

## Background

Provenance is information about the origin of data. It is one of the requirements for fair data (see [1]). Essential for reviewable data.

**Whenever data is used, provenance is needed together with the data to review it for the given usage**

- Data formats for combining provenance & data exists (for example hdf5 [2], data packages [3])
- Cloud based services exists (see for example [4])
- Metadata management systems are in use (see various MMS or PDM)

**Scenario 1: Data usage of after some years/decades (see for ex. aviation regulations for record keeping)**

- Maintenance matters
- Technology & tools will change

**Scenario 2: Data exchange between stakeholders without common infrastructure**

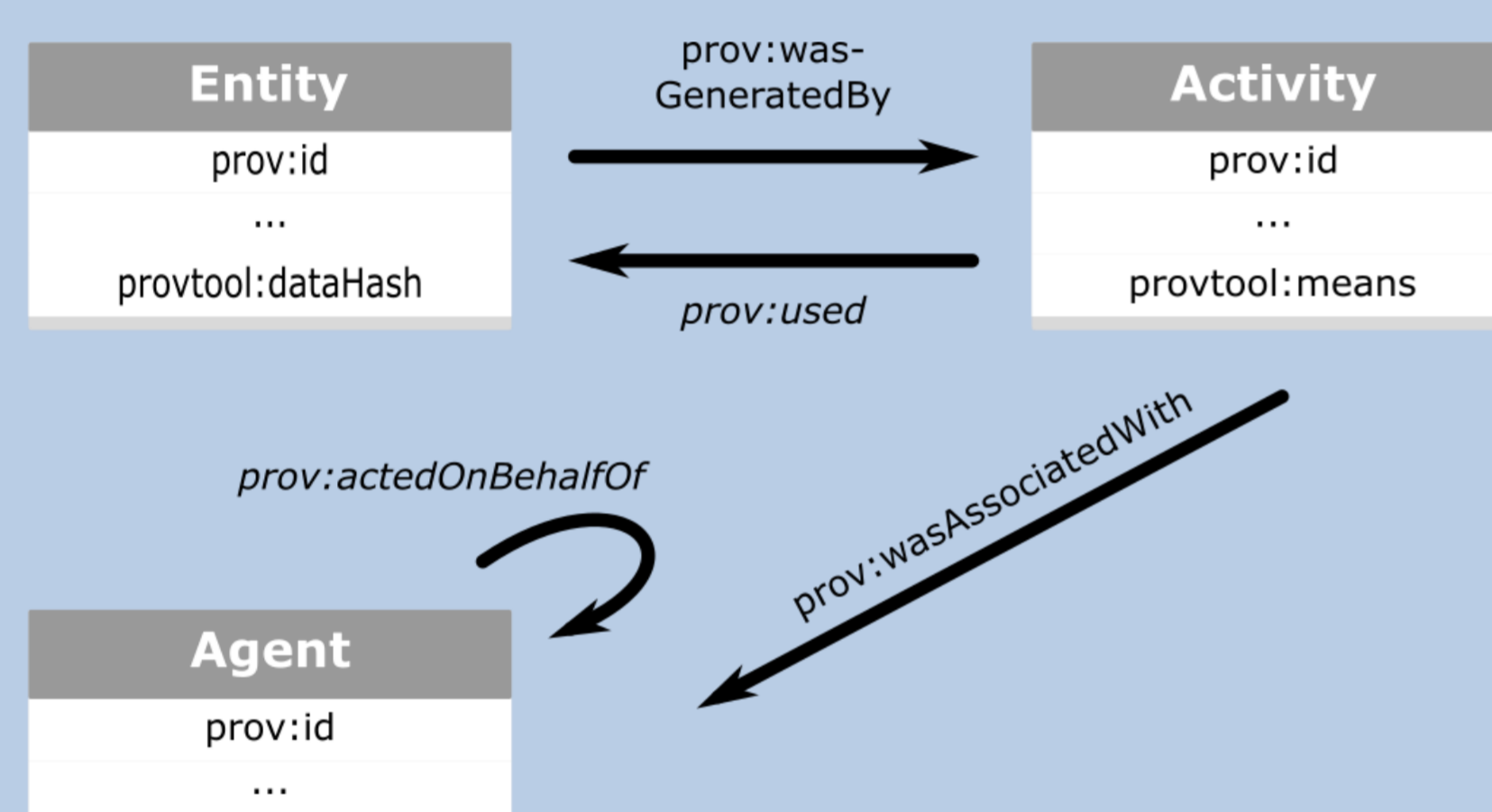
- No direct synchronization of data possible
- No direct synchronization wanted

## Provenance model used

**Model:** Reduced W3C prov model [5]: Activities, agents, entities

**Additional properties**

- provtool: datahash, :means



Reduced provenance model based on the W3C prov model.

