# Approaches to
# Making Dynamic Data Citeable
## Recommendations of the RDA Working Group

## Andreas Rauber

Vienna University of Technology
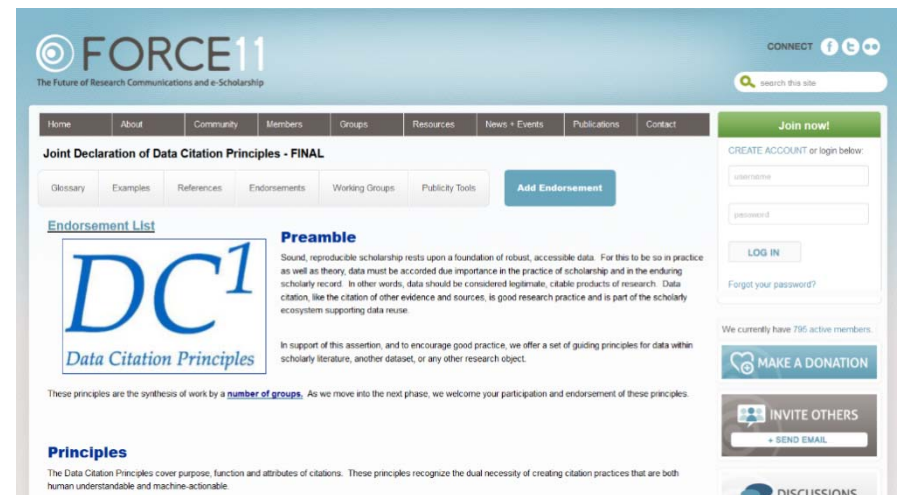rauber@ifs.tuwien.ac.at
http://www.ifs.tuwien.ac.at/~andi

# Outline

- Joint Declaration of Data Citation Principles

- Challenges in non-trivial settings

- Recommendation of the RDA Working Group

- Pilots

- Summary

# Joint Declaration of Data Citation Principles

- 8 Principles created by the Data Citation Synthesis Group

- https://www.force11.org/datacitation

- The Data Citation Principles cover purpose, function and attributes of citations

- Goal: Encourage communities to develop practices and tools that embody uniform data citation principles

FACULTY OF !NFORMATICS

1) **Importance**

   Data should be considered legitimate, citable products of research. Data citations should be accorded the <u>same importance as publications</u>.

2) **Credit and Attribution**

   Data citations should facilitate giving credit <u>and normative and legal attribution</u> to all contributors to the data.

FACULTY OF !NFORMATICS

**3) Evidence**

Whenever and wherever a <u>claim relies upon data</u>, the corresponding data should be cited.

**4) Unique Identification**

A data citation should include a <u>persistent method</u> for identification that is <u>machine actionable</u>, globally unique, and widely used by a community.

FACULTY OF !NFORMATICS

# Joint Declaration of Data Citation Principles (cont'd)

**5) Access**

Data citations should facilitate access to the data themselves and to such associated metadata, documentation, code, and other materials, as are necessary for both <u>humans and machines</u> to make <u>informed use</u> of the referenced data.

**6) Persistence**

Unique identifiers, and metadata describing the data, and its disposition, should persist - even <u>beyond the lifespan</u> of the data they describe.

FACULTY OF !NFORMATICS

## 7) Specificity and Verifiability

Data citations should facilitate identification of, access to, and verfication of the specific data that support a claim. Citations or citation metadata should include information about <u>provenance and fixity</u> sufficient to facilitate verfiying that the specific <u>timeslice, version and/or granular portion</u> of data retrieved subsequently is the same as was originally cited.

FACULTY OF !NFORMATICS

8) **Interoperability and flexibility**

Data citation methods should be sufficiently flexible to accommodate the <u>variant practices among communities</u>, but should not differ so much that they compromise interoperability of data citation practices across communities.

FACULTY OF !NFORMATICS

# Outline

- Joint Declaration of Data Citation Principles

- Challenges in non-trivial settings

- Recommendation of the RDA Working Group

- Pilots

- Summary

FACULTY OF !NFORMATICS

# Data Citation

- Citing data may seem easy
  - from providing a URL in a footnote
  - via providing a reference in the bibliography section
  - to assigning a PID (DOI, ARK, …) to dataset in a repository
- What's the problem?

