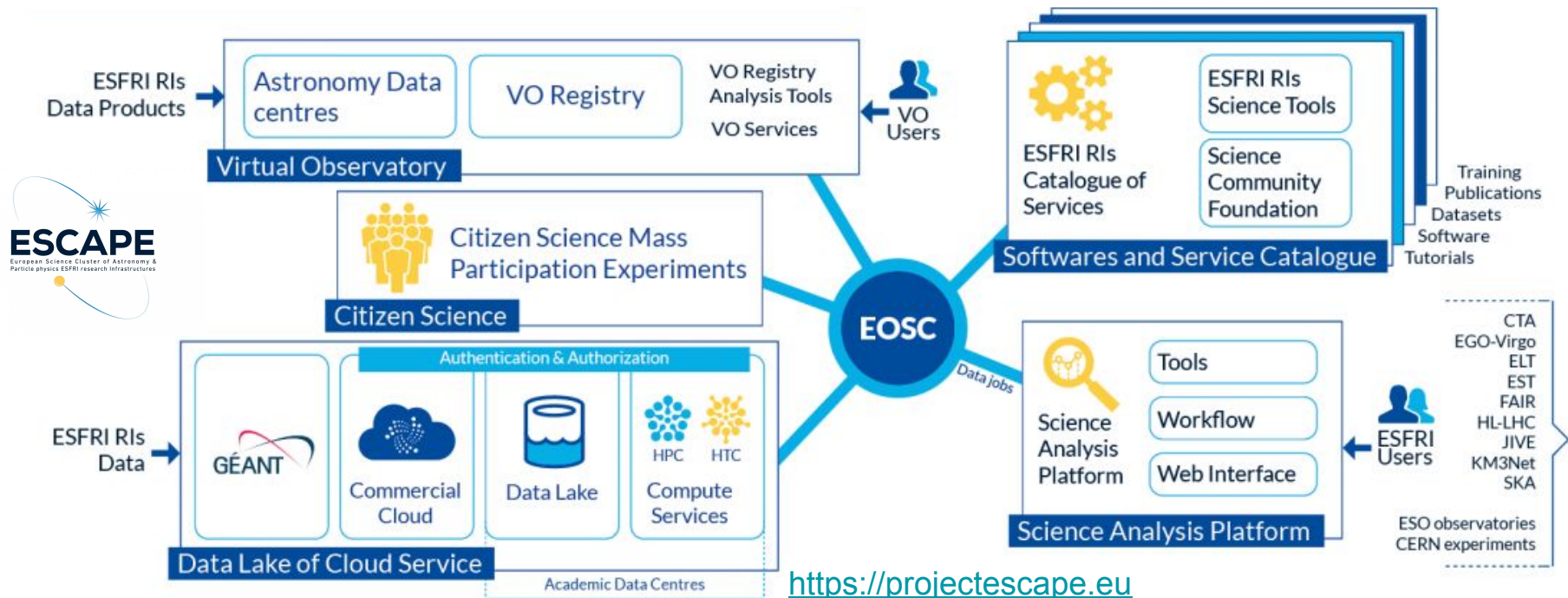# Provenance activities in the European project ESCAPE

**EOSC** – The **E**uropean **O**pen **S**cience **C**loud is a cloud for research data in Europe allowing for universal access to data;

**ESCAPE** – « The **E**uropean **S**cience **C**luster of **A**stronomy & **P**article Physics **E**SFRI Research Infrastructures » answers the EOSC ambition in bringing People, Data, Services, Training, Publications, Projects & Organisations, all together in an integrated and federated environment.



https://projectescape.eu

# Provenance 2-days workshop on use cases

The objective was to collect the requirements of ESFRI projects in order to build the road map of future developments concerning provenance.

https://indico.in2p3.fr/event/21913/page/2641-summary

- Presentation of the **model**, and the associated implementations on provenance **capture**, **storage**, **access** and **visualization**
- Presentation of **each project use case**, following the pattern:
    - Project, context, data products
    - What is the relevant provenance information kept (or to be kept)?
    - How is it kept (or will be)?
    - What provenance will/should the end user see?
- General **discussion** on questions raised during the workshop
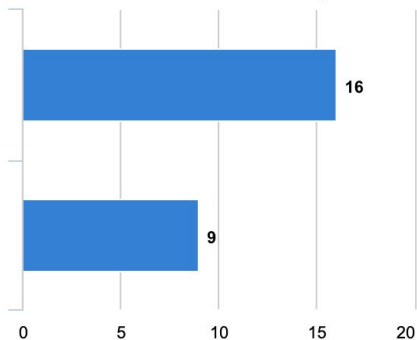
# BoF at ADASS XXX on practical provenance