## Approach

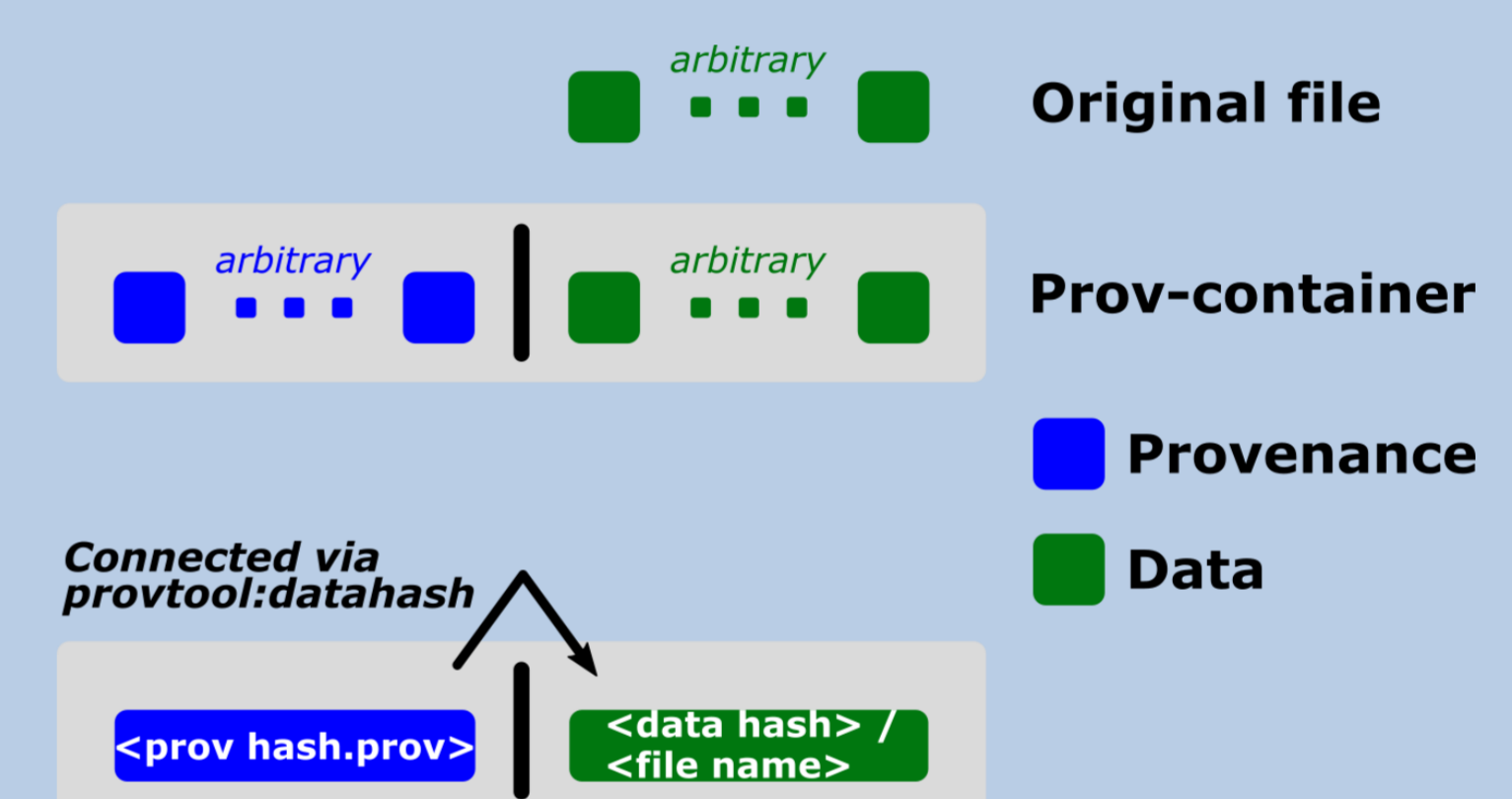
**Assumption:** Data & provenance exchange somehow possible, concrete technology unknown (time & space)

