# Digital Preservation

# Authenticity and Provenance

Stefan Pröll

# Outline

- **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

# The Concept of Trust

- 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?

# Authenticity

- What is Authenticity?


http://weird.cz/ali-g

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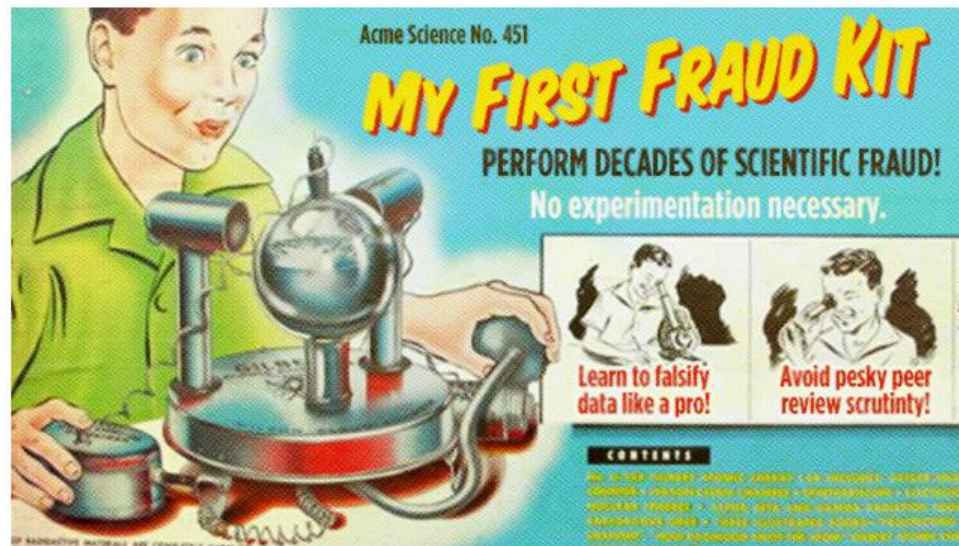- Authenticity is judged on the basis of <u>evidence</u>

# Epic Fraud



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

# Judging Authenticity

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



© www.artmarketblog.com
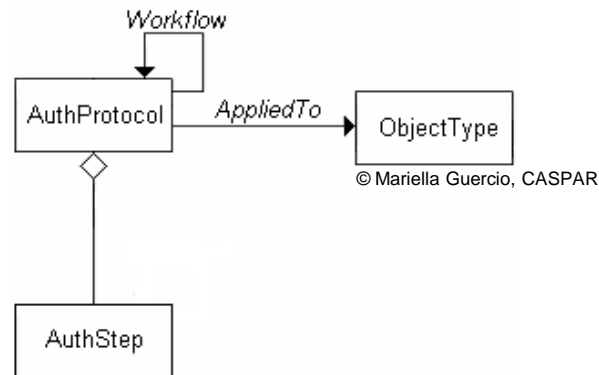
# But still...

# Judging Authenticity of Digital Data

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

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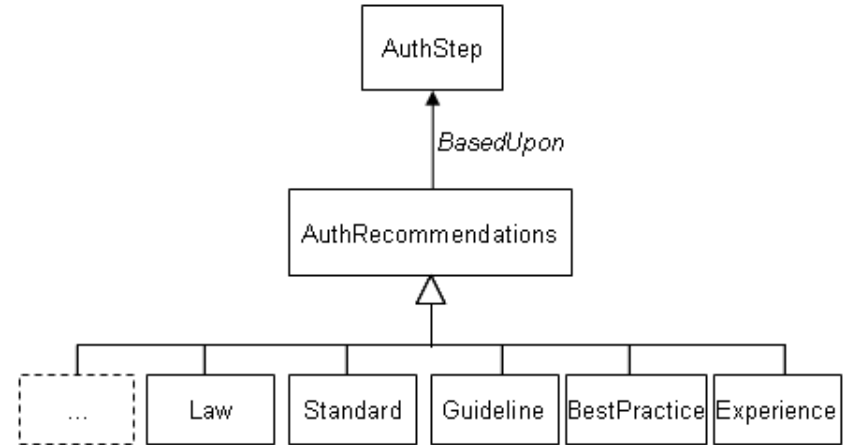
- Authenticity Protocol
  - CASPAR Project [www.casparpreserves.eu](www.casparpreserves.eu)
  - Authenticity Protocol (AP) is a <u>workflow</u> that is applied to a set of digital objects having the same features (e.g. images or documents)



