

# Digital Preservation

## Authenticity and Provenance

Stefan Pröll  
28.04.2014

- Trust in digital archives
  - Why do archives need to be trustworthy
- Authenticity
  - What does authenticity mean
- Provenance
  - How to map the genealogy of digital data
- Metadata for authenticity and provenance
- Security
  - How to secure archives and their content

- Archives store digital objects. They need to be:
  - Safe
  - Reliable
  - Trustworthy
- Trust is a fundamental property of digital archives
  - Trust is hard to establish
  - Easy to destroy
- How to establish trust?

- What is Authenticity?



<http://weird.cz/ali-g>

- Authenticity: The degree a digital object actually is what it claims to be.
- Authenticity is judged on the basis of evidence

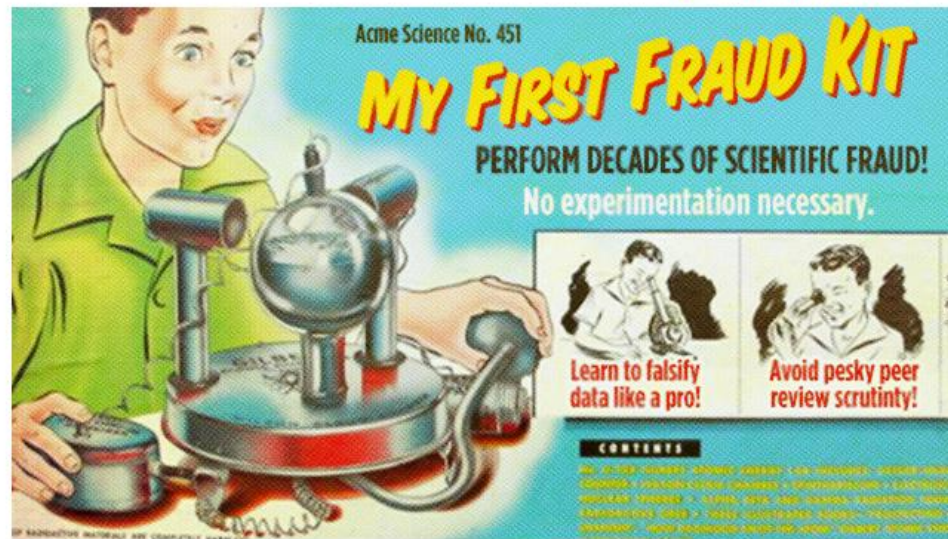
## SCIENTIFIC METHOD / SCIENCE & EXPLORATION

## Epic fraud: How to succeed in science (without doing any)

Envy those who succeed by making up their data? Here's how you can, too!

by **John Timmer** - July 19 2012, 3:30am CEST

SCIENCE POLICY AND EDUCATION 286

 Aurich Lawson

Source: <http://arstechnica.com/science/2012/07/epic-fraud-how-to-succeed-in-science-without-doing-any/>

- Well established for physical objects.
- Example: Paintings
  - Age
  - Certificates
  - Ownership and Provenance
  - Condition
  - Known style, hand writing
  - Repairs or alterations
  - Reference and context



© www.artmarketblog.com

NACHRICHTEN
reporter\*12
FREIZEIT
SCHAUFENSTER
ABO&CLUB
KARRIERE
IMMOBILIEN
SERVICE

Die Presse.com
Kultur
Kunst

Politik
Wirtschaft
Panorama
Kultur
Tech
Sport
Leben
Bildung
Wissenschaft
Gesundheit
Recht
S

## Beltracchi fälschte weit mehr Bilder als bisher bekannt

05.03.2012 | 10:17 | (DiePresse.com)

**Wegen 14 gefälschter Gemälde wurde Wolfgang Beltracchi zu sechs Jahren Haft verurteilt. Nun gab er zu, "ungefähr 50" Bilder gefälscht zu haben.**



Bild vergrößern

Drucken
Senden
+ Merken
Vorlesen
AAA Textgröße
Kommentieren

**MEHR ZUM THEMA:**

Köln: Bis zu sechs Jahre Haft für Kunstfälscher

**AUS DEM ARCHIV:**

Fälscher im Blitzlicht: "Nichts zu verbergen" (03.09.2011)

Der verurteilte Kunstfälscher Wolfgang Beltracchi hat in einem Interview mit dem Nachrichtenmagazin "**Spiegel**" mehr Fälschungen eingeräumt als bisher bekannt. In seinem Leben, so der 61-Jährige, habe er Werke von "ungefähr 50" verschiedenen Künstlern gefälscht. Die genaue Zahl und wo sich die Bilder befänden, wolle er aber nicht bekanntgeben.

Der Fälscher hatte mit seinen Fälschungen einen der größten Skandale auf dem deutschen Kunstmarkt ausgelöst und war im Oktober 2011 wegen 14 gefälschter Gemälde von Künstlern wie Campendonk, Léger und Ernst zu sechs Jahren Haft verurteilt worden. Er hatte die Sammlungen "Knops" und "Jägers" erfunden und daraus angeblich verschollene Originale der klassischen Moderne verkaufte

Beltracchi sagte, er hätte wegen der großen Nachfrage leicht "1000 oder 2000 Stück" seiner Fälschungen auf dem Kunstmarkt absetzen können.

(Ag.)

TagesAnzeiger

Front Zürich Schweiz International Wirtschaft Börse Sport **Kultur** Leben Wissen Auto Blogs Panorama Mehr ▾

Film Fernsehen Bücher Theater Kunst Musik Klassik Bestenlisten Deadline Bildstrecken

## Der Betrüger, den man mögen muss

Keiner hat die Kunstwelt so schön aufs Kreuz gelegt wie Meisterfälscher Wolfgang Beltracchi. Jetzt erzählen ein Dokumentarfilm und eine Autobiografie, wie virtuos er dabei vorging.



Ein genialer Leger falscher Spuren: Szene aus dem Film Beltracchi – «Die Kunst der Fälschung». Foto: Senator Film Verleih.



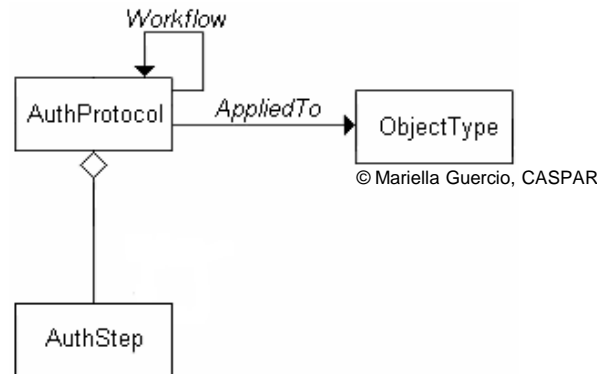
Drei echte Beltracchi aus der fiktiven Sammlung Jäger: Gefälschte Gemälde von Kees van Dongen, Max Ernst und Max Pechstein (v.l.). Foto: Katja Hoffmann (Laif).



- How to assess authenticity?
  - Analyze metadata accompanying the digital object
  - Examine checksums (fixity)
  - Compare to redundant copies
  - Investigate the context of the object
  - Track provenance (see later...)
  - Use signatures and encryption (see later...)

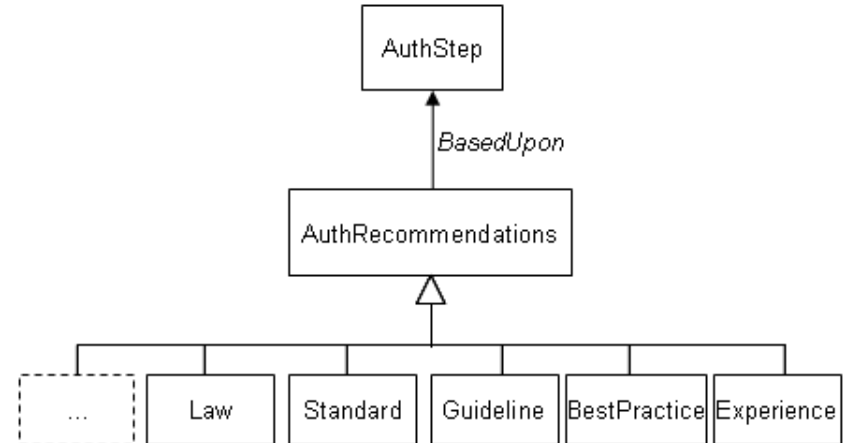
## ■ Authenticity Protocol

- CASPAR Project [www.casparpreserves.eu](http://www.casparpreserves.eu)
- Authenticity Protocol (AP) is a workflow that is applied to a set of digital objects having the same features (e.g. images or documents)



