# Functional Requirements for Information Resource Provenance on the Web

James P. McCusker, **Timothy Lebo**, Alvaro Graves, Dominic Difranzo, Paulo Pinheiro, and Deborah L. McGuinness

<Timothy Lebo>   prov:actedOnBehalfOf    <James P. McCusker> .

# What do we use a URL to identify?

(Answer: too many things)

(Problem: we're getting confused!)

# Two provenance use cases

*How do we use URLs to track the provenance of the information we get from them?*

**Image analysis:** who has reviewed this image?

- Different tools identify and use different formats and resolutions.
- How do we know when Dr. Smith reviewed the image we're seeing?

**Weather:** weather.gov offers ongoing forecasts at an unchanging URL.

# What is a URL?

Is it the data?

GET /web/tw-logo
Accept: image/jpg

GET /web/tw-logo
Accept: image/png

TWC

Use case 1: Image Analysis

# What is a URL?

~~Is it the data?~~

Is it the content?



GET /web/tw-logo
Accept: image/jpg
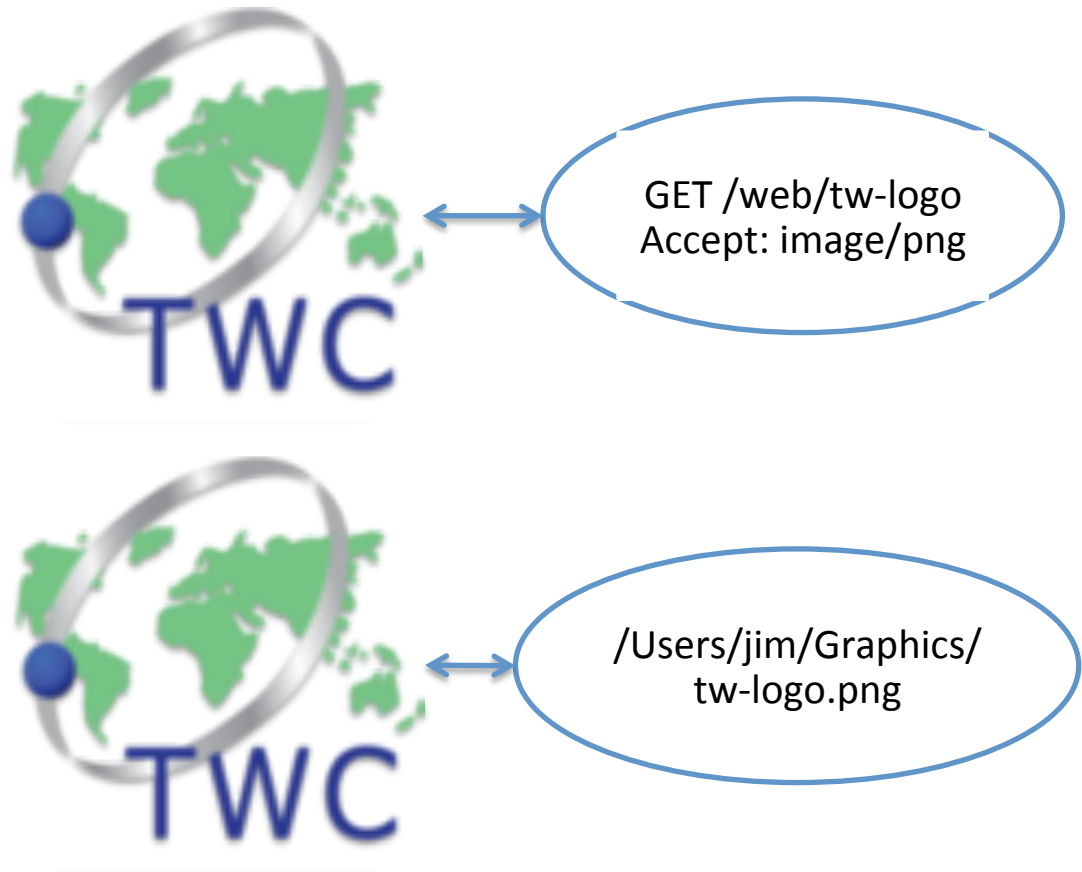Date: Today

GET /web/tw-logo
Accept: image/png
Date: Yesterday

Use case 1: Image Analysis

## What is a URL?

~~Is it the data?~~

~~Is it the content?~~

What about multiple
copies of the same file?