© Mariella Guercio, CASPAR

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

# Authenticity Steps

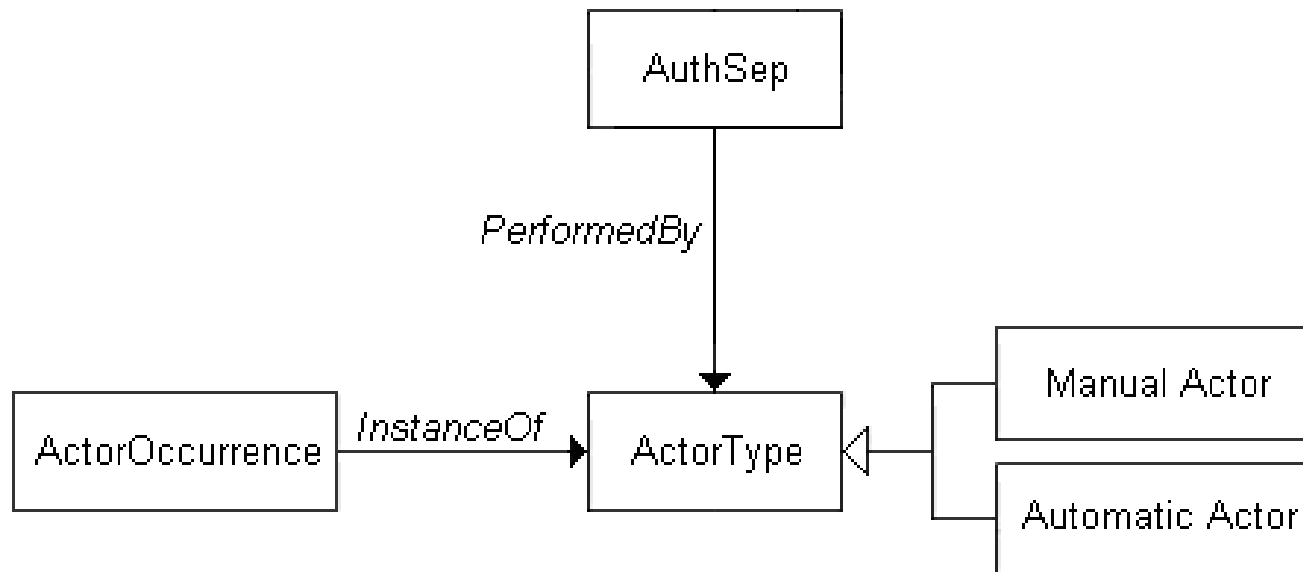- Interrelated steps to assess authenticity
- Based on:
  - Reference
  - Provenance
  - Fixity
  - Context
  - Recommendations
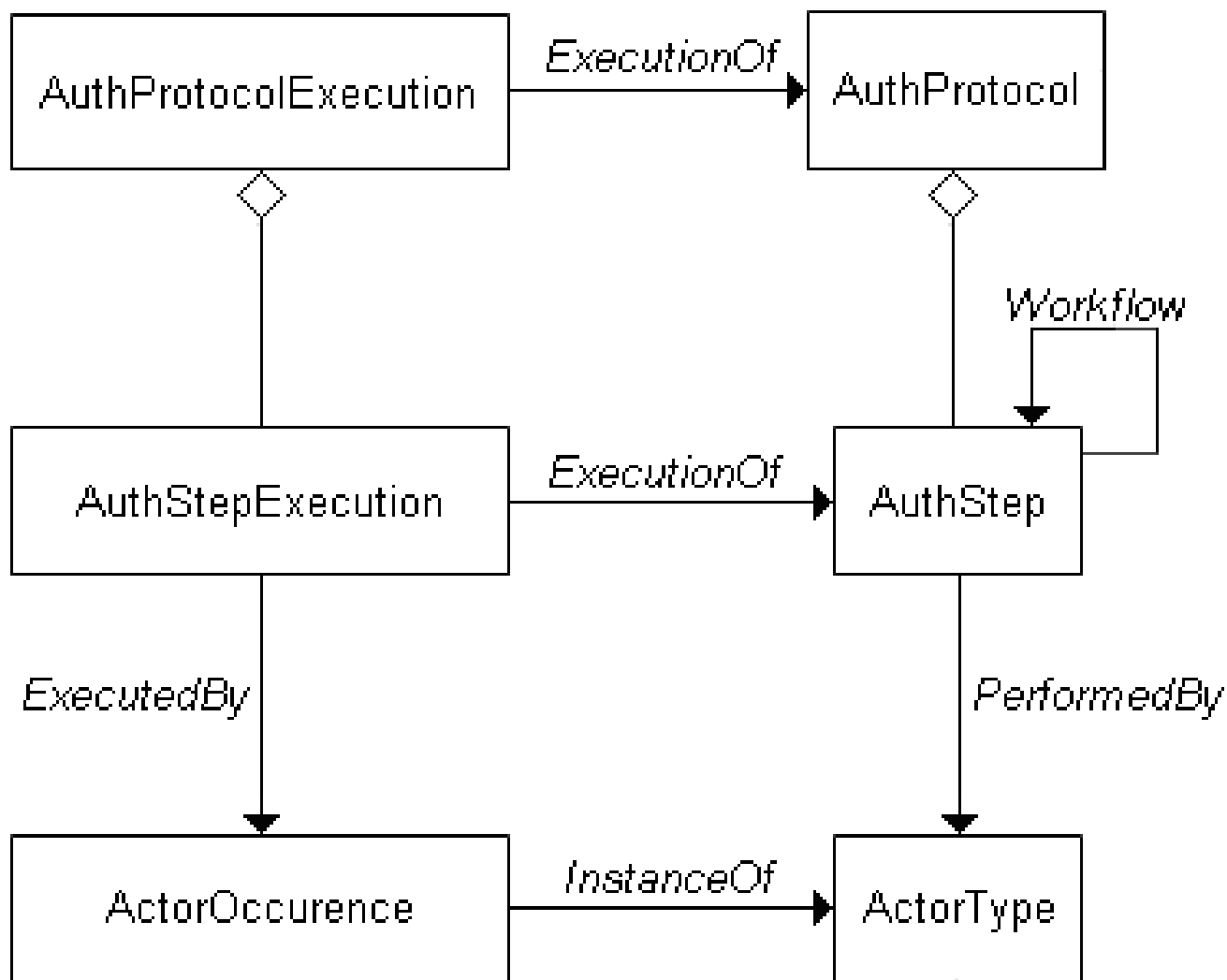