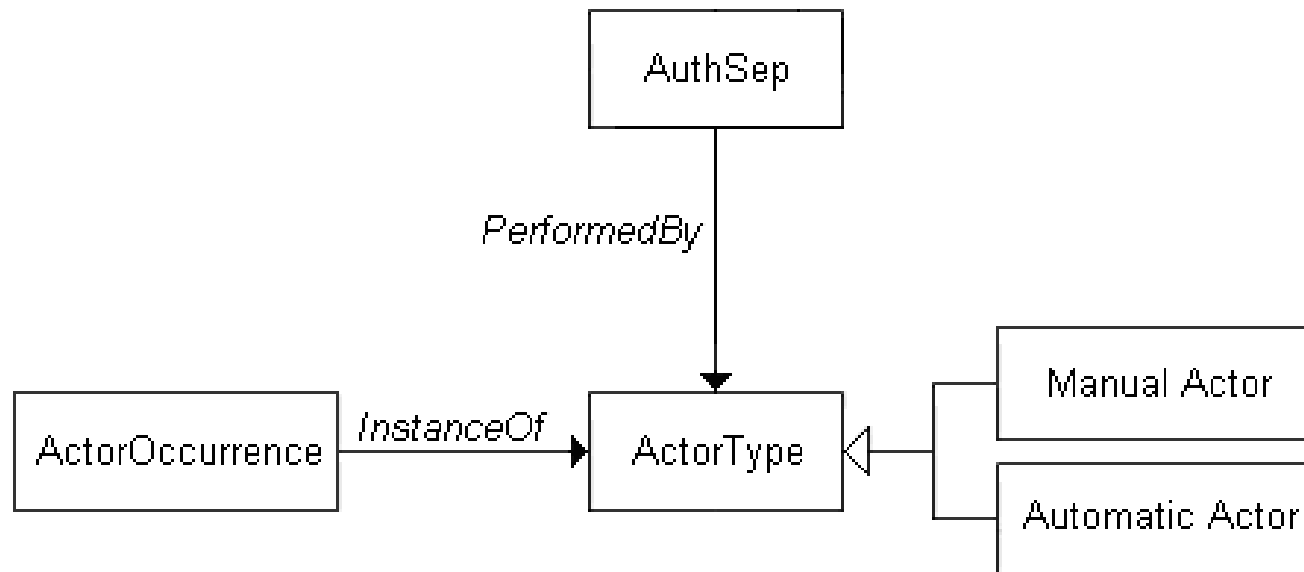
- The process itself consists of different Authenticity Steps (AS) that deal with a certain aspect of a digital object

- Interrelated steps to assess authenticity
- Based on:
  - Reference
  - Provenance
  - Fixity
  - Context
  - Recommendations