- **Introduction** to IVOA provenance
- **Questionnaire** posted and filled by attendants and others before
- **3 posters** : CTADIRAC (P9-250), VizieR (P9-216), OPUS (P9-89)
- **Discussion** (not really lively despite interest in the questionnaire...)
  - Minimum provenance: last step
  - Serialization / format: YAML, JSON are prefered ? VOTABLE (with model mapping)
  - Reproducibility can use provenance information:
    - Can we translate data flow language like CWL to IVOA provenance and reverse
    - Software parameters values may be stored with activities
    - Find out datasets produced with a given version of software
    - Proposal to use code ID to identify software in ActivityDescription
  - Visualisation using voprov library (see posters)
  - Access using VO protocols

# Some answers to the questionnaire

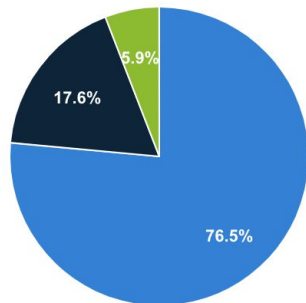## 1/ Are you a provider or a user of provenance information (or both)?

Chart options »



| provider | 16 |
|----------|-----|
| user | 9 |

## 2/ Do you already attach provenance information to the data you create or provide?

Chart options »



76.5%
17.6%
5.9%

| yes | 13 |
|-----|-----|
| no | 3 |
| not applicable | 1 |

## 3/ Do you use the IVOA Provenance data model ?

Chart options »



35.3%
64.7%

| yes | 6 |
|-----|-----|
| no | 11 |

## 5/ What is the most important goal in relation to provenance for you?

Chart options »



| Traceability | 11 |
|--------------|-----|
| Quality / Reliability | 8 |
| Acknowledgement | 1 |
| Debugging | 8 |
| Reproducibility | 14 |

# Astronomy context for data generation

Data product generation **obscure** to end user



Need **structured** and **detailed** provenance information!

# FAIR principles for data sharing

https://www.go-fair.org/fair-principles

## Findable

**F1**. (Meta)data are assigned a globally unique and persistent **identifier**
**F2**. Data are described with rich metadata
**F3**. Metadata clearly and explicitly include the **identifier** of the data they describe
**F4**. (Meta)data are **registered** or **indexed** in a searchable resource

## Accessible

**A1**. (Meta)data are retrievable by their **identifier** using a **standardised** communications **protocol**
  **A1.1**. The **protocol** is open, free, and universally implementable
  **A1.2**. The **protocol** allows for an authentication and authorisation procedure, where necessary
**A2**. Metadata are accessible, even when the data are no longer available

## Interoperable

**I1**. (Meta)data use a formal, accessible, shared, and broadly applicable **language** for knowledge representation.
**I2**. (Meta)data use **vocabularies** that follow FAIR principles
**I3**. (Meta)data include **qualified** references to other (meta)data

## Reusable  (+ Reproducible?)

**R1**. (Meta)data are richly described with a plurality of **accurate** and **relevant** attributes
  **R1.1**. (Meta)data are released with a **clear** and accessible data usage **license**
  **R1.2**. (Meta)data are associated with detailed **provenance**
  **R1.3**. (Meta)data meet domain-relevant community **standards**

# Why recording structured provenance?

- **FAIR principles** (Findable, Accessible, Interoperable, Reusable)
  - https://www.go-fair.org/fair-principles/
  - "**rich**" metadata, following standard data model, protocols and formats
  - "**detailed provenance**"
- **Quality** / **Reliability** / **Trustworthiness** of the products
- **Reproducibility requirement** in projects
  - Be able to rerun each activity (maybe testing and improving each step)
  - Not necessary to keep every intermediate file that is easily reproducible (possible gain on disk space and costs)
- **Debugging**
  - Not necessary to restart from scratch: locate in the provenance tree the faulty parts or the products to be discarded

→ We often realize too late that there are missing elements or links in the provenance. The capture of the provenance should be as detailed as possible and as naive as possible (simply record what happens).

