



Leibniz-Institut für
Astrophysik Potsdam



Provenance Data Model

Let's keep discussions going on ...

InterOp Sesto, June 2015

Kristin Riebe, GAVO

Status

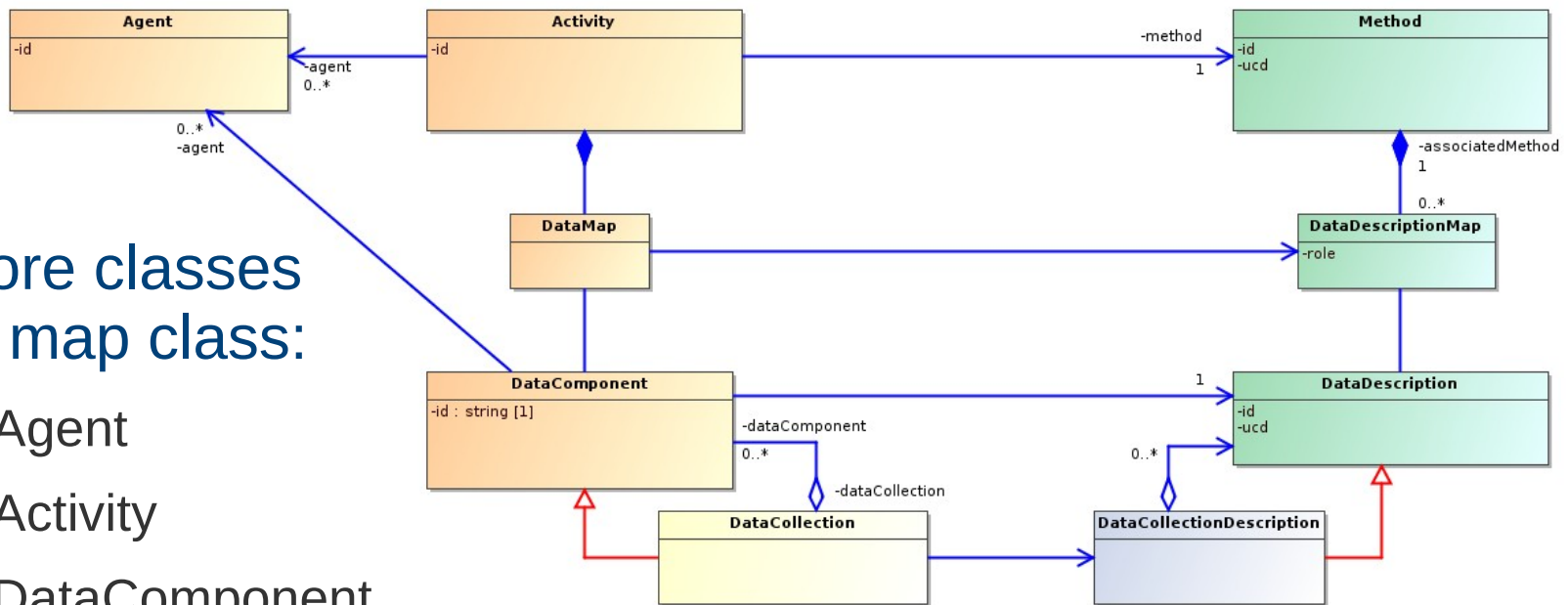
First attempt: Model with prototypes

(inspired by SimDM and W3C model)

- 3 core classes
+ 1 map class:

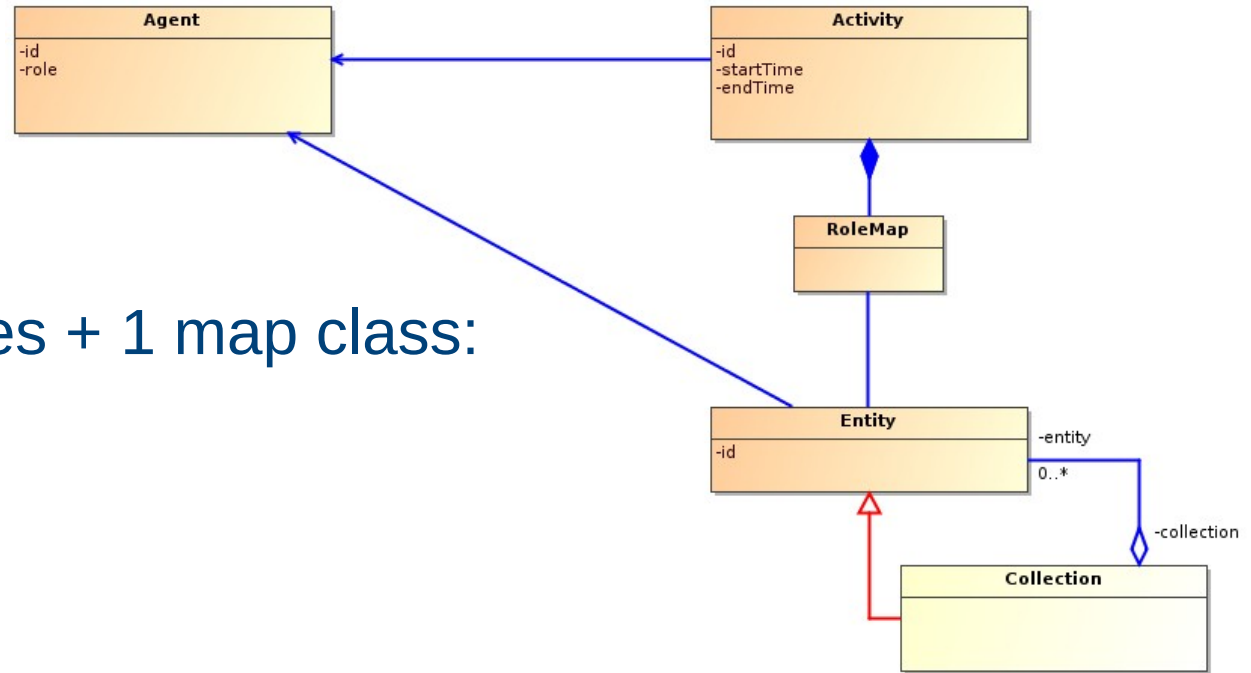
- Agent
- Activity
- DataComponent
- DataMap