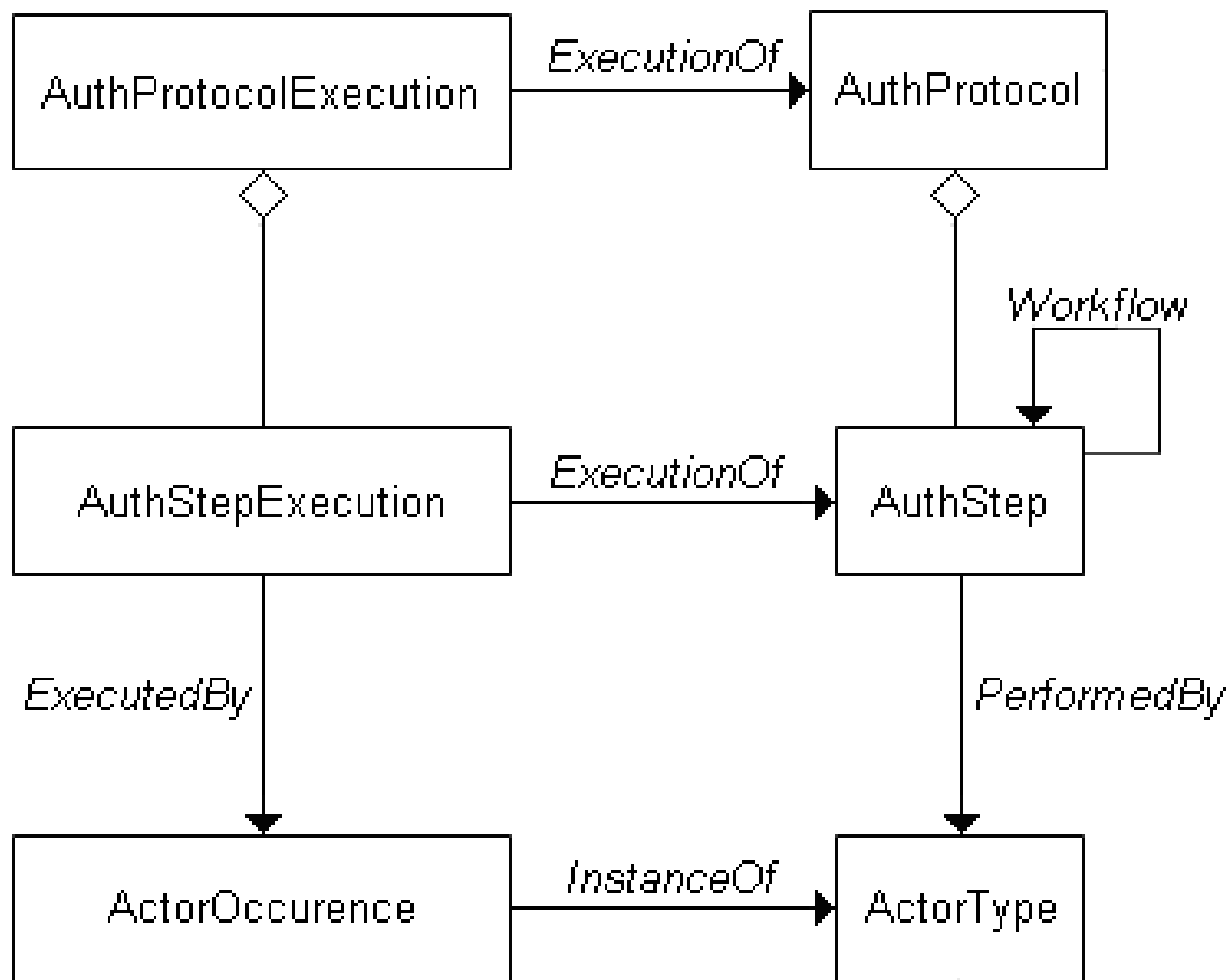


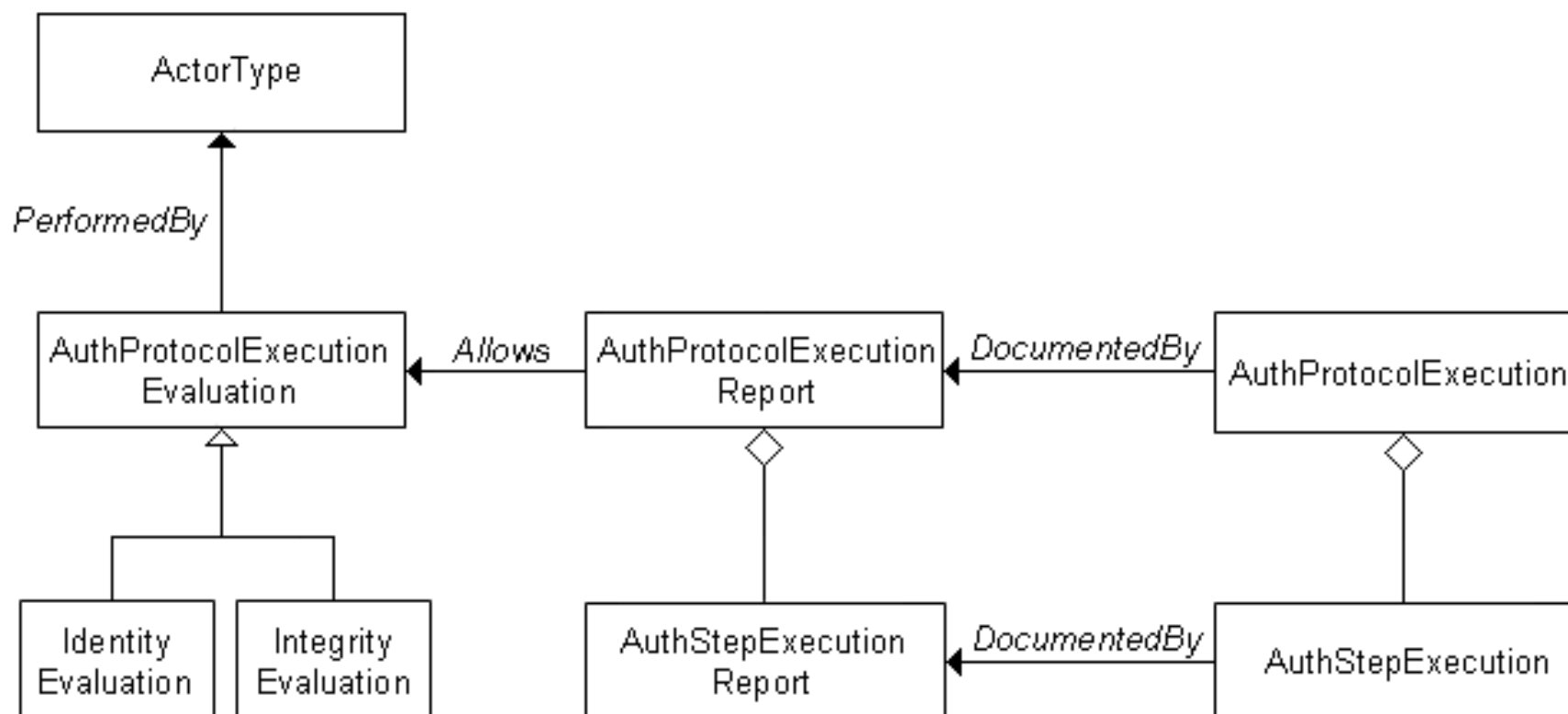
© Mariella Guercio, CASPAR

- Actors
  - Manual actors: human beings
  - Automatic actor: software

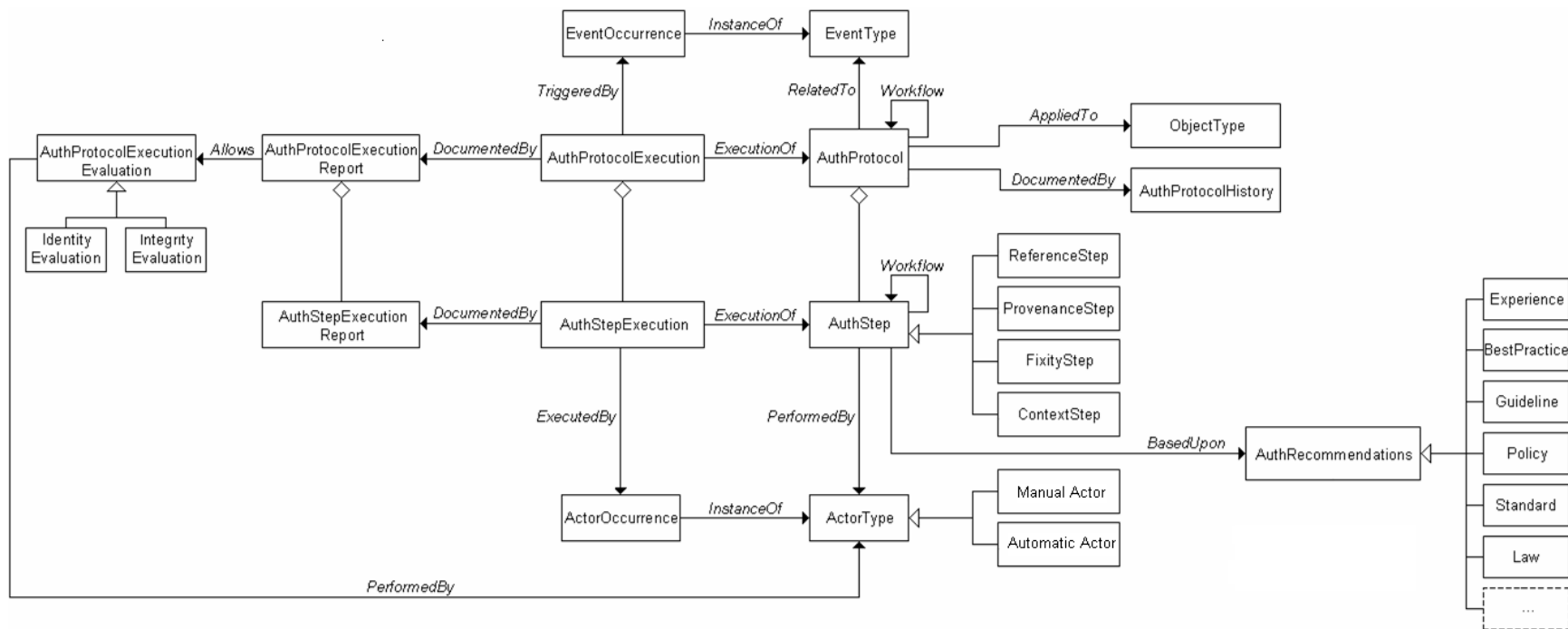


# Authenticity Protocol Execution





# The CASPAR Authenticity Protocol

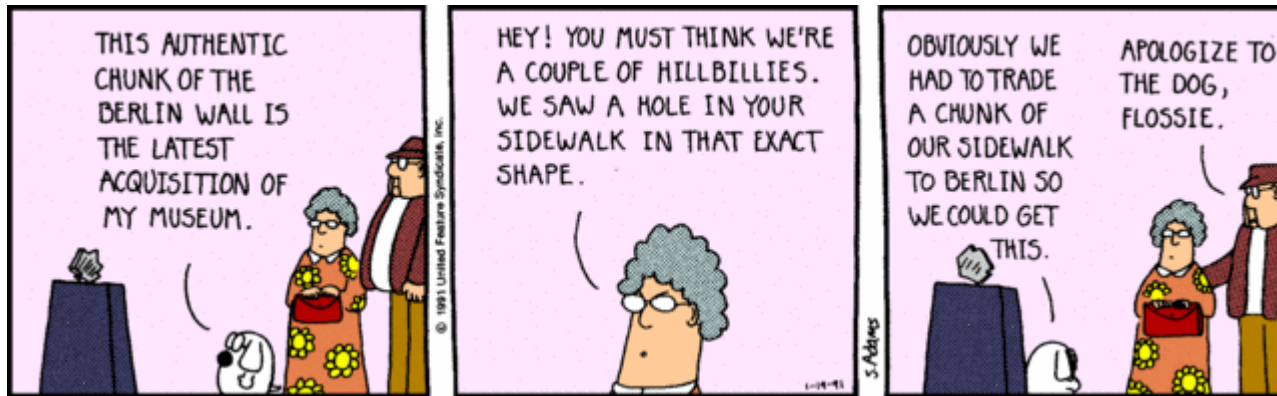


- Each AS should be supported by evidence
  - Can be technical like checksums
  - Non-technical like the reputation of administrators
- Evidence itself needs to be long term compatible
  - You can't assess authenticity if you can't read the evidence
- Comprehensive descriptions how evidence was collected are needed



- Degree to which an object is what it purports to be
  - Needs to be assessed -> Authenticity Protocol
  - Needs to be maintained
- Metadata about evidence:
  - Technical properties
    - Checksums
    - Context
    - ...
  - Social properties
    - Trust
    - Reputation

# Questions?



© Dilbert.com

# What is Provenance?



vergrößern 525x700  
Foto: sotheby's

Mit etwa 80 Millionen Dollar beziffern Sotheby's-Experten ihre Erwartungen für diese Version des "Schrei" von 1895.

SOTHEBY'S

## Munchs "Schrei" soll 80 Millionen Dollar bringen

21. Februar 2012 15:05

**Eine Version des bekanntesten Werks des norwegischen Künstlers gelangt am 2. Mai zur Versteigerung**

Edvard Munchs "Der Schrei" ist das wohl bekannteste Werk des norwegischen Künstlers und gilt als eine der Ikonen der jüngeren Kunstgeschichte. Munch variierte das Motiv zwischen 1893 und 1910 viermal. Drei Versionen befinden sich in Museumssammlungen, die vierte gelangt, wie Sotheby's New York am Dienstag in einer Aussendung bekanntgab, am 2. Mai bei der Impressionist & Modern Art Auktion zur Versteigerung. Es stammt aus dem Besitz des norwegischen Geschäftsmannes Petter Olsen und soll laut Sotheby's-Experten um die 80 Millionen Dollar einspielen.

### MEHR ZUM THEMA

OSLO: Günstig hin & retour: [austrian.com](http://austrian.com)

FRANKFURT: ab 44,99€. Jetzt buchen auf [flyniki.com](http://flyniki.com)

Werbung

Damit gilt das 1895 gemalte Pastell als Anwärter auf einen der höchstdotierten Besitzerwechsel der Auktionsgeschichte (seit Mai 2010: Pablo Picasso, "Nude, Green Leaves and Bust", 106,48 Millionen Dollar, Christie's). Zuletzt hatte die 1893 ausgeführte Variante für Aufsehen gesorgt, als sie zusammen mit einer Madonna im August 2004 aus dem Munch Museum in Oslo gestohlen wurde und erst zwei Jahre später sichergestellt werden konnte.