© Mariella Guercio, CASPAR

© Mariella Guercio, CASPAR

# Authenticity Protocol Actors

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

# Authenticity Protocol Execution

# Authenticity Report and Evaluation

# The CASPAR Authenticity Protocol

# Evidence

- 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

# Summary: Authenticity

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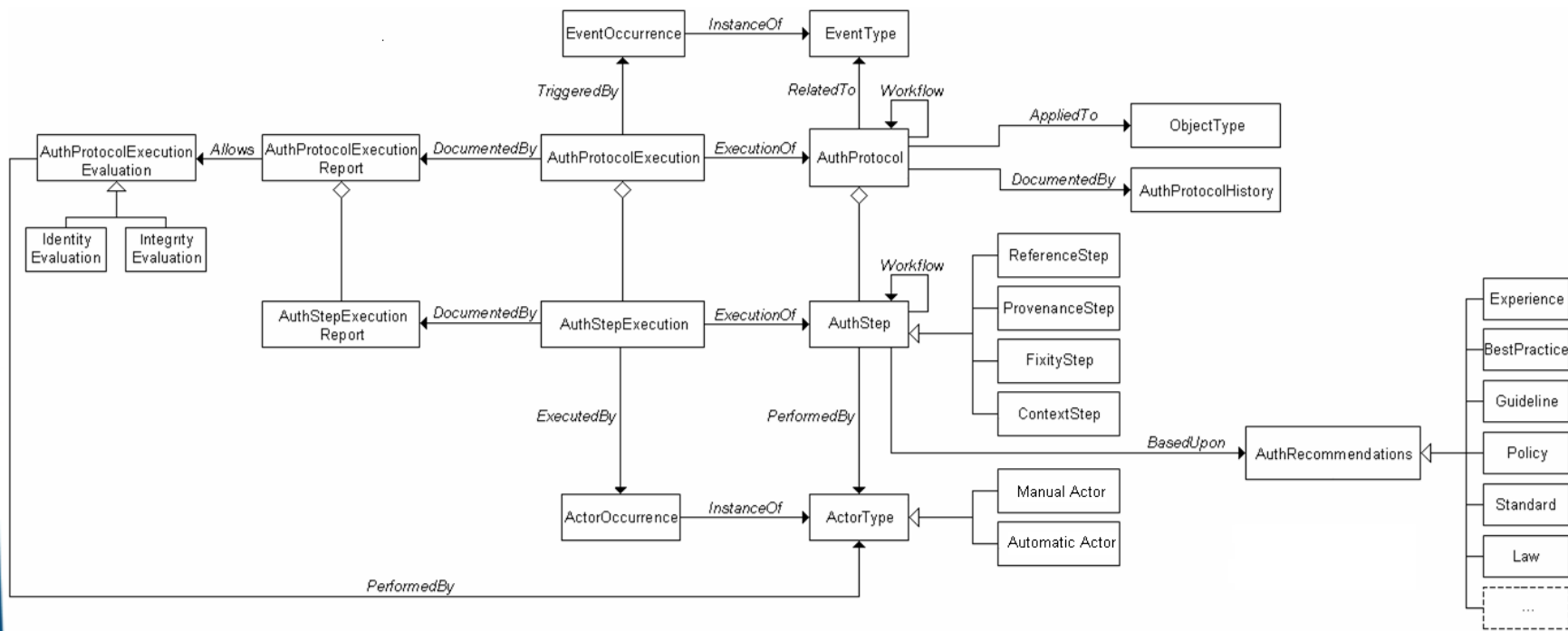  - Needs to be assessed -> Authenticity Protocol
  - Needs to be maintained
- Metadata about evidence:
  - Technical properties
    - Checksums
    - Context
    - …
  - Social properties
    - Trust
    - Reputation