- + dataDescription side => each item doubled



Status

Model without protoypes



- 3 core classes + 1 map class:

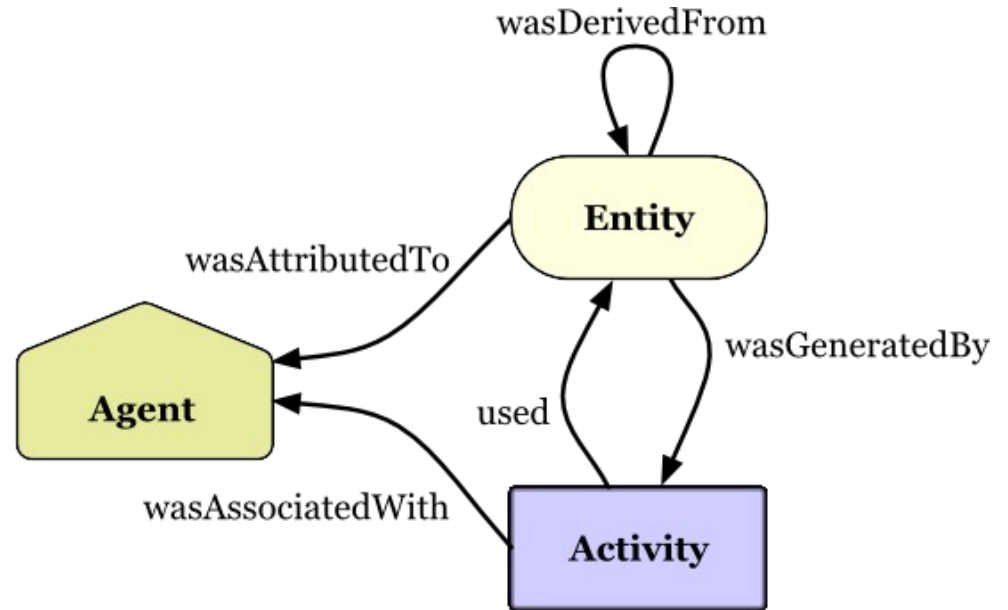
- Agent
- Activity
- Entity
- RoleMap

- Let's start with this – put descriptions into attributes, common vocabularies and see how far we get

Looking at W3C

<http://www.w3.org/TR/prov-dm/>

- 3 core classes:
 - Agent
 - Activity
 - Entity
- core relations:
 - used
 - wasGeneratedBy
 - wasDerivedFrom
 - wasAttributedTo
 - wasAssociatedWith
- + many more classes and relations