GET /web/tw-logo
Accept: image/png



/Users/jim/Graphics/
tw-logo.png

Use case 2: Weather

# A Little More on Weather

It changes over time

It has multiple representations

Copies might be saved for historical analysis

But it's all the same URL!



http://www.weather.gov/xml/current_obs/KBOS.xml

We need to solve these confusions
to get good provenance for
information on the web.

Okay, fine. Content changes.

Let's be extremely conservative:
unique identifiers to identify the content?

There are 4 steps to data sanity.

# 4 Steps to Data Sanity

**Step 1:** Distinguish resources from their representations

The Architecture of the World Wide Web [1] provides a framework for this.

URI

`http://weather.example.com/oaxaca`

*Identifies*

Resource

*Oaxaca Weather Report*

*Represents*

Representation

**Metadata:**
Content-type:
application/xhtml+xml

—————————————

**Data:**
```
<!DOCTYPE html PUBLIC "...
    "http://www.w3.org/...
<html xmlns="http://www...
<head>
<title>5 Day Forecaste for
Oaxaca</title>
...
</html>
```

[1] http://www.w3.org/TR/webarch/

# 4 Steps to Data Sanity

**Step 1:** Distinguish resources from their representations

Actually, this is very close to the semiotic triangle.

We can now separate out symbols, referents, and representations.

This leads to fragmentation as identifiers proliferate.

Thought

symbolizes

refers to

stands for

Symbol (Sem) ────────→ Referent (Sem)

correspondence

URL (WWW) ────────→ Resource (WWW)

