IVOA Provenance Data Model Version 1.0

IVOA Working Draft 2018-05-30



Working group DM This version http://www.ivoa.net/documents/ProvenanceDM/20180530 Latest version http://www.ivoa.net/documents/ProvenanceDM Previous versions WD-ProvenanceDM-1.0-20170921.pdf WD-ProvenanceDM-1.0-20161121.pdf ProvDM-0.2-20160428.pdf ProvDM-0.1-20141008.pdf Author(s) Kristin Riebe, Mathieu Servillat, François Bonnarel, Anastasia Galkin, Mireille Louys, Markus Nullmeier, Florian Rothmaier, Michèle Sanguillon, Ole Streicher, IVOA Data Model Working Group Editor(s)Kristin Riebe, Mathieu Servillat

Objectives and context



- Data product generation obscure to end user
- Quality, reliability, trustworthiness?
- Usefulness of the data?

Need structured and detailed provenance information

Use cases

- CTA pipeline and data access
- RAVE (Radial Velocity Experiment)
- POLLUX (synthetic stellar spectra service)
- SVOM gamma ray burst / transients
- TAP-based API for images in an archive @CDS
- APPLAUSE database
- MUSE WISE pipeline
- ⇒ Different aspects of Provenance
 - How to **collect** the provenance information
 - How to store this information
 - How to access and visualize the provenance

Draft content

1 Introduction

- 1.1 Goal of the provenance model
- 1.2 Minimum requirements for provenance
- 1.3 Role within the VO architecture
- 1.4 Previous efforts

2 The provenance data model

- 2.1 Overview: Conceptional UML class diagram and introduction to core classes
- 2.2 Model description
 - 2.2.1 Class diagram and VO-DML compatibility

separate section

separate section

- 2.2.2 Entity (and EntityDescription) hereine sub-entities discussed by the second seco
- 2.2.3 Collection
- 2.2.4 Activity and ActivityDescription
- 2.2.5 ActivityFlow _____ link to
- ProvONE?
- 2.2.6 Entity-Activity relations
- 2.2.7 Parameters

2.2.8 Agent

3 Serialization of the provenance data model

3.1 Introduction

3.2 Serialization formats: PROV-N, PROV-JSON and

PROV-XML



- 3.3 PROV-VOTable format
- 3.4 Serialization of description classes for web services
- 3.5 W3C PROV-DM compatible vor alignication fit in W3C serializations?

4 Accessing provenance information

- Appendix B Links to other data models
 - B.1 Links with Dataset/ObsCore Model
 - B.2 Links with Simulation Data Model

Goals

A: Tracking the production history

Find out which steps were taken to produce a dataset and list the methods/tools/software that was involved.

B: Attribution and contact information

Find the people involved in the production of a dataset, that need to be cited or can be asked for more information.

C: Locate error sources

Find the location of possible error sources in the generation of a dataset.

D: Quality assessment

Judge the quality of an observation, production step or dataset.

E: Search in structured provenance metadata

This would allow one to also do a "forward search", i.e. locate derived datasets or outputs.

What is provenance?



- Provenance = Identify how a data product was produced
- Configuration = Identify what detailed options were used
- Contextual information:
 - Instrument Configuration
 - Ambient Conditions
 - Software environment

Core Provenance Data Model



- Core concepts from the W3C PROV recommendations
 - Entity Activity Agent
 - **Relations** and **roles** = provenance information
 - W3C PROV has many more relations
 - IVOA Provenance connected to VO concepts and astronomy needs

Concepts In Astronomy

- Entities: datasets composed of VOTables, FITS files or database tables, or files containing logs, values (spectra, lightcurves), parameters, etc.
- Activities: an observation, a simulation, or processing steps (image stacking, object extraction, etc.).
- Agents: the people involved can be individual persons (observer, publisher...), groups or organisations.
- Connections to existing VO concepts
 - Entity <—> Dataset (Curation, DataID), ObsCore, SimDM DataObject
 - Activity <---> SimDM (Resource, Experiment)
 - Agent <--> Party, Contact
- Connections to external concepts (PROV, ProvONE, DOI, ORCID, ...)