Comparison with W3C

<http://www.w3.org/TR/prov-dm/>

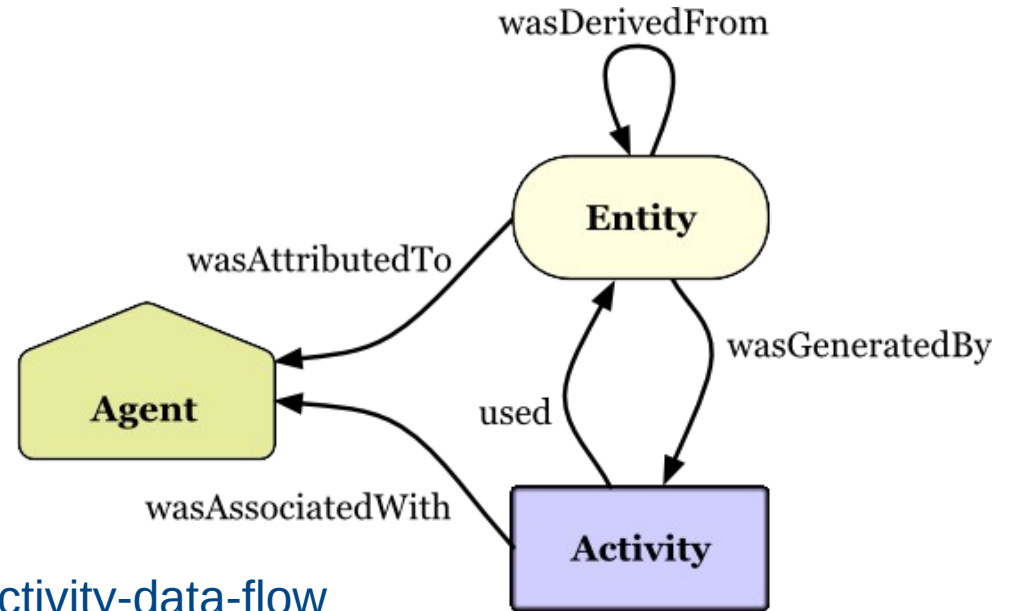
- 3 core classes:

- Agent
- Activity
- Entity

- core relations:

- used
 - wasGeneratedBy
 - wasDerivedFrom
 - wasAttributedTo
 - wasAssociatedWith
- } activity-data-flow
- data-flow
- } responsibility view

- + many more classes and relations



Comparison with W3C

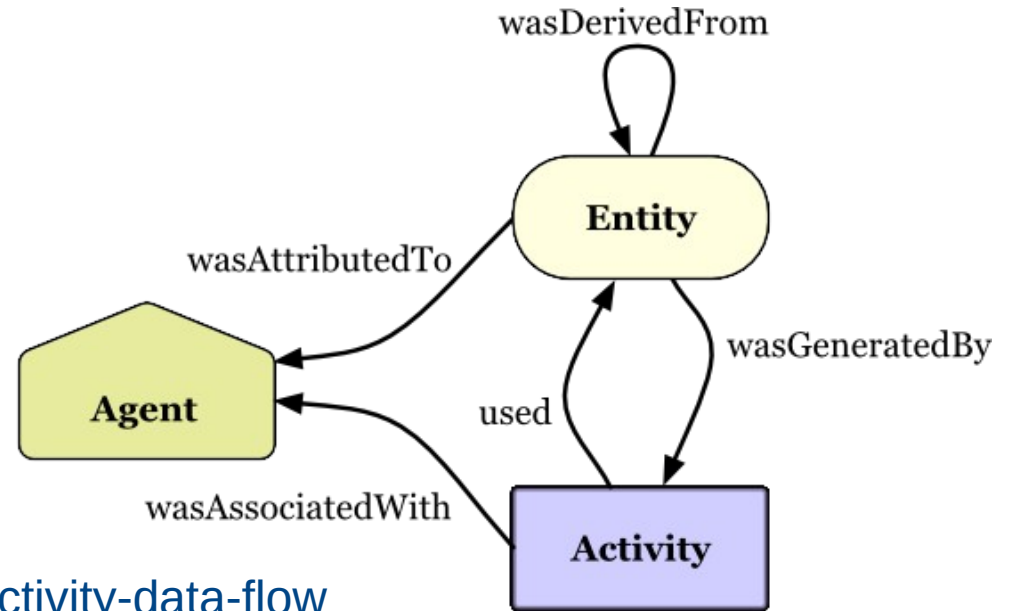
<http://www.w3.org/TR/prov-dm/>

- 3 core classes:

- Agent
- Activity
- Entity

- core relations:

- used
 - wasGeneratedBy
 - wasDerivedFrom
 - wasAttributedTo
 - wasAssociatedWith
- } activity-data-flow
- data-flow
- } responsibility view



=> Should be sufficient to give us what we need!

- + many more classes and relations

Example: Reduced RAVE-fits file

- PROV-N notation

- 2 files

```
entity(rave:0645m522I0049.wav.fits, [prov:type = 'std:fits']  
entity(rave:0645m522I0049.fits, [prov:type = 'std:fits']
```

- 2 agents

- 2 activities

- relations

Example: Reduced RAVE-fits file

- PROV-N notation

- 2 files

```
entity(rave:0645m522I0049.fits, [prov:type = 'std:fits']  
entity(rave:0645m522I0049.wav.fits, [prov:type = 'std:fits']
```

- 2 agents

```
agent(aao:Paul_Cass, [prov:type='prov:Person'])  
agent(rave:Alessandro_Siviero, [prov:type='prov:Person'])
```

- 2 activities

- relations

Example: Reduced RAVE-fits file

- PROV-N notation

- 2 files

```
entity(rave:0645m522I0049.fits, [prov:type = 'std:fits']  
entity(rave:0645m522I0049.wav.fits, [prov:type = 'std:fits']
```