# **I**nternational **V**irtual **O**bservatory **A**lliance

## IVOA Documents



http://www.ivoa.net/documents/ProvenanceDM/

## IVOA Provenance Data Model
## Version 1.0

### IVOA Recommendation 11 April 2020

**Interest/Working Group:**
> http://www.ivoa.net/twiki/bin/view/IVOA/IvoaDataModel

**Author(s):**
> **Mathieu Servillat, Kristin Riebe, Catherine Boisson, François Bonnarel, Anastasia Galkin, Mireille Louys, Markus Nullmeier, Nicolas Renault-Tinacci, Michèle Sanguillon, Ole Streicher**
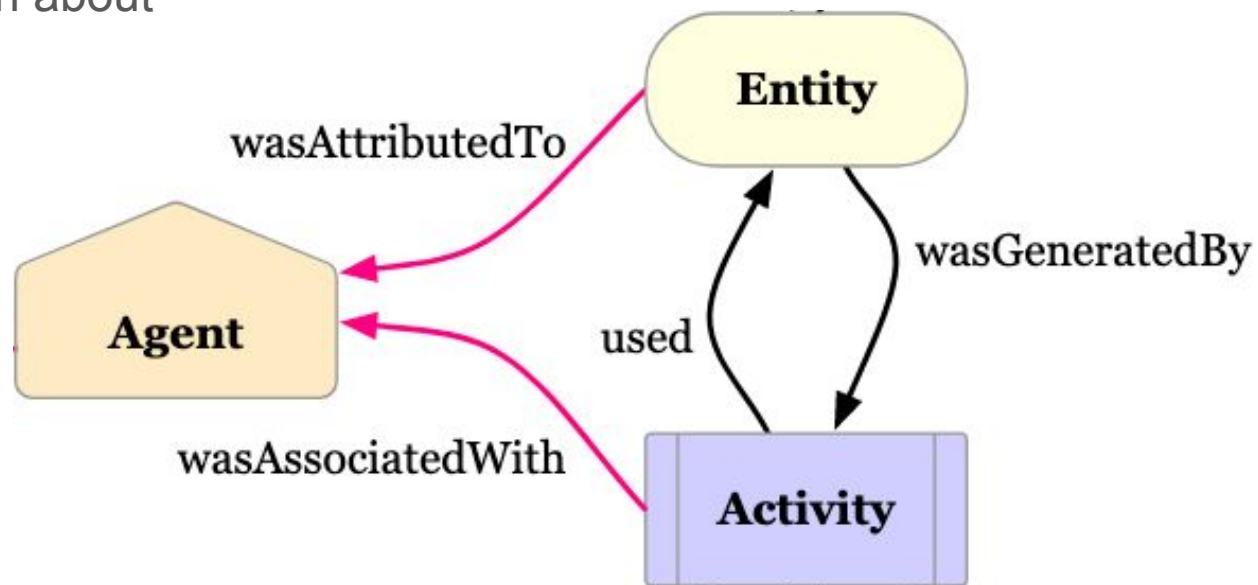
**Editor(s):**
> **Mathieu Servillat**