# Questions?



© Dilbert.com

# What is Provenance?



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
FRANKFURT: ab 44,99€. Jetzt buchen auf 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.
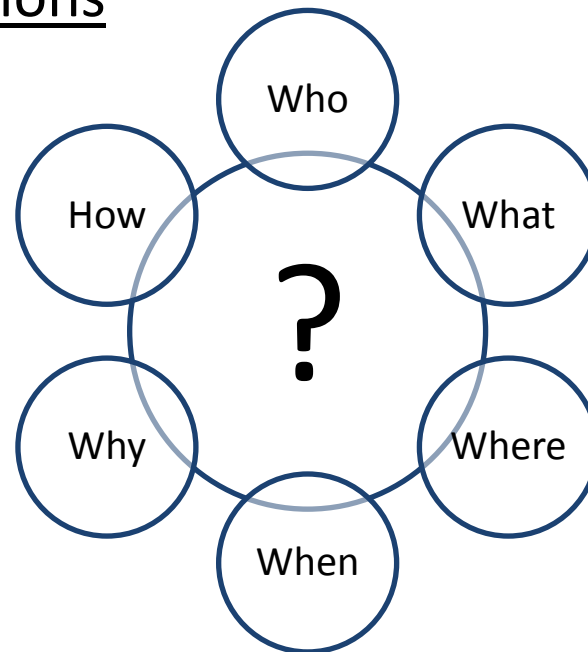
Aktuell widmet die Schirn Kunsthalle in Frankfurt dem Bahnbrecher des Expressionismus mit "Edvard Munch. Der moderne Blick" (bis 13. Mai) eine Ausstellung, die einen neuen Blick auf sein Schaffen bieten will. (kron, derStandard.at, 21.2.2012)

vergrößern 525x700
Foto: sotheby's

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

# 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

# Provenance

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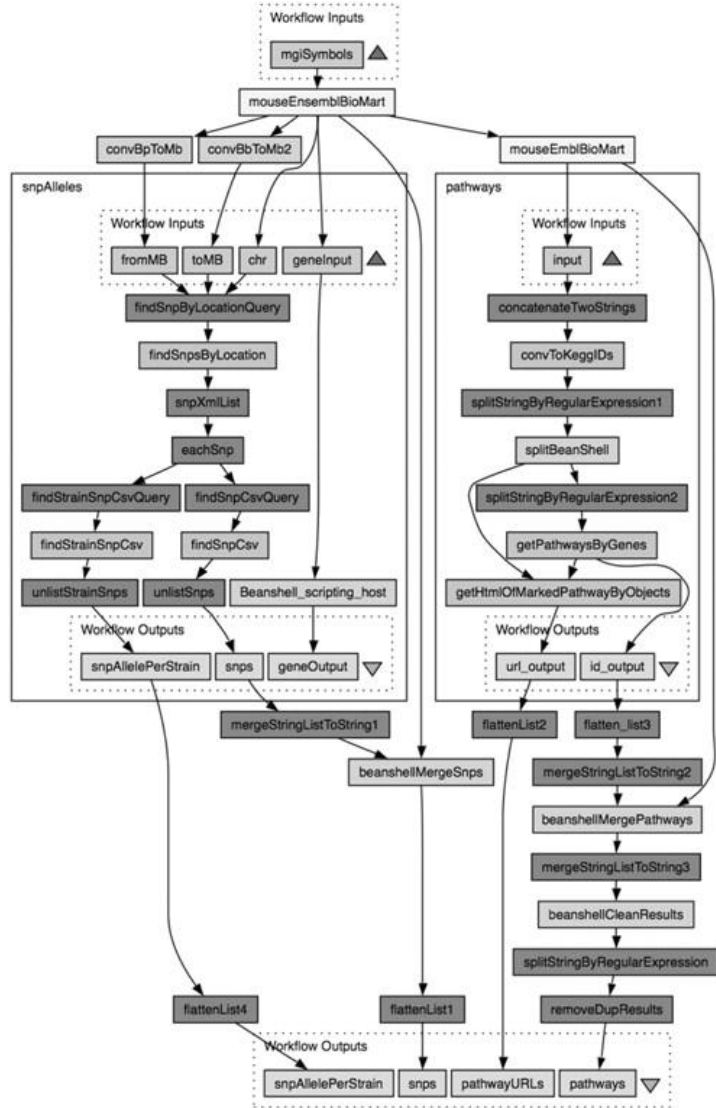- <u>Answers six questions</u>