- 2 agents

```
agent(aao:Paul_Cass, [prov:type='prov:Person'])  
agent(rave:Alessandro_Siviero, [prov:type='prov:Person'])
```

- 2 activities

```
activity(rave:act_observation, 2008-02-16T13:25:24, -,  
        [ prov:type = 'obs:Observation' ] )  
activity(rave:act_irafReduction, 2008-03-04T09:46:57, -,  
        [ prov:type = 'std:reduction' ] )
```

- relations

Example: Reduced RAVE-fits file

- PROV-N notation

- 2 files

```
entity(rave:0645m522I0049.fits, [prov:type = 'std:fits']  
entity(rave:0645m522I0049.wav.fits, [prov:type = 'std:fits']
```

- 2 agents

```
agent(aao:Paul_Cass, [prov:type='prov:Person']  
agent(rave:Alessandro_Siviero, [prov:type='prov:Person'])
```

- 2 activities

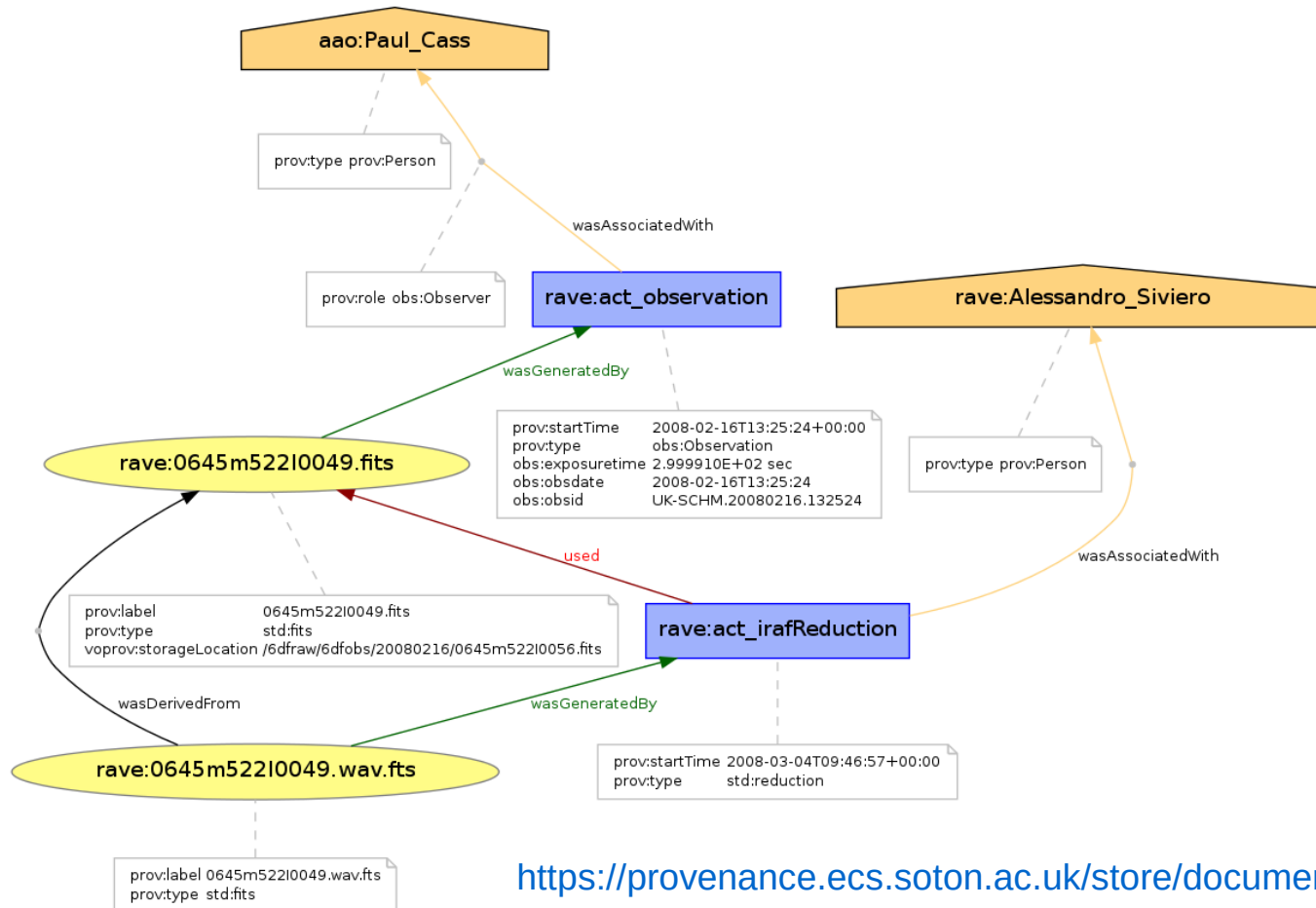
```
activity(rave:act_observation, 2008-02-16T13:25:24, -,  
[ prov:type = 'obs:Observation' ] )  
activity(rave:act_irafReduction, 2008-03-04T09:46:57, -,  
[ prov:type = 'std:reduction' ] )
```