Minimum requirements

- 1. Provenance information must be stored in a **standard model**, with **standard serialization formats**.
- 2. Provenance information must be **machine readable**.
- 3. Provenance data model classes and attributes should be **linked to other IVOA concepts** when relevant (DatasetDM, ObsCoreDM, SimDM, VOTable, UCDs...).
- 4. Provenance information should be **serializable into the W3C provenance standard formats** (PROV-N, PROV-XML, PROV-JSON) with minimum information loss.
- 5. Provenance metadata must contain information to find immediate **progenitor(s)** (if existing) for a given entity, i.e. a dataset.
- 6. An entity must point to the activity that generated it (if the activity is recorded).
- 7. Activities must point to input entities (if applicable).
- 8. Activities may point to output entities.
- 9. Provenance information should make it possible to derive the **chronological** sequence of activities.
- 10. Provenance information can only be given for **uniquely identifiable entities**, at least inside their domain.
- 11. Released entities should have a main contact.
- 12. It is recommended that all activities and entities have **contact information** and contain a (short) **description** or link to a description.

IVOA Provenance Data Model diagram



Parameter (section 2.2.7)



Specialization of entities



- Entity: An entity is a physical, digital, conceptual, or other kind of thing with some fixed aspects; entities may be real or imaginary. (W3C PROV definition)
- **Data**: digital, machine-readable information in some content that will be used/transformed/analyzed... Could be a cell or a column in a table, a file, an image, a cube... The **value** attribute could directly transport this information, or the **location** attribute directs to it.
- **Document**: information presented in a human readable form.
- Thing: a physical object, such as a tool, an instrument, a telescope...
- ActivityDescription/Plan: explanations on the activity, main contact(s), inputs expected, output produced...
- **ConfigFile:** file containing configuration information for the activity.
- **Parameter:** configuration of the activity before execution as a key=value parameter in the IVOA framework.
- AmbientConditions: common, prevailing, and uncontrolled atmospheric and weather conditions in a room or place that influence the activity
- InstrumentConfiguration: description of how the structure of the instrument and how it is set up before an activity starts.
- **ExecutionEnvironment**: describes a particular execution platform, such as an operating system or a database management system. Execution environments are used to describe the context in which the execution of an activity takes place. Execution environments could also describe the computing hardware of a system.