# W3C Provenance definition
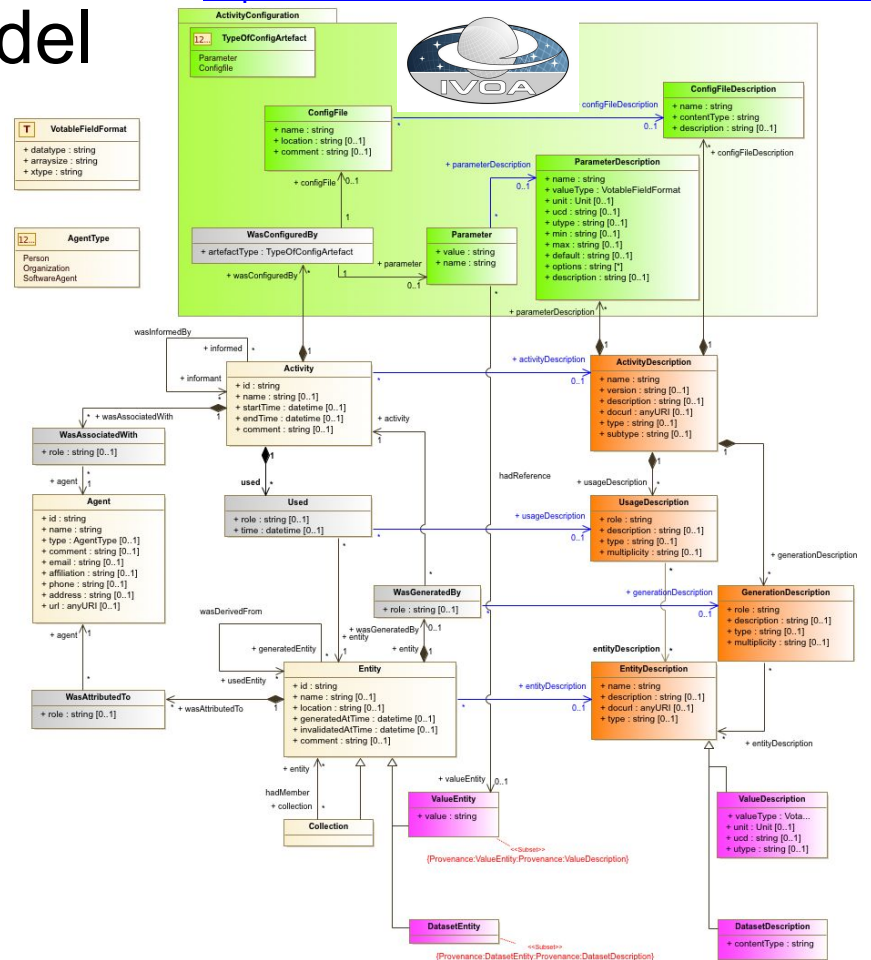
**W3C PROV (PROV-DM, 2013)**

Provenance is information about
**entities**, **activities**,
and people (**agents**)
involved in producing a
piece of data or thing,
which can be used to
form assessments
about its **quality**,
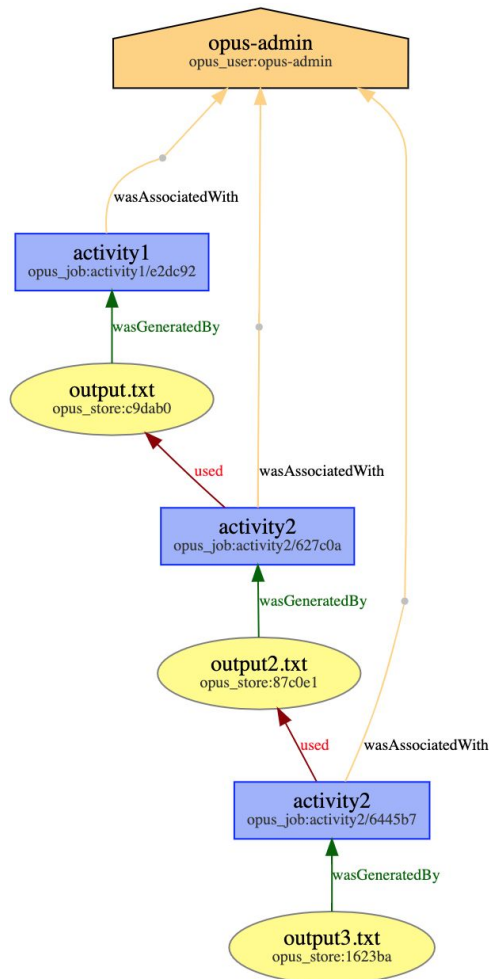**reliability** or
**trustworthiness**.