- relations

```
wasAssociatedWith(rave:act_observation, aao:Paul_Cass, -,  
[ prov:role = 'obs:Observer' ] )  
wasAssociatedWith(rave:act_irafReduction, rave:Alessandro_Siviero, -)  
wasGeneratedBy(rave:0645m522I0049.fits, rave:act_observation, -)  
used(rave:act_irafReduction, rave:0645m522I0049.fits, -)  
wasGeneratedBy(rave:0645m522I0049.wav.fits, rave:act_irafReduction, -)  
wasDerivedFrom(rave:0645m522I0049.wav.fits, rave:0645m522I0049.fits)
```

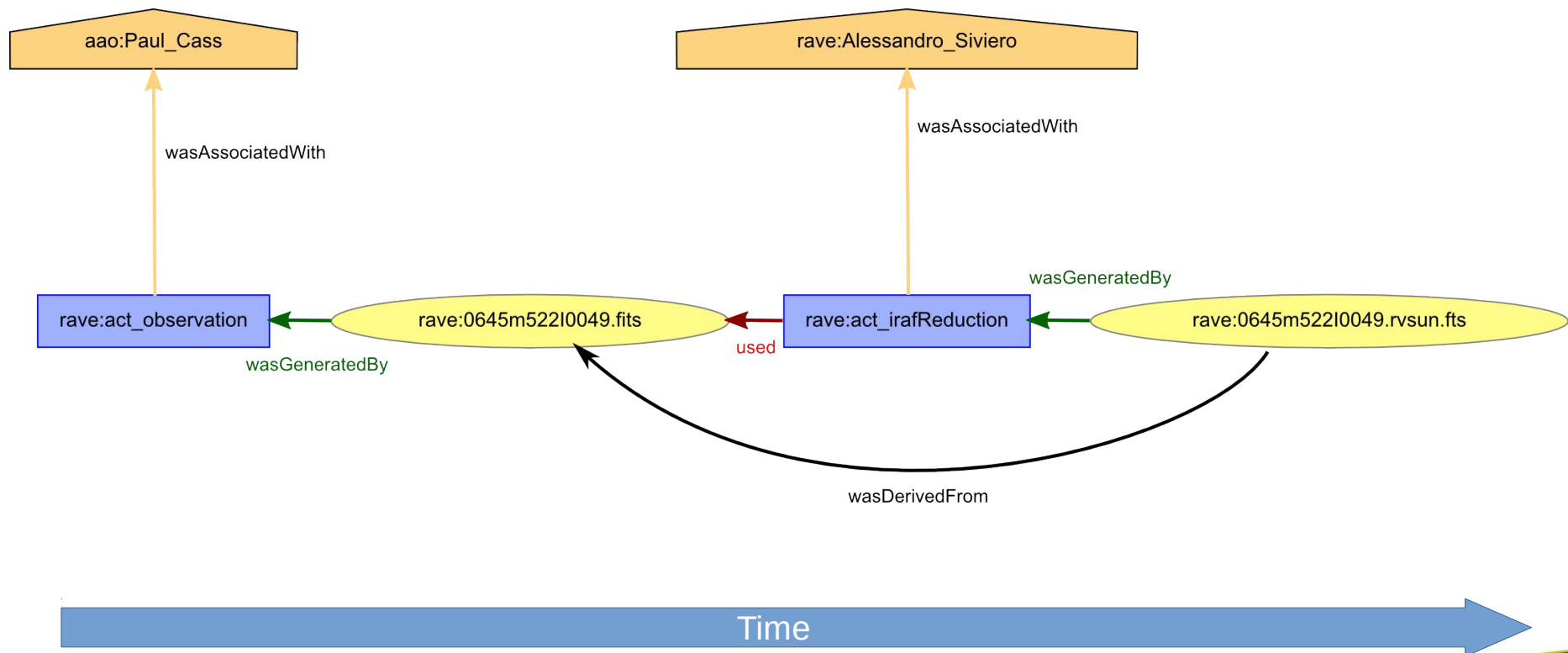
Example: Reduced RAVE-fits file

- Graph produced with ProvStore (using GraphViz):