Specialized relations



Serializations - W3C PROV formats

{

```
<previdocument xmlns:ctadata="ivo://vopdc.obspm/cta#" xmlns:ctajob</pre>
  <prev:activity prov:id="ctajobs:ctbin">
    cprov:startTime> 2016-03-13T23:44:46 </prov:startTime>
    corov:endTime> 2016-03-13T23:44:56 
  </prov:activity>
  sect prov:id="cta:consortium">
    sprov:type xsi:type="xsd:string">Organization </prov:type>
  </proviagent>
  v:wasAssociatedWith>
    <proviactivity proviref="ctajobsictbin" />
    cyrov:agent prov:ref="cta:consortium" />
  </prov:wasAssociatedWith>
  prov:entity prov:id="uwsdata:parameters/inobs" />
  <prov:used>
    <prev:activity prov:ref="ctajobs:ctbin" />
    sprov:entity prov:ref="uwsdata:parameters/inobs" />
  </proviused>
  ov:entity prov:id="uwsdata:results/outcube" />
  prov:wasGeneratedBy>
    sentity prov:ref="uwsdata:results/outcube" />
    oprov:activity prov:ref="ctajobs:ctbin" />
  </prov:wasGeneratedBy>
  prov:wasDerivedFrom>
    sprov:generatedEntity prov:ref="uwadata:results/outcube" />
    <previusedEntity proviref="uwsdata:parameters/inobs" />
  </prov:wasDerivedFrom>
  entity prov:id="uwsdata:results/logfile" />
  prov:wasGeneratedBy>
    <previentity proviref="uwsdata:results/logfile" />
    viactivity proviref="ctajobs:ctbin" />
  </prov:wasGeneratedBy>
  <previgeneratedEntity proviref="uwsdata:results/logfile" />
    <prev:usedEntity prov:ref="uwsdata:parameters/inobs" />
  </prov:wasDerivedFrom>
</providocument>
```

```
- wasAssociatedWith: {
   - :idl: {
         prov:agent: "cta:consortium",
         prov:activity: "cta:anactools v1.1"
     }
  },
- agent: (
   - cta:consortium: {
         prov:type: "Organization"
  1,
- entity: {
     uwsdata:results/fit_results: { },
     uwsdata:results/configfile: { },
     uwsdata:results/butterfly: { },
     uwsdata:results/spectrum plot: { },
     uwsdata:results/spectrum: { }
  1,
- prefix: {
     uwsdata: "https://voparis-uws-test.obspm.fr/rest
     cta: "http://www.cta-observatory.org#",
     voprov: "http://www.ivoa.net/ns/voprov#"
  1,
- activity: {
   - cta:anactools vl.1: {
         prov:startTime: "2016-04-07T00:26:00",
         prov:endTime: "2016-04-07T00:27:15"
     ŀ
  1,
- wasGeneratedBy: {
   - :id5: {
         prov:entity: "uwsdata:results/butterfly",
         prov:activity: "cta:anactools v1.1"
     ł,
   - :id4: {
         prov:entity: "uwsdata:results/fit results",
         prov:activity: "cta:anactools v1.1"
     },
```

Serializations - VOTable

```
<?xml version="1.0" encoding="UTF-8"?>
<VOTABLE version="1.2" xmlns="http://www.ivoa.net/xml/VOTable/v1.2"
    xmlns:ex="http://www.example.com/provenance"
    xmlns:ivo="http://www.ivoa.net/documents/rer/ivo/"
    xmlns:voprov="http://www.ivoa.net/documents/dm/provdm/voprov/"
    xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
    xsi:schemaLocation="http://www.ivoa.net/xml/VOTable/v1.2 http://www.ivoa.net/xml/VOTable/VOTable-1.2.xsd">
<RESOURCE type="provenance">
    <DESCRIPTION>Provenance V0Table</DESCRIPTION>
    <TABLE name="Usage" utype="voprov:used">
        <FIELD arraysize="*" datatype="char" name="activity" ucd="meta.id" utype="voprov:Usage.activity"/>
        <FIELD arraysize="*" datatype="char" name="entity" ucd="meta.id" utype="voprov:Usage.entity"/>
        <DATA>
            <TABLEDATA>
                <TR>
                    <TD>ex:Process1</TD>
                    <TD>ivo://example#DSS2.143</TD>
                </TR>
            </TABLEDATA>
        </DATA>
    </TABLE>
    <TABLE name="Generation" utype="voprov:wasGeneratedBy">
        <FIELD arraysize="*" datatype="char" name="entity" ucd="meta.id" utype="voprov:Generation.entity"/>
        <FIELD arraysize="*" datatype="char" name="activity" ucd="meta.id" utype="voprov:Generation.activity"/>
        <DATA>
            <TABLEDATA>
                <TR>
                    <TD>ivo://example#Public NGC6946</TD>
                    <TD>ex:Process1</TD>
                </TR>
            </TABLEDATA>
        </DATA>
    </TABLE>
    <TABLE name="Activity" utype="voprov:Activity">
        <FIELD arraysize="*" datatype="char" name="id" ucd="meta.id" utype="voprov:Activity.id"/>
        <FIELD arraysize="*" datatype="char" name="name" ucd="meta.title" utype="voprov:Activity.name"/>
        <FIELD arraysize="*" datatype="char" name="start" ucd="" utype="voprov:Activity.startTime"/>
        <FIELD arraysize="*" datatype="char" name="stop" ucd="" utype="voprov:Activity.endTime"/>
        <DATA>
            <TABLEDATA>
                <TR>
                    <TD>ex:Process1</TD>
```