W3C
Word Wide Web Consortium

# IVOA Provenance Data Model

**Recommendation 1.0 in April 2020**

- Adds "**Description**" classes
- Adds "**Configuration**" classes
- Plugged in with
  - **VO** data models and concepts
    (UCD, VOUnit, VOTable…)
  - **VO** access protocols
    (ProvTAP, ProvSAP)
- Serializations
  - W3C PROV (XML, JSON, SVG…)
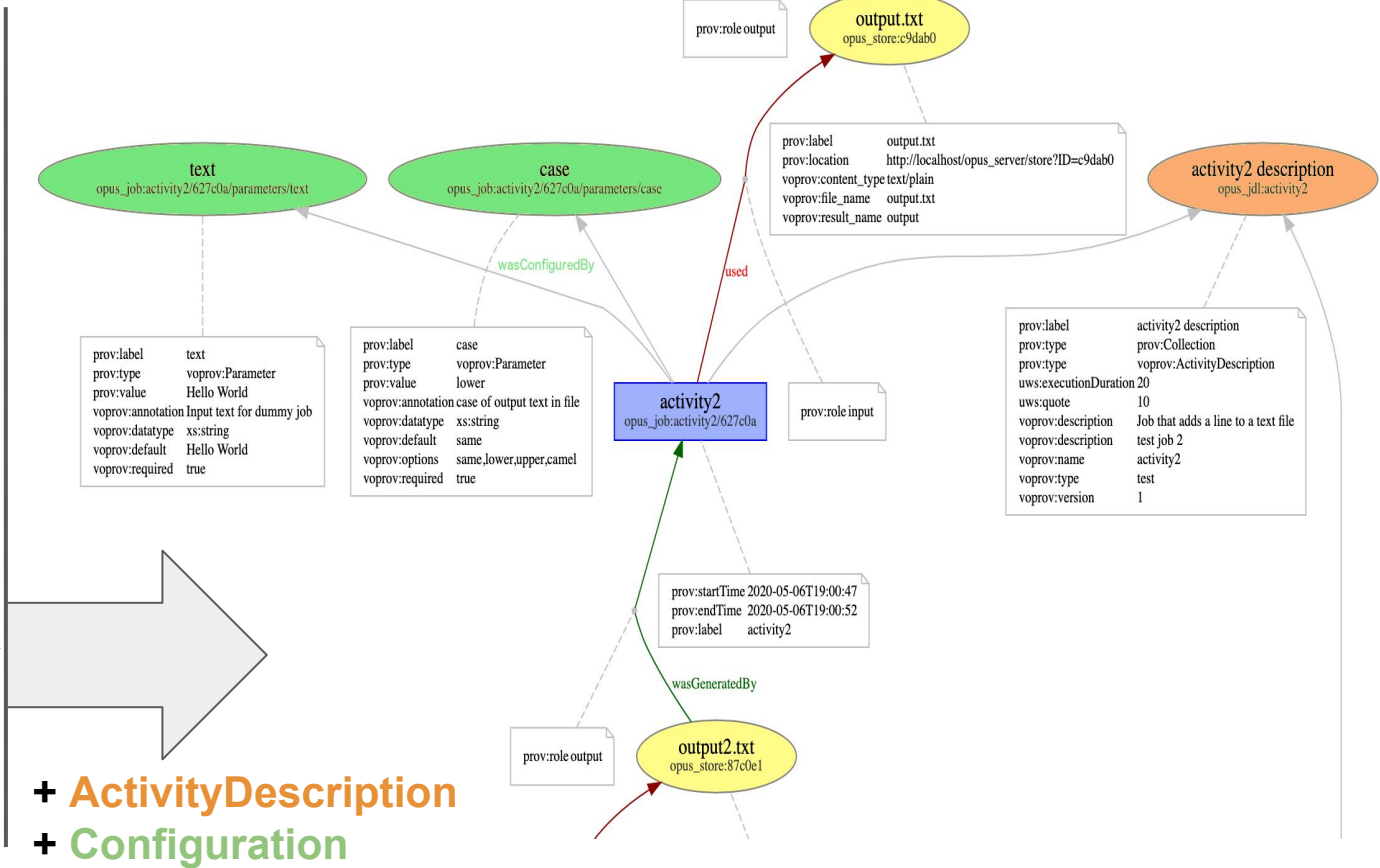  - **VO** specific (VOTable)

# Provenance graph

Provenance is :
- a **chain** of activities and entities (used and generated)
- that occured in the **past**