identifies

represents

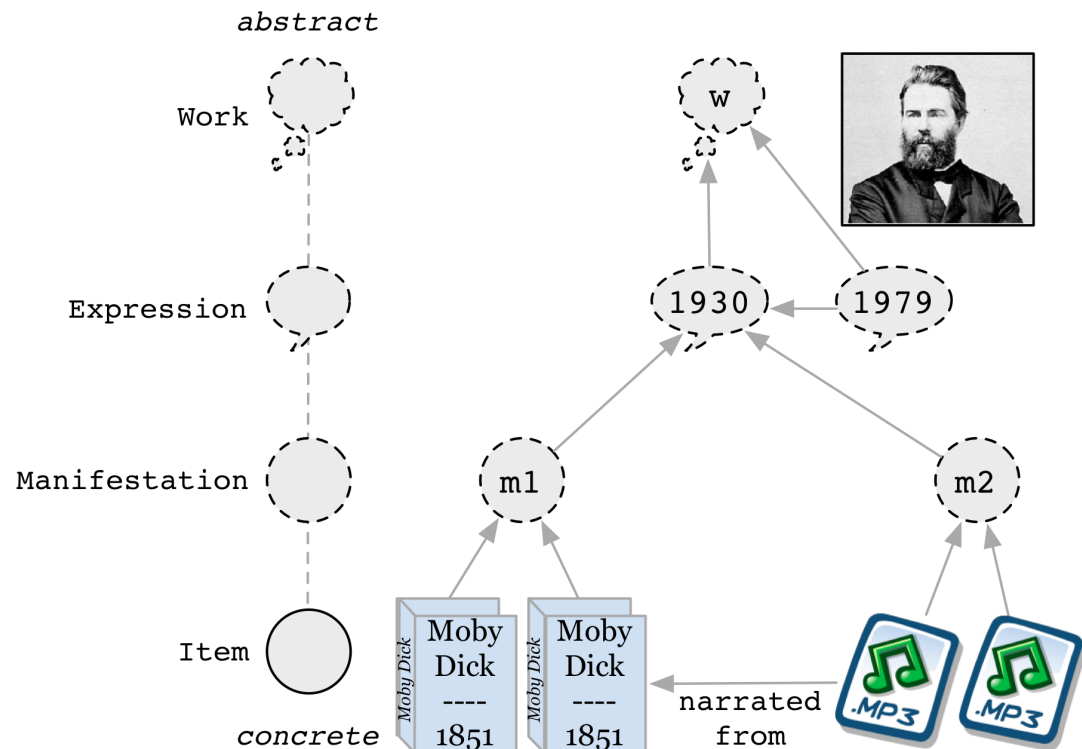Representation

# Wait, what is FRBR?

**Functional Requirements for Bibliographic Records**

**Work:** a "distinct intellectual or artistic creation."[1]

**Expression:** "the specific intellectual or artistic form that a work takes each time it is 'realized.'"[1]

**Manifestation:** "the physical embodiment of an expression of a work. As an