# Provenance Challenges

- 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

# Data Provenance

- 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 <u>fundamental principle of archiving</u>
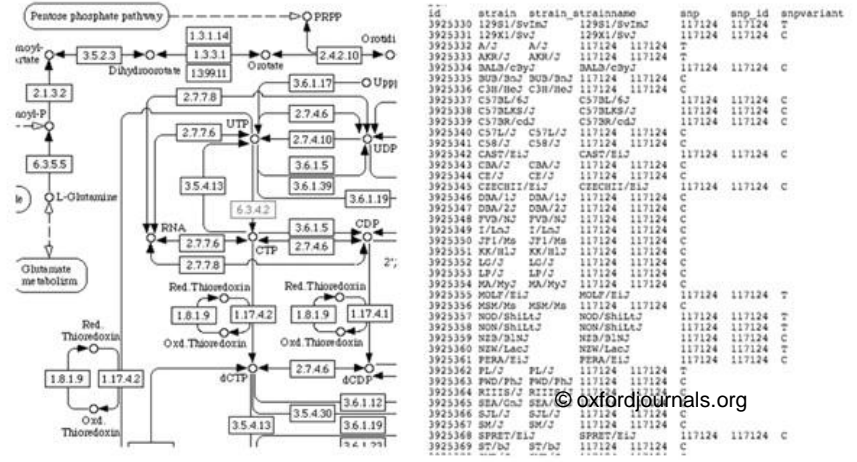- Secure Provenance: Protecting the Genealogy of Bits, R. Hasan et. al.

# Example

© oxfordjournals.org

# Capturing Provenance

- Complex experiments and business workflows
  - Data sources and data flows have to be <u>traced</u> 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)

# The Open Provenance Model

- 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

# OPM Nodes

- 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

A

P

Ag

# OPM Edges

# OPM: A Simple Example

# The Open Provenance Model Extensibility



¹ prefix time: http://www.w3.org/2006/time#

Object properties implementing OPM
Object properties not as exactly defined in OPM
rdfs:subClassOf relationships

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

# Generating Provenance Data

- 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?

# Relation of Authenticity and Provenance

- 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

# Preservation Metadata

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

# PREMIS

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



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

# 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/

# Security Considerations

- 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

# Security Requirements

- 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

# Security Requirements

- 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

# Information Security Aspects of Authenticity and Provenance

# Threat Model for Provenance

- 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

# Secure Provenance

- 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

# Secure Logging

- 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

# Append-Only Signatures

- 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 appended to the end
    - Append-only signatures aggregate signatures from previous records -> no intermediate records can be inserted
    - Forward-secure signatures

# Signature Chains



Provenance Chain: $P_1$ $P_2$ $P_3$ $P_4$ ... $P_{n-1}$ $P_n$

Provenance Record: $U_3$ $W_3$ $Hash(D_3)$ $K_3$ $C_3$ $Pub_3$

Example Subfields: $Uid_3$ $Pid_3$ $Host_3$ $IP_3$ $time_3$

Hasan et. al. [2]

# Secure Logging - Challenges

- Preservation and encryption
  - Cryptographic methods have to be observed for obsoletion
  - 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

# Secure Storage

- 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?

# Literature and Links

- CASPAR Project
  - www.casparpreserves.eu
  - 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://twiki.ipaw.info/bin/view/OPM/
- Data provenance – the foundation of data quality. P. Buneman, S. Davidson, University of Edinburgh. Edinburgh, UK

# Thank you for your attention.

sproell@sba-research.org