Using the **core data model**, some goals are achieved:
- Unique identifiers
- Traceability
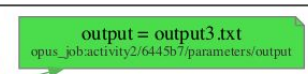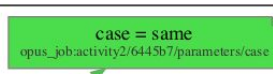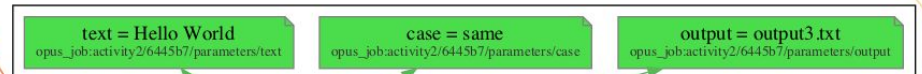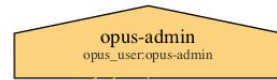- Contact and Acknowledgement

By using the **full IVOA data model**, more questions are answered:
- What **happened** during each **activity**?
- How was the **activity tuned** to be executed properly?
- What **kind of content** is in the **entities** ?

opus-admin
opus_user:opus-admin

wasAssociatedWith

activity1
opus_job:activity1/e2dc92

wasGeneratedBy

output.txt
opus_store:c9dab0

used

activity2
opus_job:activity2/627c0a

wasGeneratedBy

output2.txt
opus_store:87c0e1

used

wasAssociatedWith

activity2
opus_job:activity2/6445b7

wasGeneratedBy

output3.txt
opus_store:1623ba

text
opus_job:activity2/627c0a/parameters/text

case
opus_job:activity2/627c0a/parameters/case

wasConfiguredBy

| prov:label | text |
|---|---|
| prov:type | voprov:Parameter |
| prov:value | Hello World |
| voprov:annotation | Input text for dummy job |
| voprov:datatype | xs:string |
| voprov:default | Hello World |
| voprov:required | true |

| prov:label | case |
|---|---|
| prov:type | voprov:Parameter |
| prov:value | lower |
| voprov:annotation | case of output text in file |
| voprov:datatype | xs:string |
| voprov:default | same |
| voprov:options | same,lower,upper,camel |
| voprov:required | true |

prov:role output

output.txt
opus_store:c9dab0

wasGeneratedBy

| prov:label | output.txt |
|---|---|
| prov:location | http://localhost/opus_server/store?ID=c9dab0 |
| voprov:content_type | text/plain |
| voprov:file_name | output.txt |
| voprov:result_name | output |

used

activity2
opus_job:activity2/627c0a

prov:role input

activity2 description
opus_jdl:activity2

| prov:label | activity2 description |
|---|---|
| prov:type | prov:Collection |
| prov:type | voprov:ActivityDescription |
| uws:executionDuration | 20 |
| uws:quote | 10 |
| voprov:description | Job that adds a line to a text file |
| voprov:description | test job 2 |
| voprov:name | activity2 |
| voprov:type | test |
| voprov:version | 1 |

| prov:startTime | 2020-05-06T19:00:47 |
|---|---|
| prov:endTime | 2020-05-06T19:00:52 |
| prov:label | activity2 |

wasGeneratedBy

prov:role output

output2.txt
opus_store:87c0e1
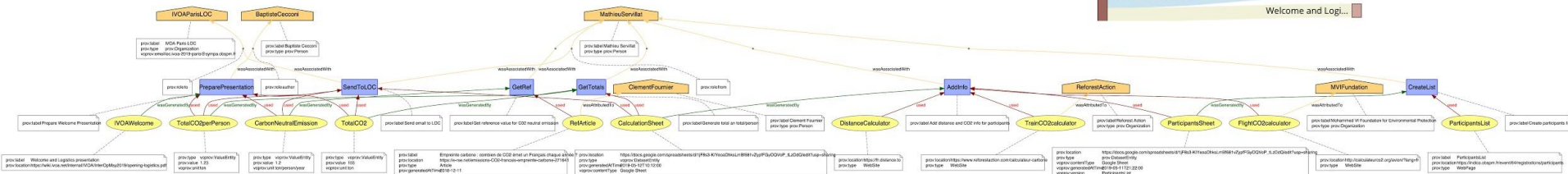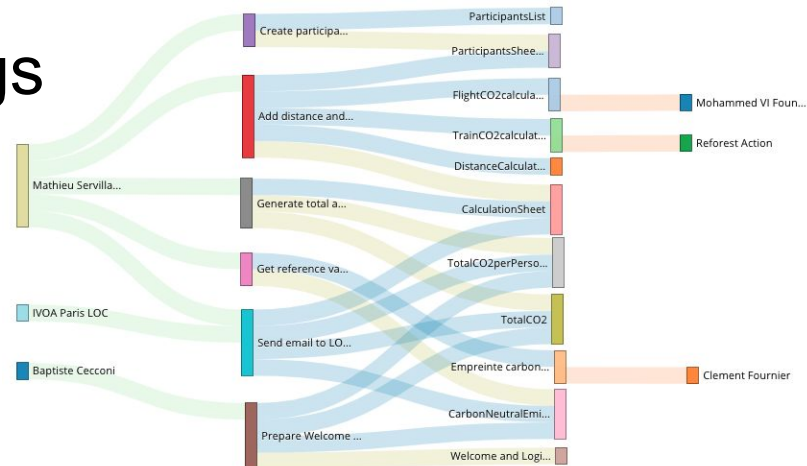
**+ ActivityDescription**
**+ Configuration**