Serialization context



Serializations - ActivityDescription

<resource id="gammapy_maps" name="gammapy_</th><th><pre>maps" type="meta" utype="voprov:ActivityDescription"></resource>			
<pre><description>Use gammapy to generate a <!-- Service Descriptor--> <param datatype="char</pre></th><th>count map from a list of observations</DESCRIPTION></th><th>mmany mans" name="access!RL"/></description></pre>			
<pre><param arraysize="*" datatype="char" name="standardID" value="ivo://ivoa.net/std/SODA#1.0"/> <!-- Activity Description--></pre>			
<param arraysize="*" datatype="char" name="type" utype="voprov:ActivityDescription.type" value="None"/> <param arraysize="*" datatype="char" name="subtype" utype="voprov:ActivityDescription.subtype" value="None"/> <param arraysize="*" datatype="char" name="annotation" utype="voprov:ActivityDescription.version" value="None" version"=""/> <param arraysize="*" contact_name"="" datatype="char" name="doculink" utype="voprov:Agent.name" value="Julien Lefaucheur"/>			
<pre><param arraysize="*" datatype="cnar" name="contact_email" value=" utype=" voprov:agent.email"=""/> </pre>			
UWS parameters (Provenance Entities or Parameters)		DataLink Service Descriptor	
<pre><group name="InputParams"> <param arraysize="*" datatype="double" id="obs_ids" n<br="" name="obs_ids" ra"="" value="47802 47803 47804 <pre></pre> <pre></pre> <pre></pre> </pre> </pre></th><th>DWS Job Description Language
Provenance ActivityDescription</th></tr><tr><th></PARAM>
<PARAM ID="/><param [<="" datatype="double" id="Dec" th=""/><th>ame="RA" value="329.7169379" unit="deg"></th><th></th></group></pre>		ame="RA" value="329.7169379" unit="deg">	
<pre><param <="" arraysize="*" da="" id="nxpix" pre=""/> <pre><pre><pre><pre>Council</pre> </pre> <pre><pre><pre><pre><pre><pre><pre><pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre>	<pre><!-- Used Entities--> </pre> <pre> <pre> <pre> <pre> </pre></pre></pre></pre>		
<values> <min value="0"></min> <max value="1000"></max></values>	<pre><param arraysize="*" datatype="char" name="role" utype="voprov:UsedDescription.role" value="DL3"/> <param <="" <param="" arraysize="*" datatype="char" name="content_type" pre="" utype="voprov:EntityDescription.content_type" v="" value=""/></pre>		
 <param arraysize="*" da<="" id="nypix" th=""/> <th> </th> <th></th>	 		
<param <br="" datatype="float" id="binsz"/>	<pre><!-- Generated Entities / UWS results--> <group name="Generated" utype="voprov:WasGeneratedBy"></group></pre>		
	<pre><description>Count map</description> <param arraysize="*" datatype="char" name="role" utype="voprov:UsedDescription.role" value="DL4 image"/> <param arraysize="*" datatype="char" name="default" utype="voprov:Entity.id" value="count_map.fits"/> <param <="" arraysize="*" datatype="char" name="content_type" pre="" utype="voprov:EntityDescription.content_type" v=""/></pre>		
<pre> </pre>			

Conclusions



- Many sections in the draft are stable
- ProvSAP and ProvTAP moved to DAL drafts
 - Implementation note based on many use cases
- Still some open questions:
 - modeling of Parameter

 - relations with Description classes
 - mapping for valid W3C serialization without loss

Next steps

move to RFC track before next IVOA meeting