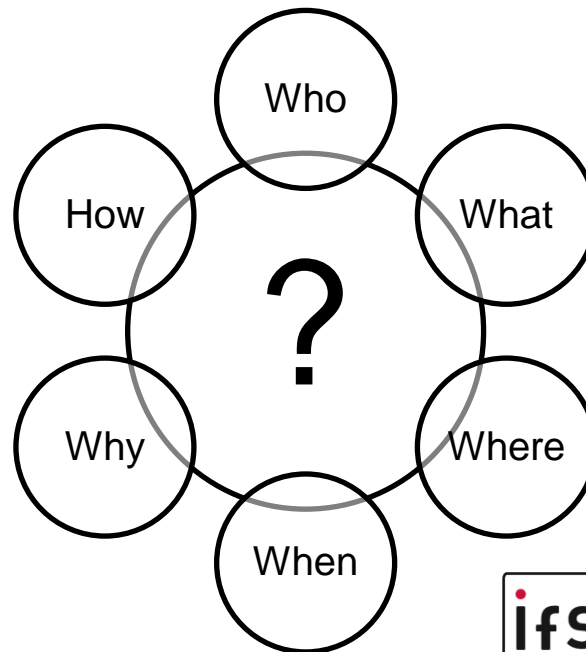
Aktuell widmet die Schirn Kunsthalle in Frankfurt dem Bahnbrecher des Expressionismus mit "Edvard Munch. Der moderne Blick" (bis 13. Mai) eine Ausstellung, die einen neuen Blick auf sein Schaffen bieten will. (kron, [derStandard.at](http://derStandard.at), 21.2.2012)

# What is Provenance?

---

- Provenance of objects describes:
  - Origin
  - Lineage and chain of ownership
  - Chronology of important events
- Relation to Digital Preservation:
  - Provenance documents the history of digital objects
  - Digital archives need to collect and maintain provenance information

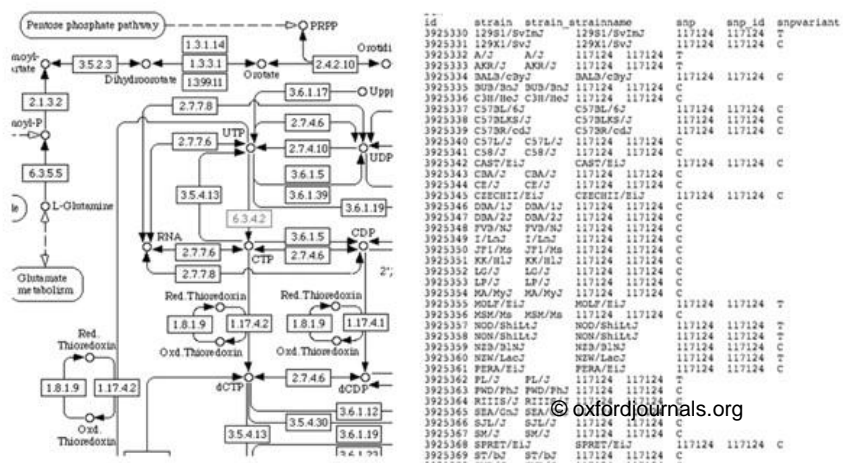
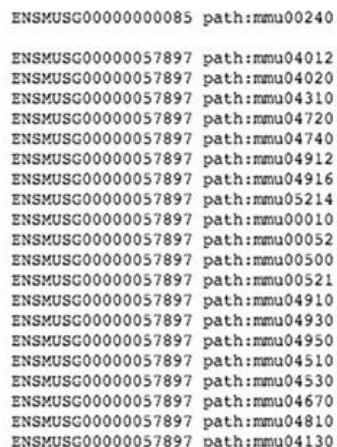
- Metadata describes the complete history of digital information
  - Includes creators, authors, timestamps,....
  - Transactions, modifications, contributions...
- Answers six questions



- There is no “original” of a digital object like in the physical world
- Any copy has the exact same attributes
  - A copy does not destroy the provenance of a digital object
- Digital objects have to be transformed
  - Transformations have to be tracked in the provenance information
  - Provenance is part of the workflow

- Essential for science, governance and commerce
  - Whenever evidence for documents and information is needed
- Many applications
  - Verification of scientific experiments
  - Financial transactions
  - Information flows
  - Drug trials
  - ...
- Provenance is a fundamental principle of archiving
- Secure Provenance: Protecting the Genealogy of Bits, R. Hasan et. al.



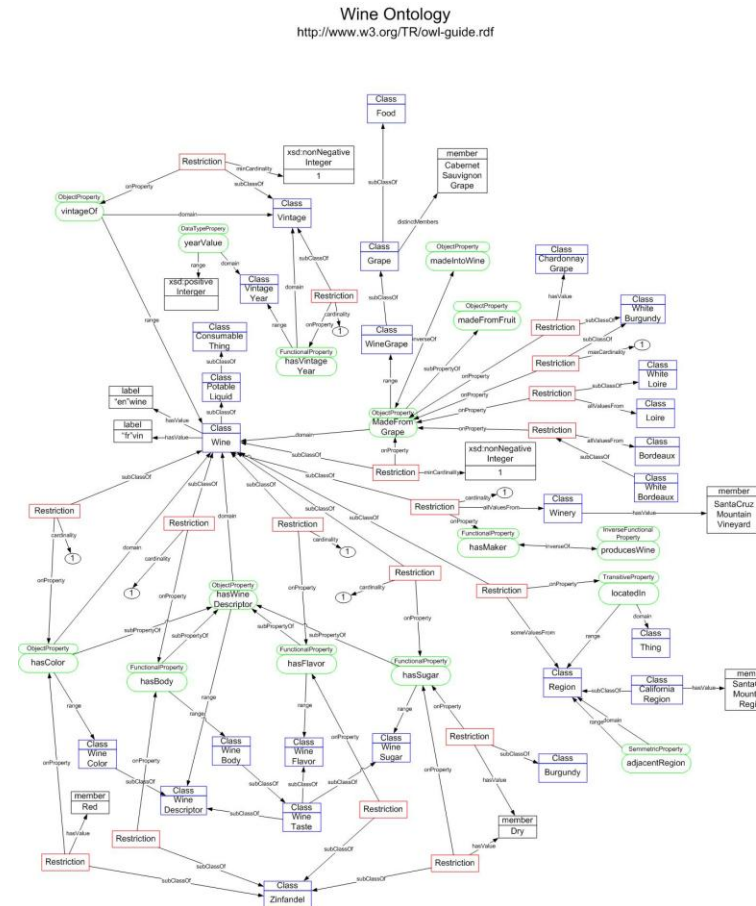




- Complex experiments and business workflows
  - Data sources and data flows have to be traced in order to be used later
  - Monitor data leakage, redundancy and efficiency of information flows
- Benefits of capturing provenance data:
  - Security
    - Detect anomalies
    - Prevent fraud
  - Data quality
    - Assess quality of stored data assets
    - Poor quality will have serious consequences
    - Horror example: U.S. bombing of the Chinese embassy in Belgrade 1999 caused by outdated map (<http://www.defense.gov/transcripts/transcript.aspx?transcriptid=536>)

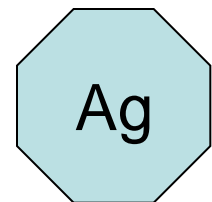
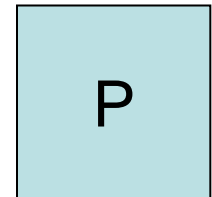
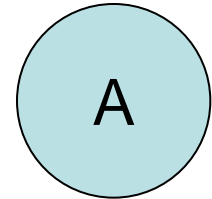
# A Short Excursion: Ontologies

- Formalize Knowledge
- Provide controlled vocabulary
- Specify concepts of a Domain
- Allow to derive knowledge
- Used in Semantic Web ...
- Can be used to describe provenance information