# Carbon footprint of IVOA meetings



## IVOA Paris 2019 - 132 participants

**155 tons of $CO_2$** for travels only

Check the detailed provenance of the calculations:

https://frama.link/CO2prov

## IVOA Sydney 2020 virtual - 207 participants

**300 grams of $CO_2$** for visioconferences

(~1500 h.person of visio estimated)

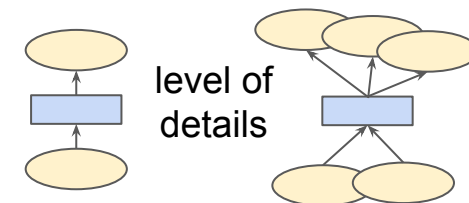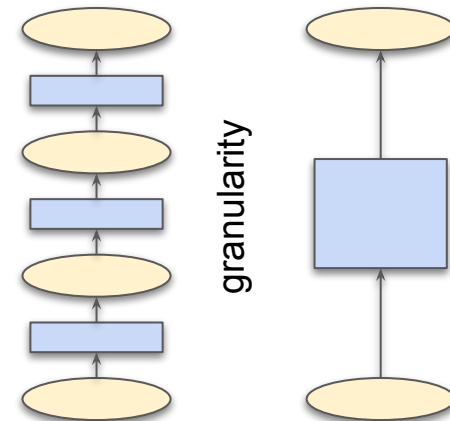https://greenspector.com/en/which-video-conferencing-mobile-application-to-reduce-your-impact/

# Applying the model

## Different contexts in use cases

- Two flavours:
  - **on-top** (data products/collection already exist)
  - **inside** (save provenance information during the processing)
- **Identifiers:** unique and without meaning
- **Granularity** (what steps? what objects?)
- **Level of details** (descriptions? configuration?)

## Different steps in provenance management

- How to **capture** the provenance information
- How to **store** this information
- How to **access** it
- How to **visualize** a provenance graph

# Some terminology

- **full provenance**: graph/tree/chain of activities and entities up to the raw data. This information is not hosted by the entities themselves (stored on an external server? as separate files?)

- **minimum provenance**: attached to an entity, list of keywords that gives some context and info on **last activity** (general process/workflow, software versions, contacts...).
  *Note: it would be interesting to include used entities, so that a full provenance may be reconstructed from each minimum provenance. However, such information on what was used may not be kept, or may not be complete.*

- **end-user/specific "provenance"**: attached to an entity, list of keywords or data that provides **key information to use/analyse** the entity (e.g. for CTA: event class, event type, telescope configuration, sky conditions, reco method...)
  *Note: may be extracted from full provenance (some parameters or entities generated at a given step), but it is considered as **data** here. Reversely, this specific "provenance" information may be a source of information to be mapped in the standard in order to fill the full provenance graph with more details.*

1. Defining the content of a **minimum provenance**

   - List of keywords related to the last activity and context

2. **Serializing** provenance

   - Both human and machine readable
   - Explore W3C formats, and YAML / VOTABLE / VOEvent formats

3. Provenance and **workflows**

   - Workflow information simply attached to provenance as used entities
   - Links with CWL, mapping

4. From provenance "**on-top**" to provenance "**inside**"

   - How to map specific provenance information into the IVOA model?
   - How to introduce provenance **capture** inside a pipeline?

5. Provenance **storage**

   - Database, interface and ingestion

6. Provenance **exploration** and **visualization**

   - Access protocols (ProvTAP, ProvSAP)
   - `voprov` Python package, or other tools