entity, manifestation represents all the physical objects that bear the same characteristics, in respect to both intellectual content and physical form."[1]

**Item:** "a single exemplar of a manifestation. The entity defined as item is a concrete entity."[1]



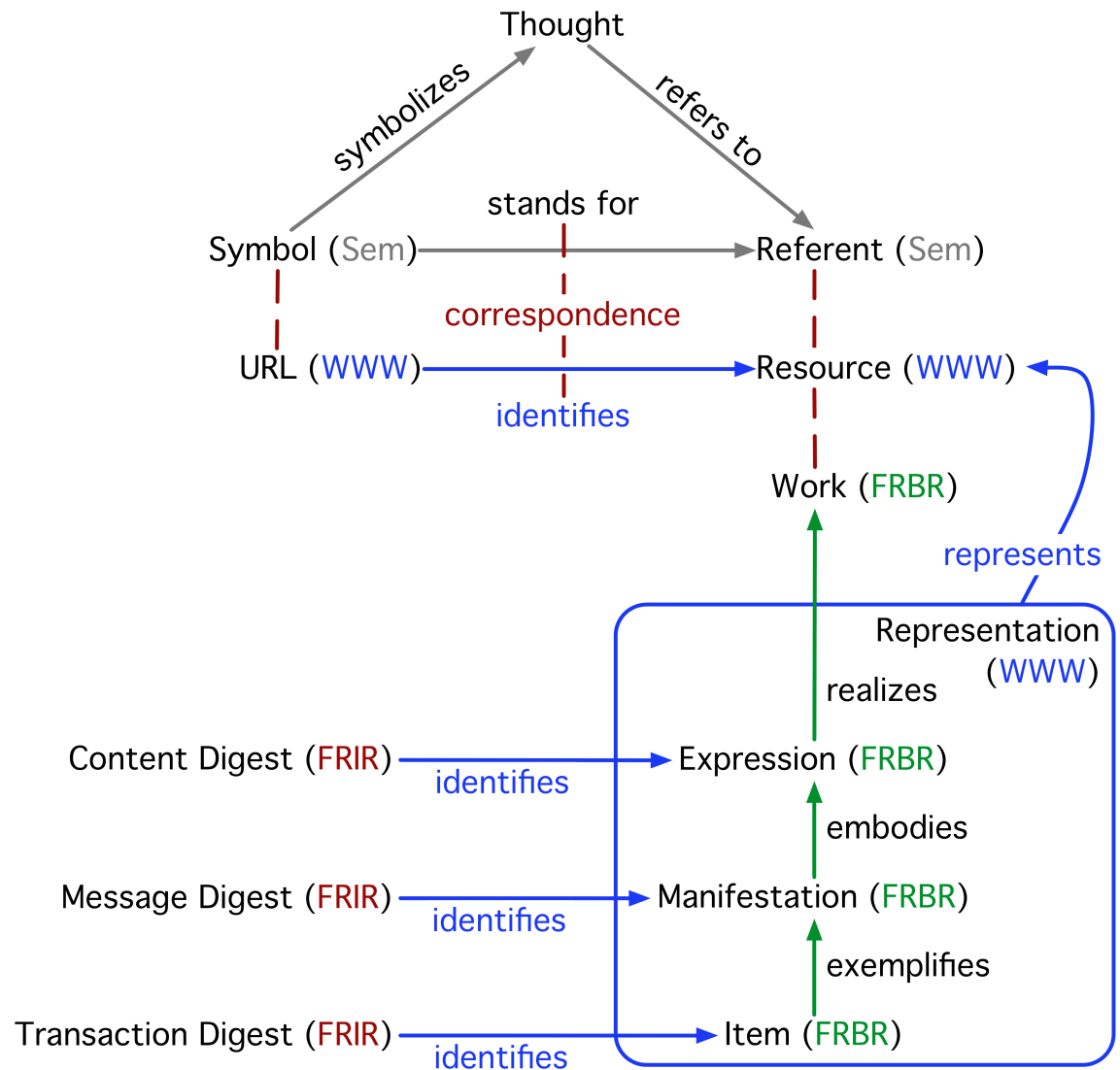1. Functional Requirements for Bibliographic Records: http://archive.ifla.org/VII/s13/frbr/frbr1.htm#3.2