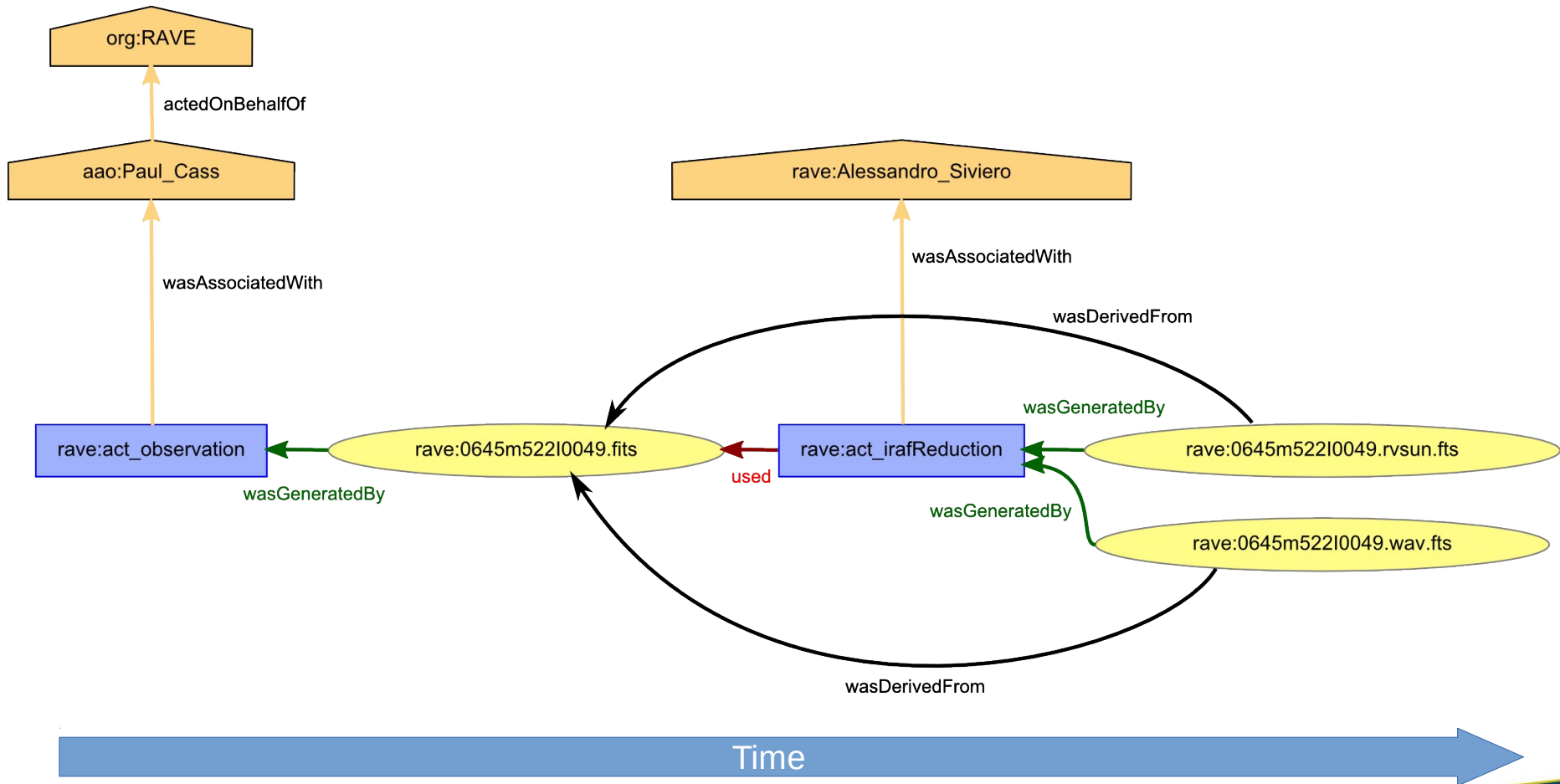
Example: Reduced RAVE fits-file

- Graph reordered, attributes hidden:

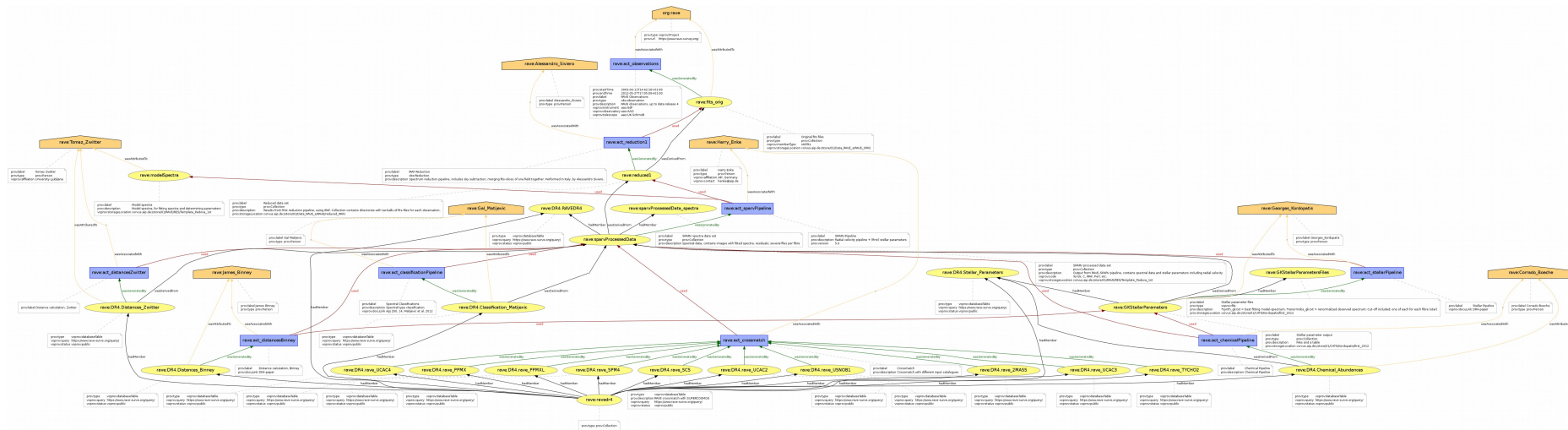


Example: Reduced RAVE fits-file

- Graph reordered, attributes hidden:

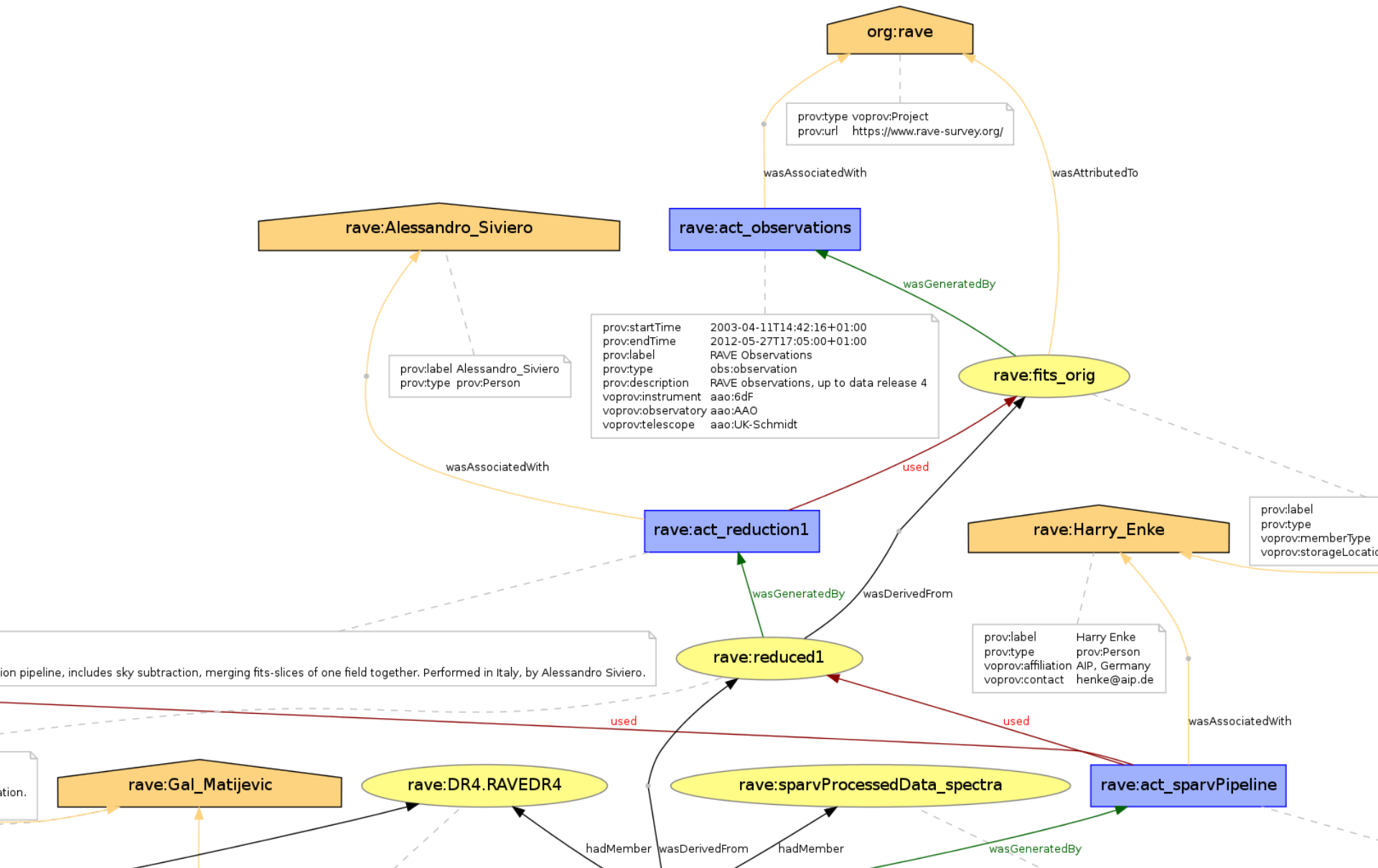


Example: RAVE database tables (nearly complete history)