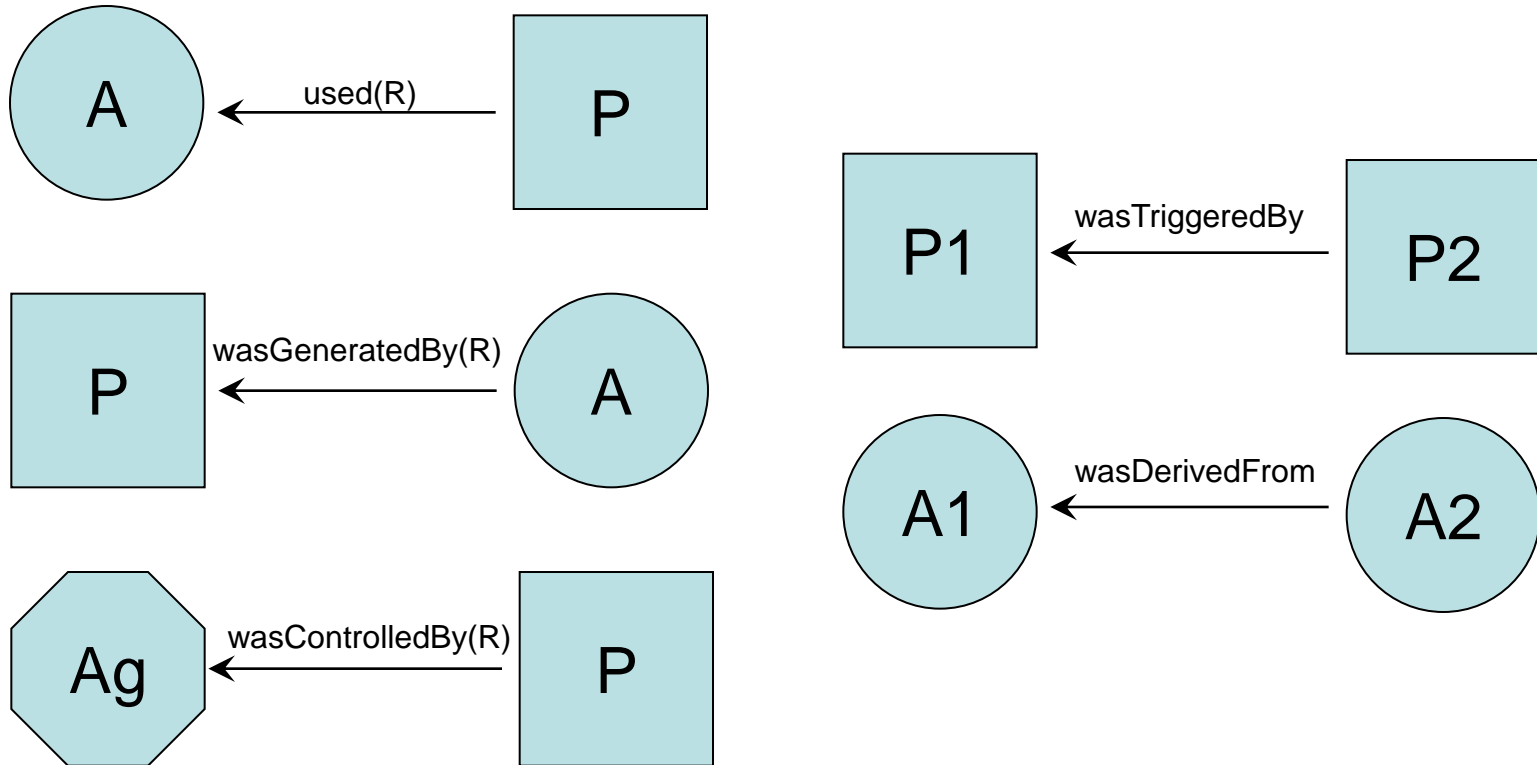


- The Open Provenance Model
  - Can be used for digital and physical objects
  - Is a reference model
  - Defines a core set of inference rules
  - Enables interoperability between different systems
- Goals:
  - Controlled vocabulary
  - Serialization formats
  - APIs

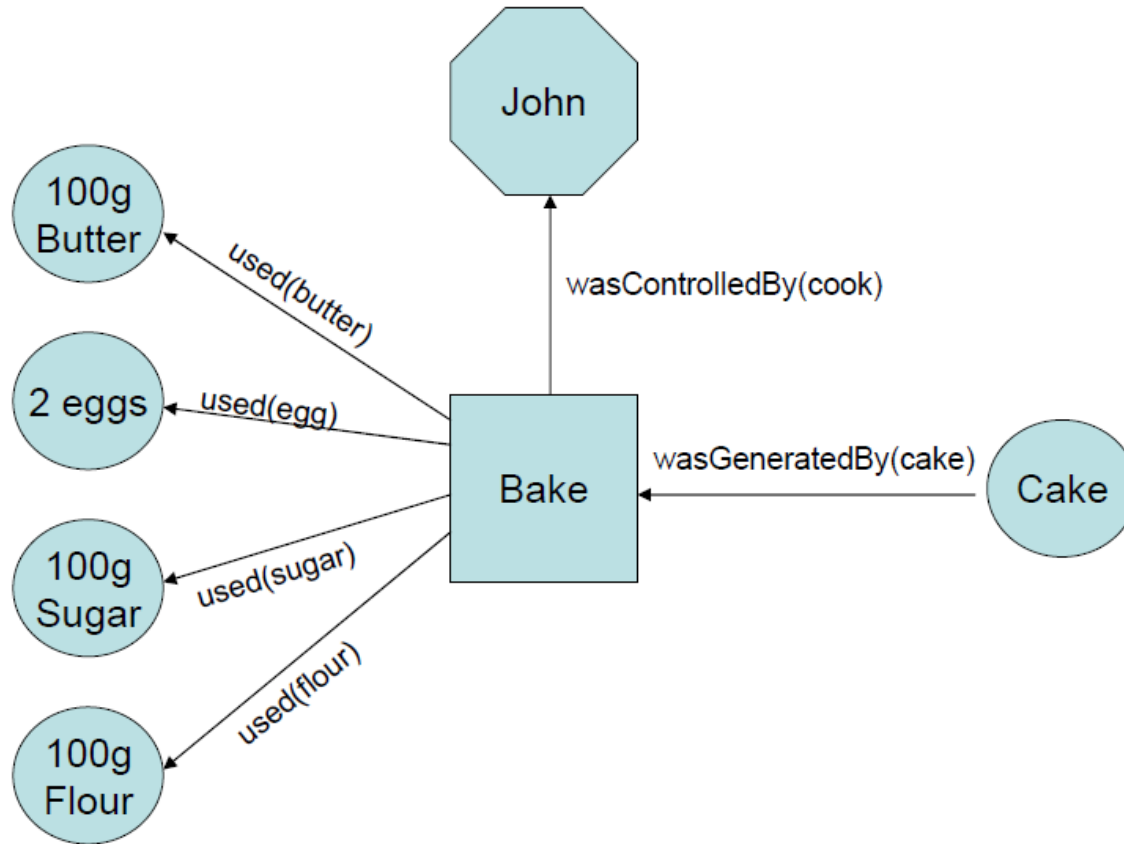
- **Artifact**
  - Object (digital or physical)
  
- **Process**
  - Action or series of actions
  - performed on or caused by artifacts
  - Results in new artifacts
  
