*
5. *myExperiment as user interface*

26