# 4 Steps to Data Sanity

**Step 1:** Distinguish resources from their representations.

**Step 2:** Avoid fragmentation by using FRBR and message/content digests.

# 4 Steps to Data Sanity

**Step 1:** Distinguish resources from their representations.

**Step 2:** Avoid fragmentation by using FRBR and message/content digests.
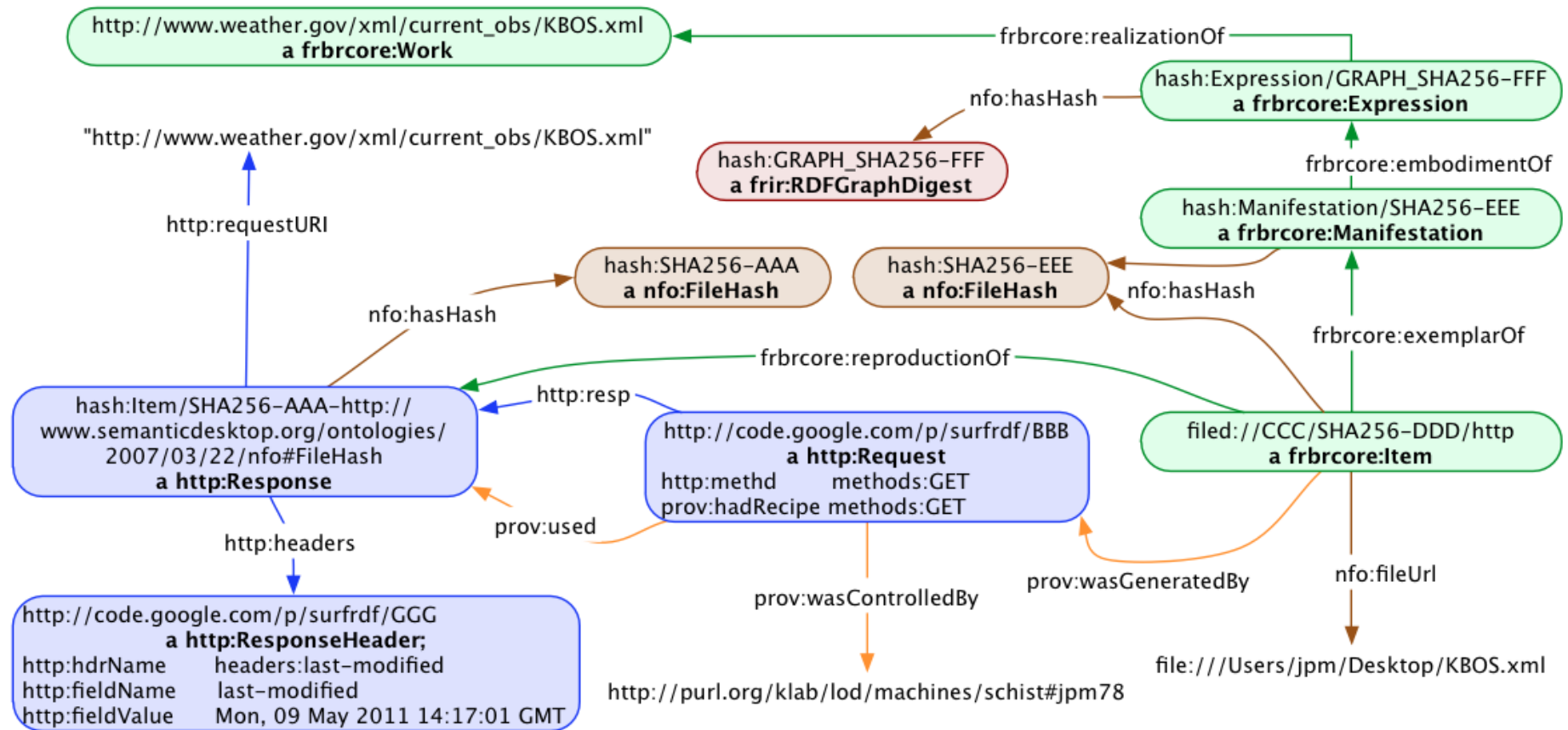
**Step 3:** Applying PROV