**Scope:**

- No additional services/tools needed (but support possible)
- Data as is
- Not necessarily local
- Usable by nearly any existing DMS
- Lookup and search out of scope
- Provenance entity id = provenance hash

**Provenance container**

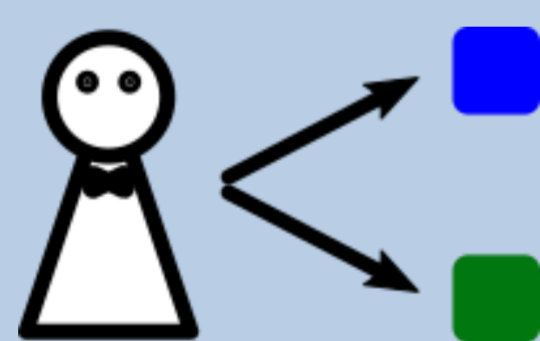
- Data kept as binary
- Provenance in utf-8 encoded json according to [6] and reduced model for current activity only
- datahash as reference from provenance to data
- Serialized independently with hash sum as id & name



## Usage

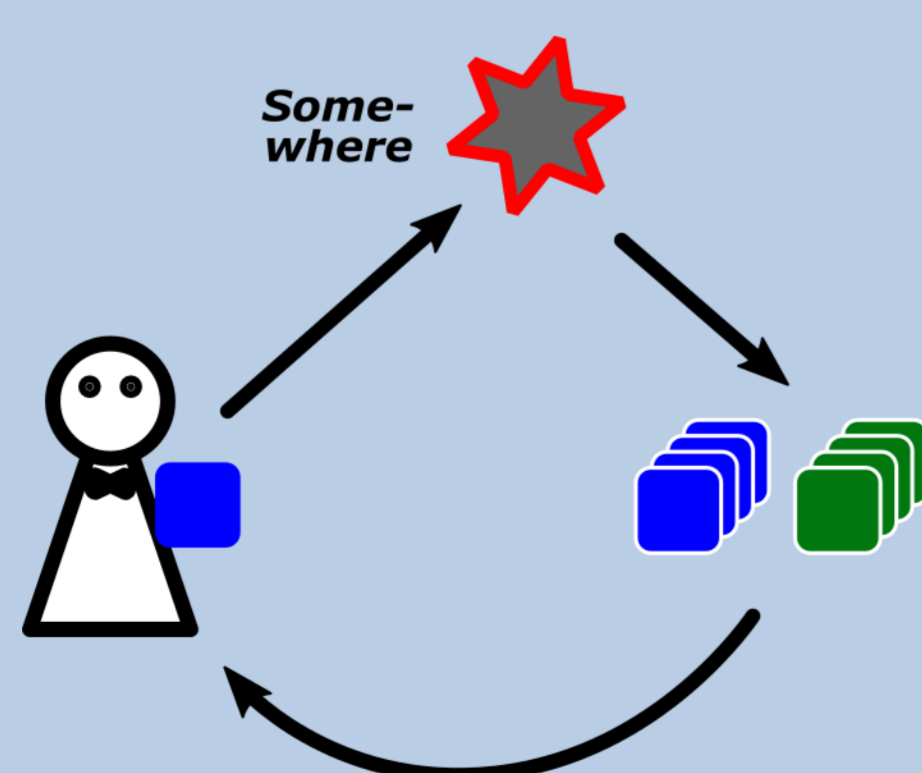
**Production**

- Agent generates data and provenance in text format



**Usage**

- User starts with provenance
- Queries used



## Implications

**Hash as id: Content addressable**

- Storage location is irrelevant
- Nearly any storage technology will do

**Effective Immutability**

- Modification detectable due to hash id
- Data generation using an existing container seals the previous container

**Self describing within model**

- Provenance container  $\triangleq$  Entity

**Arbitrary DMS possible**

- Use already: Shepard [7, 8]
- File storage & indexing with lucene

**Reference to primary sources**

- Build & traverse provenance graph
- Additional systems may provide easy lookup

**Exchange w/o technical debts**

- Only existing formats and utf8-text
- Nearly any existing exchange format will do

**Provenance first! (see [4] for similar id-based approach)**

**Con:**

Two artefacts instead of one

**Purely organizational effort to use it: Can start today**