- **Agent**
  - Responsible for process execution

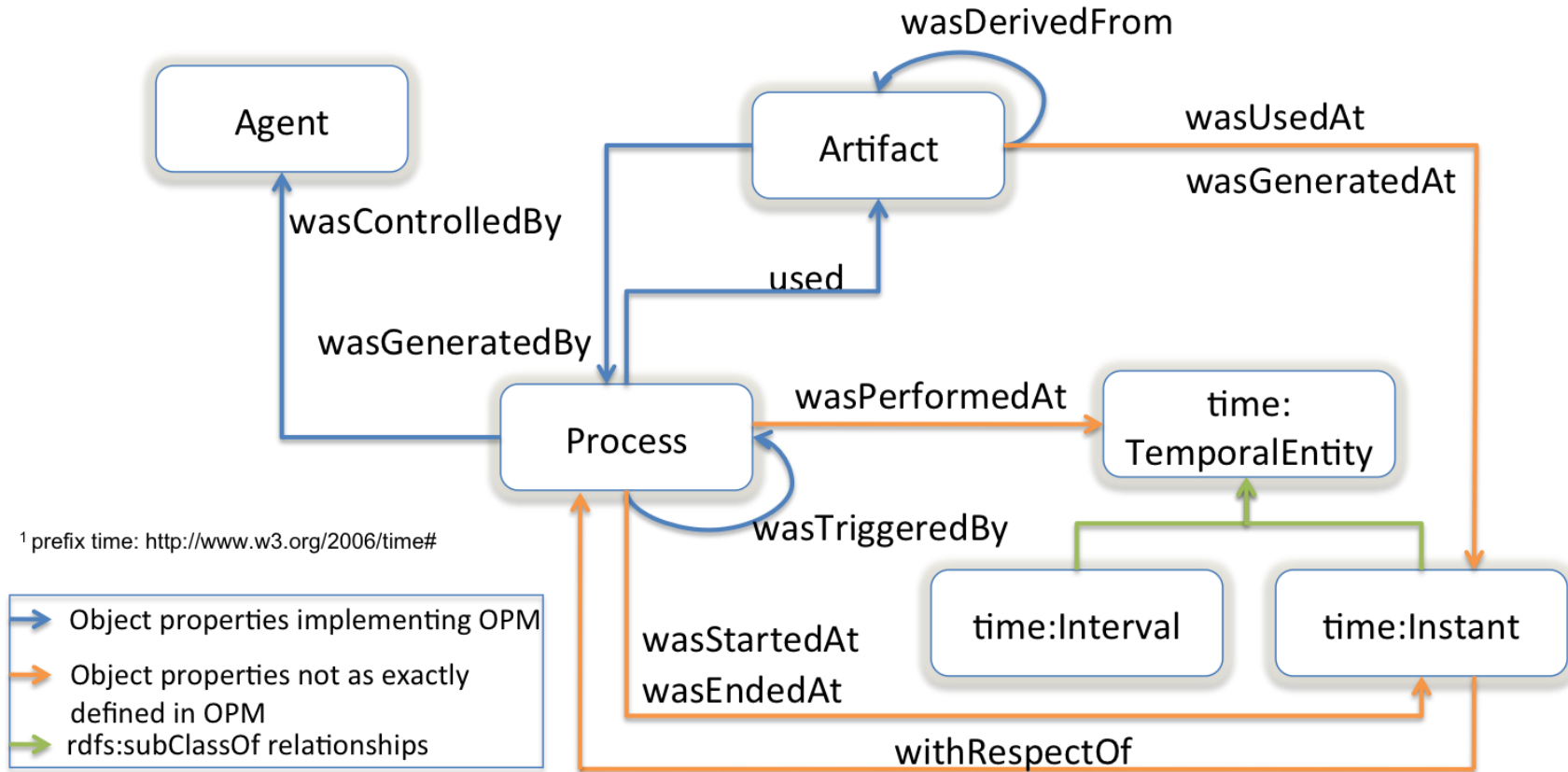


# OPM Edges



# OPM: A Simple Example



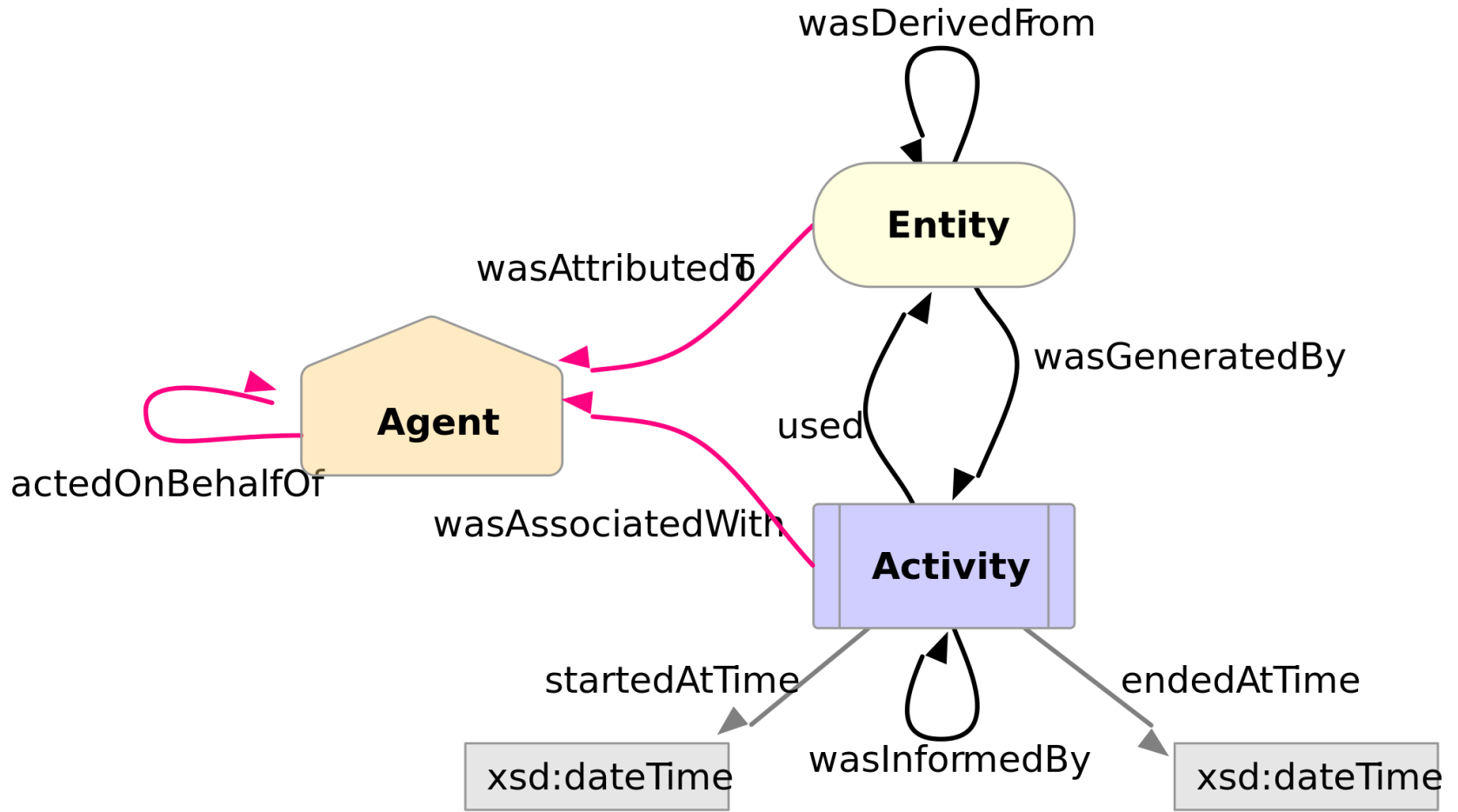


© <http://open-biomed.sourceforge.net/opmv/ns.html>

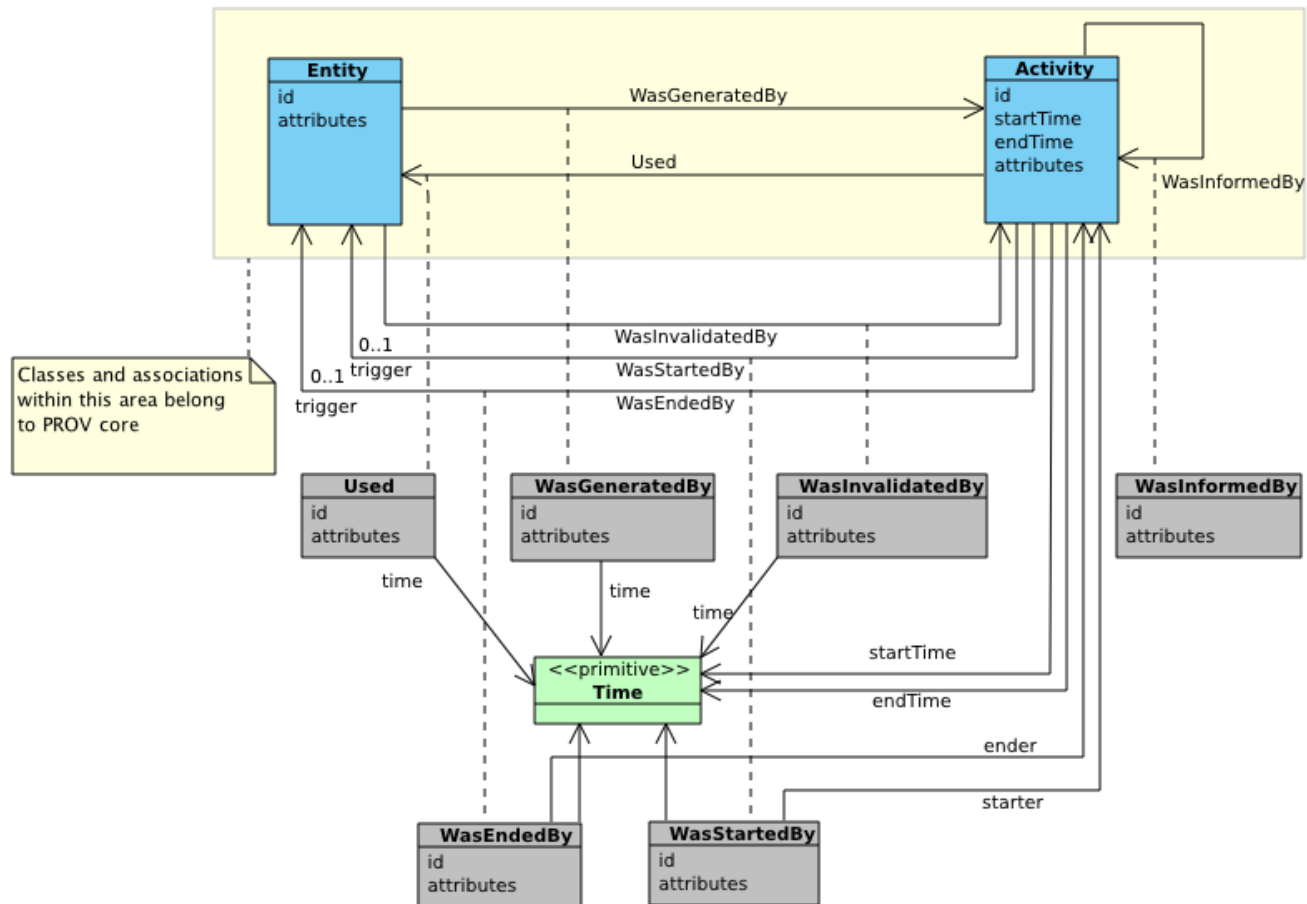
- The current definition of the term provenance by the W3C with reference to data:
  - “Provenance refers to the sources of information, such as entities and processes, involved in producing or delivering an artifact.”
- PROV Ontology (PROV-O) defines the OWL2 Web Ontology Language encoding of the PROV Data Model
- Describes the set of classes, properties, and restrictions that constitute the PROV Ontology.
- The ontology specification provides the foundation to implement provenance applications in different domains that can represent, exchange, and integrate provenance information generated in different systems and under different contexts.