<https://provenance.ecs.soton.ac.uk/store/documents/84064/>

Example: RAVE database tables



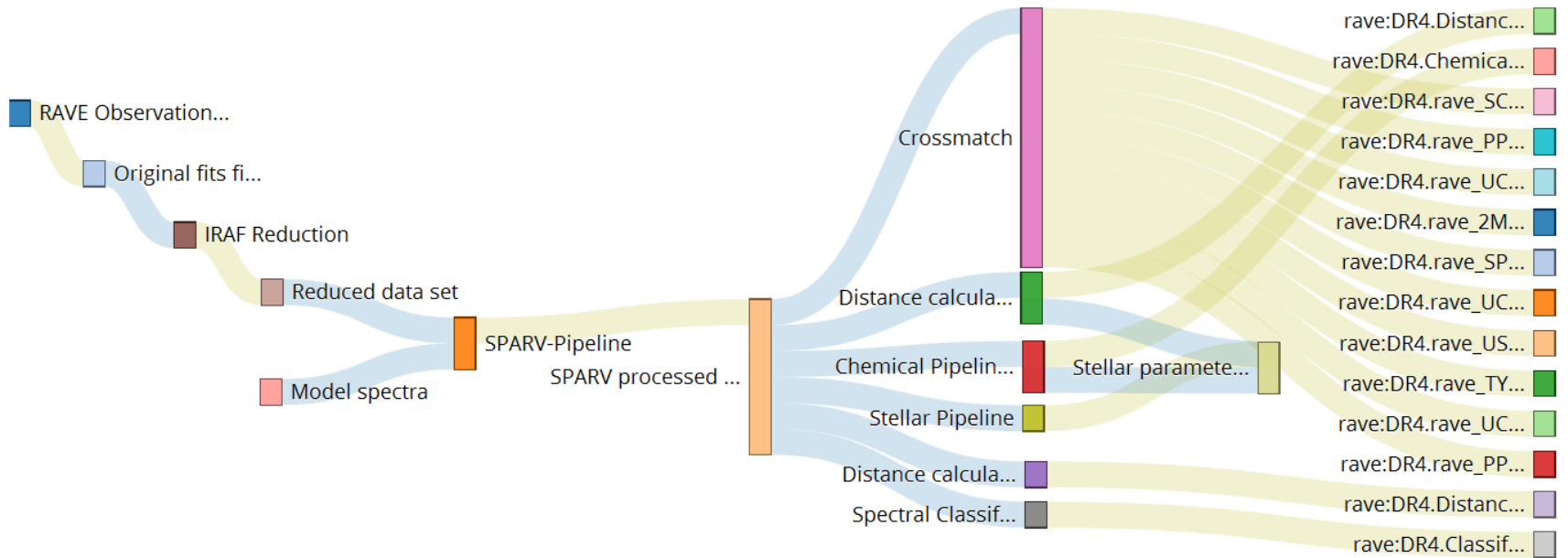
ravedr4 >

Visualizations

Sankey

Select Visible Relations

Created on 09 Jun 2015 at 14:30 by kristinriebe 23 views



Example: RAVE

- RAVE data flow can be modeled in principle using W3C Prov
- RAVE data providers could provide web service:
 - takes id of entry in database table
 - returns complete history of that entry, tracing back through the different steps until the original files from the observation
 - using id as parameter

=> could be queryable via VO services as well?

Discussion

- Is W3C enough? Is it too abstract?
 - In short: Do you agree that we should adopt this?
 - Many implementations already exist, also see:
 - Southampton Provenance Suite, <https://provenance.ecs.soton.ac.uk/> includes validator, converter, visualisation tools
 - Prov Implementation report: <http://www.w3.org/TR/prov-implementations/>
- VO:
 - know most common processes
 - => could predefine input/output of activities (roles)
e.g. image stacking needs n fits-images as input,
one fits-image as output
 - => could predefine standard entities (fits-files, VO-tables, ...)
 - => use PROV-Template system or similar?

Discussion

- Which non-core relations should we include?
 - e.g. actedOnBehalfOf (between agents), wasInformedBy (between activities), ...?
- Collections
 - Can we allow entities to belong to >1 collection?
- What about accessibility of intermediate data products?
 - Ignore non-public data? Flag them?
 - Provide only their header/metadata, more details on request?
- Instrument characteristics + ambient conditions
 - should be modeled elsewhere
 - use provenance data model only for relations between data/activities/agents
- More test cases!

Discussion

- Which non-core relations should we include?
 - e.g. actedOnBehalfOf (between agents), wasInformedBy (between activities), ...?
- Collections
 - Can we allow entities to belong to >1 collection?
- What about accessibility of intermediate data products?
 - Ignore non-public data? Flag them?
 - Provide only their header/metadata, more details on request?
- Instrument characteristics + ambient conditions
 - should be modeled elsewhere
 - use provenance data model only for relations between data/activities/agents
- More test cases!

More discussions

- Discussion session on Wednesday afternoon, 4 pm
- More on W3C: TaPP workshop in Edinburgh,
<http://workshops.inf.ed.ac.uk/tapp2015/>