This is **Functional Requirements for Information Resources (FRIR).**

We found that 14 of 18 frbr:relatedEndeavour subproperties mapped to one or more PROV properties.

| Subclass | Superclass |
|---|---|
| frbr:Event | prov:Activity |
| frbr:ResponsibleEntity | prov:Agent |
| frbr:Endeavour | prov:Entity |
| nie:DataObject | prov:Entity |

| Subproperty | wasDerivedFrom | alternateOf | specializationOf |
|---|---|---|---|
| frbr:adaptionOf | X | | |
| frbr:imitationOf | X | | |
| frbr:reconfigurationOf | X | | |
| frbr:transformationOf | X | | |
| frbr:abridgementOf | X | X | |
| frbr:arrangementOf | X | X | |
| frbr:reproductionOf | X | X | |
| frbr:summarizationOf | X | X | |
| frbr:translationOf | X | X | |
| frbr:alternateOf | | X | |
| frbr:revisionOf | | X | |
| frir:redirectsToTransitive | | X | |
| frbr:embodimentOf | | | X |
| frbr:exemplarOf | | | X |
| frbr:realizationOf | | | X |

# Explaining a HTTP transaction for the Weather.

**Green**: Information aspects of information accessed from the URL

**Blue**: The HTTP request and response

**Brown**: Content and message digest hash entities that identify the aspects.

## 4 Steps to Data Sanity

**Step 1:** Distinguish resources from their representations.

**Step 2:** Avoid fragmentation by using FRBR and message/content digests.

**Step 3:** Use PROV

**Step 4:** Use our tools.

**pcurl.py:** provenance-enabled curl.

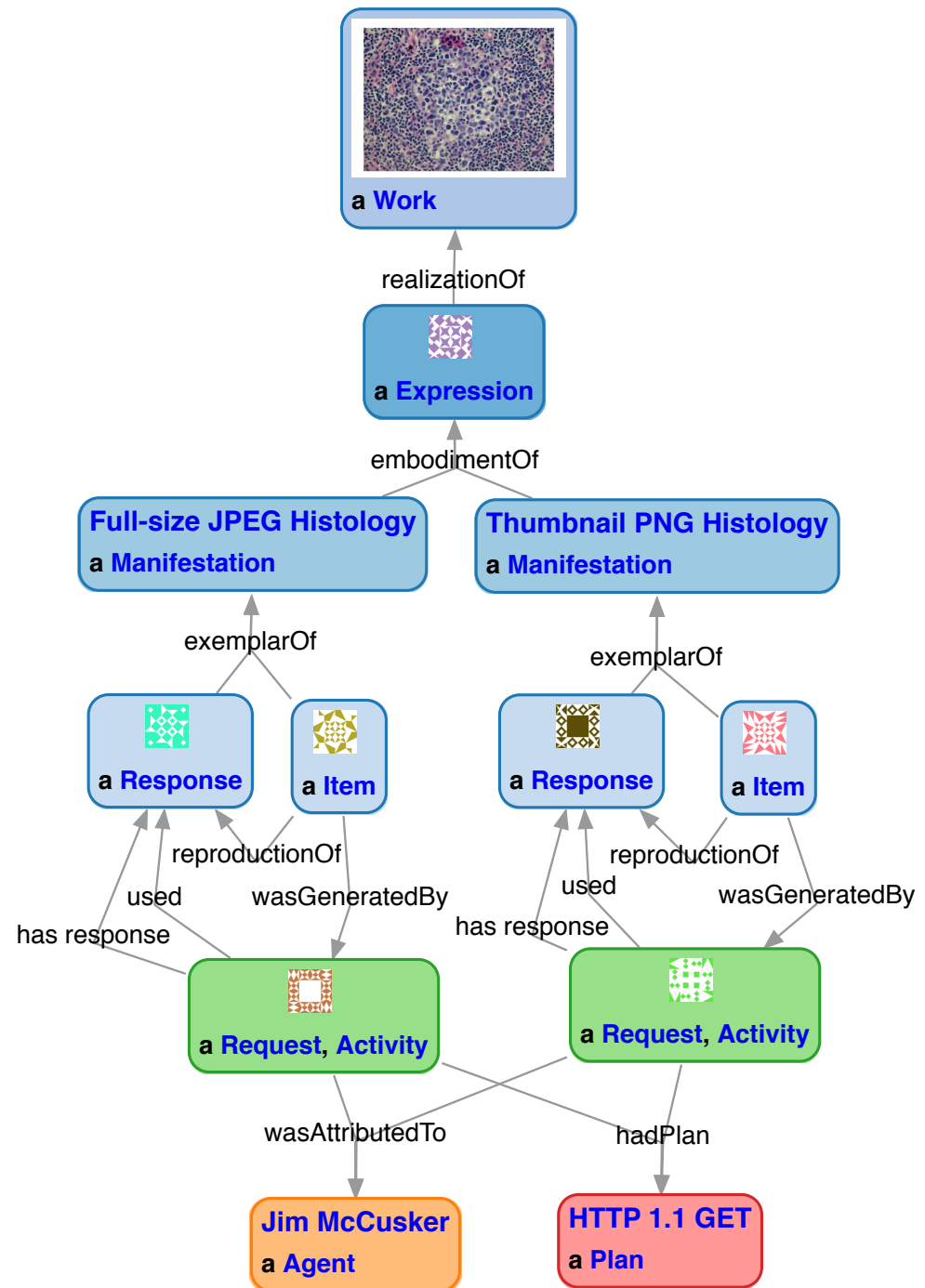**fstack.py:** build a FRBR stack for any file.

pcurl.py and fstack.py:

- Use cryptographic identities based on message digests and content digests.

- Have Java and Python implementations.

- Are open source.

  https://github.com/timrdf/csv2rdf4lod-automation/wiki

# Example: Files as different formats

Two files with the same image at different resolutions in different formats still have the same Expression and Work.

We can now link the high-resolution file used by the pathologist to the low resolution summary image the patient sees.

# It also works for HTTP POST

(See the paper for details)

## 4 Steps to Data Sanity

**Step 1:** Distinguish resources from their representations.

**Step 2:** Avoid fragmentation by using FRBR and message/content digests.

**Step 3:** Use PROV

**Step 4:** Use our tools.

# Conclusions

- Know what is being identified.

- Separate out content, data, and work.

- Relate them to each other using FRBR/FRIR

- Walk that abstraction hierarchy to talk about things in the detail you need.

Avoid confusion and Avoid fragmentation!

# Questions?

(Thank you)