# PROV-O Core Concepts



# PROV-O Example in UML



Source: <http://www.w3.org/TR/prov-dm/#data-model-components>

# PROV-O in Protege

prov-o-20130430 (<http://www.w3.org/ns/prov-o-20130430>) : [/home/stefan/subversion/SBA-repo/DP-Vorlesungsvortrag/DP2014/2014\_Authenticity\_Provenence/prov-o.owl] - + x

File Edit View Reasoner Tools Refactor Window Help

prov-o-20130430 (<http://www.w3.org/ns/prov-o-20130430>) Search for entity

Active Ontology Entities Classes Object Properties Data Properties Annotation Properties Individuals OWLViz DL Query OntoGraf Ontology Differences SPARQL Query

Class hierarchy Class hierarchy (inferred)

Class hierarchy: Entity

- Thing
  - Activity
  - Agent
    - Organization
    - Person
    - SoftwareAgent
  - Entity
    - Bundle
    - Collection
    - EmptyCollection
    - Plan
    - Influence
      - ActivityInfluence
      - Communication
      - Generation
      - Invalidation
    - AgentInfluence
      - Association
      - Attribution

Class Annotations Class Usage

Annotations: Entity

Annotations +

- label Entity
- category starting-point
- component entities-activities
- constraints [type: anyURI] <http://www.w3.org/TR/2013/REC-prov-constraints-20130430/#prov-dm-constraints-fig>
- definition [language: en] An entity is a physical, digital, conceptual, or other kind of thing with some fixed aspects; entities may be real or imaginary.
- dm [type: anyURI] <http://www.w3.org/TR/2013/REC-prov-dm-20130430/#term-entity>

Description: Entity

Equivalent To +

SubClass Of +

SubClass Of (Anonymous Ancestor)

Members +

Target for Key +

Disjoint With +

Individuals by type Annotation property hierarchy Datatypes

Object property hierarchy Data property hierarchy

Object property hierarchy:

topObjectProperty

No Reasoner set. Select a reasoner from the Reasoner menu ☒ Show Inferences

- Provenance data should be generated
  - Automatically (i.e. log files)
  - Machine readable
  - In suitable granularity
  - Securely

# Provenance- Questions?

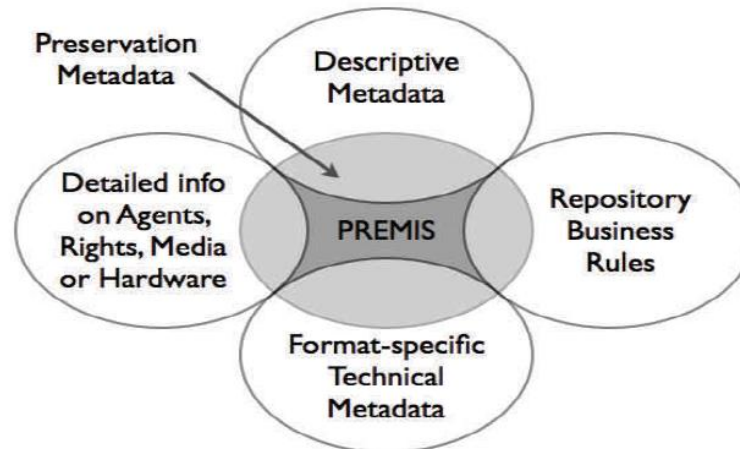
---

- Why is provenance data important?
- What are metadata?
- How can provenance be modelled?
- What are the core concepts of OPM and PROV-O?

- Provenance can serve as evidence for authenticity
  - If the full provenance traces are available, the degree of authenticity is higher
- Authenticity has to be maintained along the object lifecycle
  - Each event that interacts with the object needs to trigger an authenticity protocol execution
  - Provenance metadata keeps track of this events and the actors involved

- Metadata describes data and events in a precise way
  - Needed for authenticity and provenance information
  - Collect all metadata an archive needs for supporting digital preservation processes
- Different Metadata standards exist for various purposes
  - Level of granularity

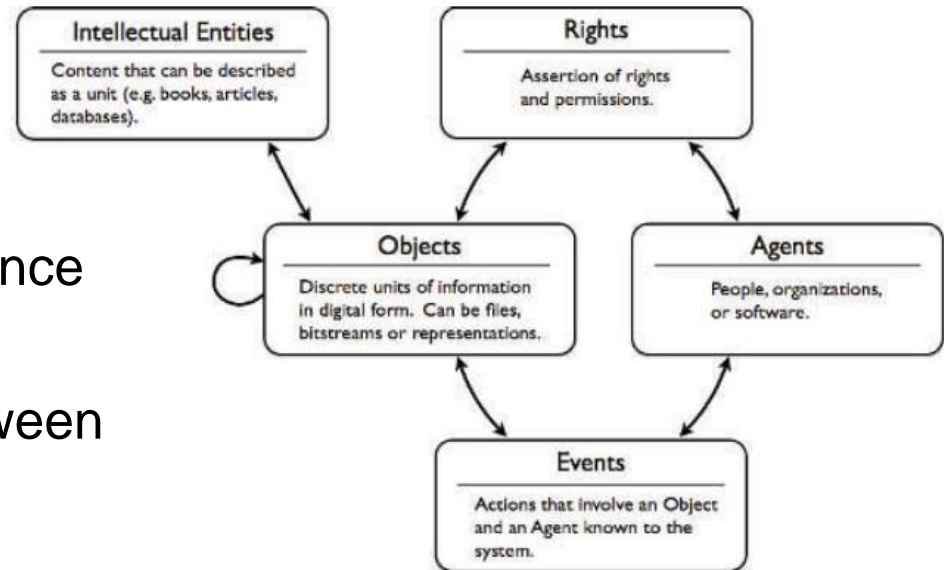
- Preservation Metadata: Implementation Strategies (PREMIS)
  - Data dictionary
  - XML Schema



© <http://www.loc.gov/standards/premis/>



- PREMIS data model
  - Events capture relevant actions and map provenance information.
  - Models relationships between the objects
- Provides OWL ontology
  - Defines clear semantics of the metadata elements
  - Reasoning

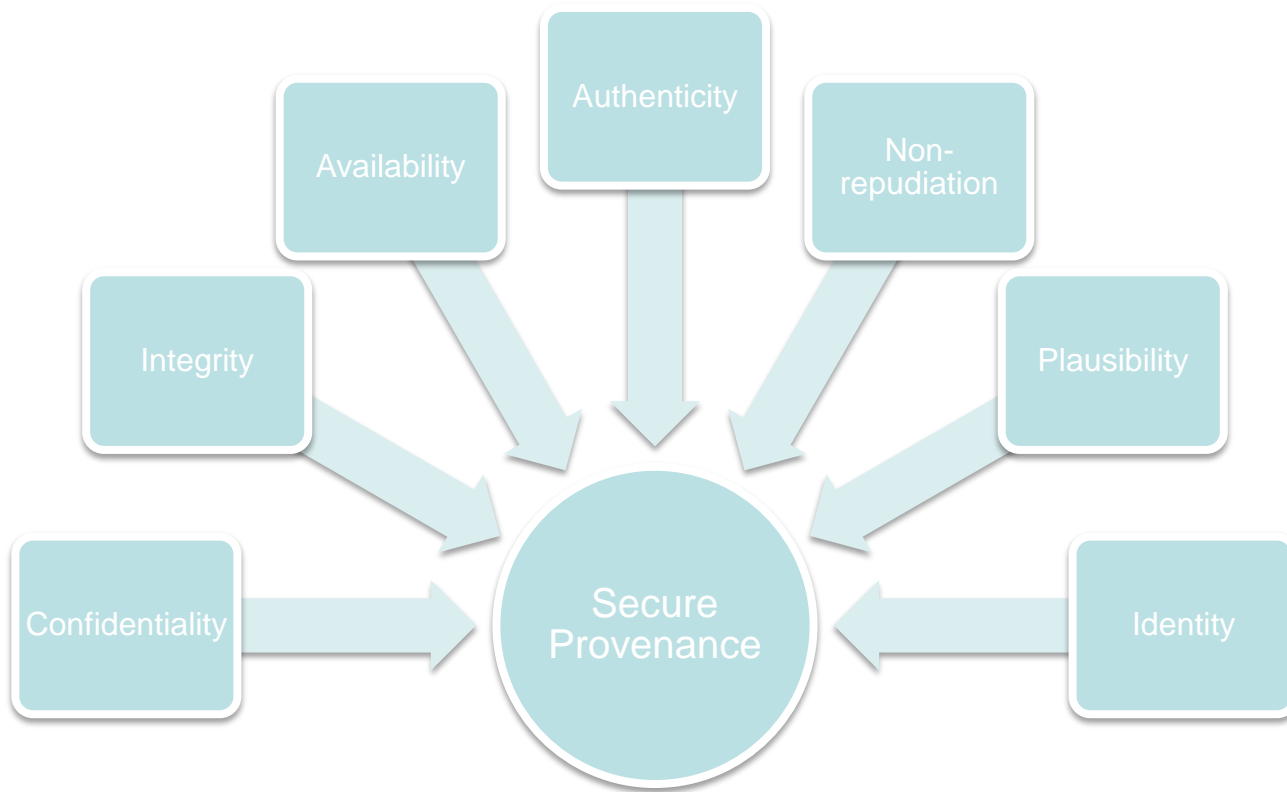


© <http://www.loc.gov/standards/premis/>

- Authenticity and provenance data needs to be protected from manipulation
  - No tampering
  - No insertions, updates or deletes
- Provenance data can be sensitive
  - Privacy considerations
  - Espionage of critical business processes

- Confidentiality
  - Sensitive data must be protected
  - Cryptography
  - Policies
- Integrity
  - Completeness and wholeness in all the significant properties of a digital object
  - Not only on bit-level, but on intellectual form
- Availability
  - Data must be available when they are needed
  - Information must be protected

- **Authenticity**
  - Degree to which digital objects are what they seem to be
- **Non-Repudiation**
  - The participation of an activity can not be denied
  - Events are verifiable
- **Plausibility**
  - Occurring events to not contradict logical assumptions
- **Identity**
  - Uniqueness
  - Distinguishable from other digital objects



- Attackers could try to
  - Delete provenance records
    - Remove incriminating evidence
  - Add fake entries
    - Claim authorship of data
    - „Enhancements“ -> Scientific fraud
  - Manipulate records and alter history
    - Cover tracks
  - Hide contributions
    - Deny responsibility

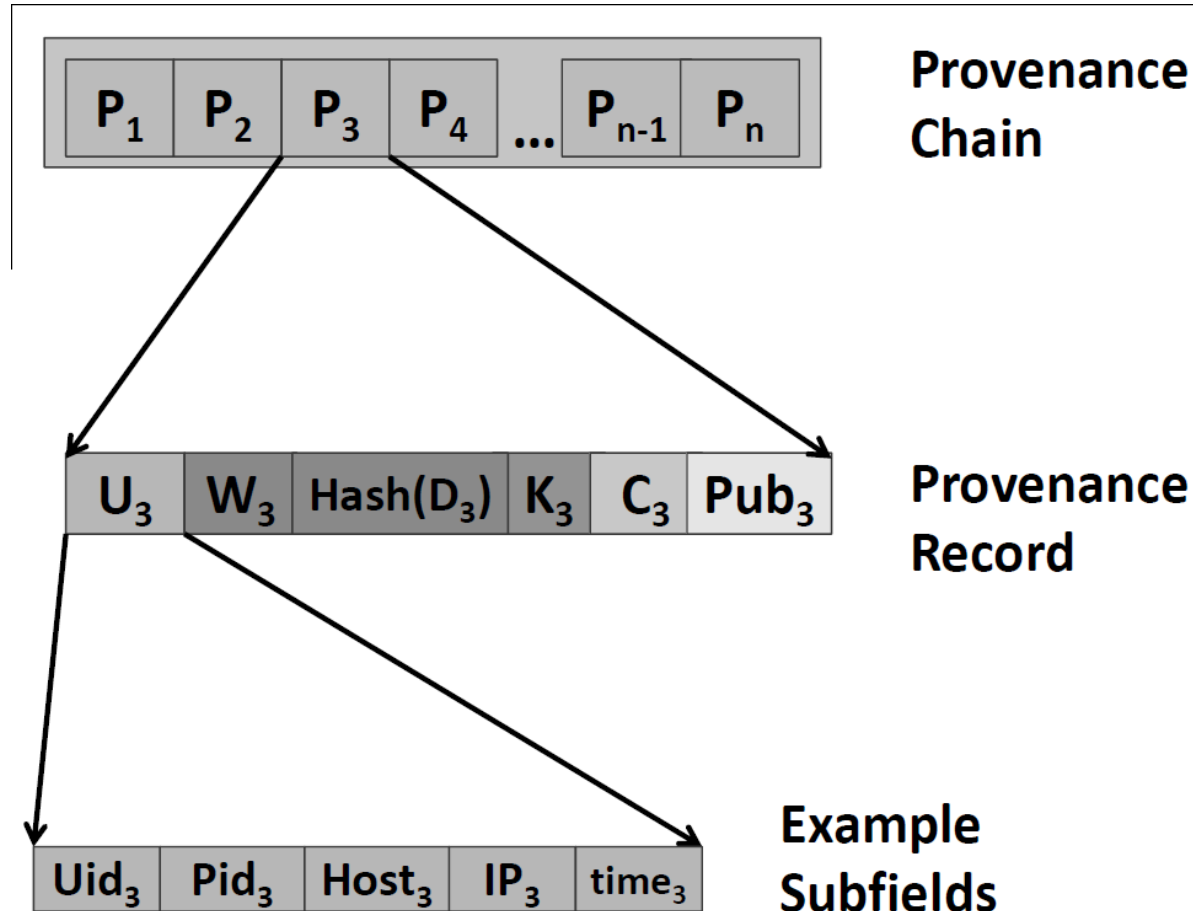
- Logging mechanisms
  - Provenance data can be collected by logs
  - Logs capture relevant events automatically
  - Granularity of logs from high level to system calls
- Log architectures
  - Sender transmits event notification to relay or collector
  - Scalability is important
- Logs have to be auditable
  - Event chain has to be reproducible

- Log file content is highly sensitive
  - Have to be safeguarded against unauthorized access and manipulation
  - Sender and receiver have to agree on a shared cryptographic protocol
  - Signatures ensure integrity
  - Encryption hides content from intruders



- Signature chains
  - Ensure fixity by signing all occurred events
  - Signatures have to be protected from manipulation
    - Only allow new records to the provenance chain to be appended to the end
    - Append-only signatures aggregate signatures from previous records -> no intermediate records can be inserted
    - Forward-secure signatures

# Signature Chains



Hasan et. al. [2]

- Preservation and encryption
  - Cryptographic methods have to be observed for obsolescence
  - What if the keys are lost?
  - XML Advanced Electronic Signatures (XAdES )
- Audits
  - Frequent audits are necessary
  - Auditors are only allowed to read relevant pathways in the graph

- Specialized Provenance Stores
  - Provenance Aware Storage System (PASS)
  - Provenance Data Store (PDS)
- Write Once Read Many
  - Prevents records from being altered or deleted
  - Various systems available

# Security Questions

---

- Why is provenance data a profitable goal for attacks?
- What are the requirements for secure logs?
- What solutions are available?

- CASPAR Project
  - [www.casparpreserves.eu](http://www.casparpreserves.eu)
  - [www.casparpreserves.eu/Members/metaware/Events/training/newsletter-december-2008/training-presentations/michetti-guercio.pdf](http://www.casparpreserves.eu/Members/metaware/Events/training/newsletter-december-2008/training-presentations/michetti-guercio.pdf)
- Secure Provenance: Protecting the Genealogy of Bits, R. Hasan et. al.
- The OPM Provenance Model
  - [www.openprovenance.org](http://www.openprovenance.org)
  - <http://twiki.ipaw.info/bin/view/OPM/>
- Data provenance – the foundation of data quality. P. Buneman, S. Davidson, University of Edinburgh. Edinburgh, UK
- PROV-O:
  - <http://www.w3.org/TR/prov-o/>
  - <http://www.edbt.org/Proceedings/2013-Genova/papers/edbt/a80-missier.pdf>

Thank you for your attention.

[sproell@sba-research.org](mailto:sproell@sba-research.org)