FACULTY OF !NFORMATICS

# Citation of Dynamic Data

- Citable datasets have to be static
  - Fixed set of data, no changes:
    no corrections to errors, no new data being added
- But: (research) data is **dynamic**
  - Adding new data, correcting errors, enhancing data quality, …
  - Changes sometimes highly dynamic, at irregular intervals
- Current approaches
  - Identifying entire data stream, without any versioning
  - Using "accessed at" date
  - "Artificial" versioning by identifying batches of data (e.g. annual), aggregating changes into releases (time-delayed!)
- Would like to cite precisely the **data as it existed at certain point in time**, without delaying release of new data

FACULTY OF !NFORMATICS

# Granularity of Data Citation

- What about the **granularity** of data to be cited?
  - Databases collect enormous amounts of data over time
  - Researchers use specific subsets of data
  - Need to identify precisely the subset used
- Current approaches
  - Storing a copy of subset as used in study -> scalability
  - Citing entire dataset, providing textual description of subset -> imprecise (ambiguity)
  - Storing list of record identifiers in subset -> scalability, not for arbitrary subsets (e.g. when not entire record selected)
- Would like to be able to cite precisely the **subset of (dynamic) data used** in a study

FACULTY OF !NFORMATICS

# Data Citation – Requirements

- Dynamic data
  - corrections, additions, …
- Arbitrary subsets of data (granularity)
  - rows/columns, time sequences, …
  - from single number to the entire set
- Stable across technology changes
  - e.g. migration to new database
- Machine-actionable
  - not just machine-readable,
    definitely not just human-readable and interpretable
- Scalable to very large / highly dynamic datasets
  - but should also work for small and/or static datasets

# RDA WG Data Citation

- Research Data Alliance
- WG on **Data Citation:
  Making Dynamic Data Citeable**
- WG officially endorsed in March 2014
    - Concentrating on the problems of
      **large, dynamic (changing) datasets**
    - Focus!
      Not: PID systems, metadata, citation string, attribution, …
    - Liaise with other WGs and initiatives on data citation
      (CODATA, DataCite, Force11, …)

    - https://rd-alliance.org/working-groups/data-citation-wg.html

FACULTY OF !NFORMATICS

# Outline

- Joint Declaration of Data Citation Principles

- Challenges in non-trivial settings

- **Recommendation of the RDA Working Group**

- Pilots

- Summary

# Making Dynamic Data Citeable

## Data Citation: Data + Means-of-access

- Data → time-stamped & versioned (aka history)

Researcher creates working-set via some interface:

- Access → **assign PID to QUERY**, enhanced with
  - **Time-stamping** for re-execution against versioned DB
  - **Re-writing** for normalization, unique-sort, mapping to history
  - **Hashing** result-set: verifying identity/correctness

  leading to landing page

S. Pröll, A. Rauber. **Scalable Data Citation in Dynamic Large Databases: Model and Reference Implementation.** In IEEE Intl. Conf. on Big Data 2013 (IEEE BigData2013), 2013
http://www.ifs.tuwien.ac.at/~andi/publications/pdf/pro_ieeebigdata13.pdf

FACULTY OF !NFORMATICS

# Data Citation – Deployment

- ██████████████████ ████ subset of data
- Upon executing selection („download") user gets
  - Data (pac...
  - PID (e.g. ...
  - Hash val...
  - Recommended citation text (e.g. bTeX)

- PID resolves to landing page
  - Provides detailed metadata, link parent data set, subset,…
  - Option to retrieve **original data** OR **current version** OR **changes**

- Upon activating PID associated with a data citation
  - Query is re-executed against time-stamped and versioned DB
  - Results as above are returned

**Note: query string provides excellent provenance information on the data set!**

**This is an important advantage over traditional approaches relying on, e.g. storing a list of identifiers/DB dump!!!**

# Data Citation

**System set-up to support dynamic data:**

1. **Ensure data is time-stamped**
   i.e. any additions, deletions are marked with a timestamp
   (*optional, if data is dynamic*)

2. **Ensure data is versioned**
   i.e. updates not implemented as overwriting an earlier value, but as *marked-as-deleted* and *re-inserted* with new value, both time-stamped
   (*optional, if data is dynamic and access to previous versions is desired*)

3. **Create a query store for queries and metadata**

FACULTY OF !NFORMATICS

# Data Citation – Recommendations
## *(Draft, 2/4)*

**When a specific subset of data needs to be persistently identified** (i.e. not necessarily for all subsets!):

1. **Re-write the query to a normalized form** (*optional*)
2. **Specifically: re-write the query to ensure unique sort of the result set** (*optional*)
3. **Compute a hash key** of the normalized query to identify identical queries (*optional*)
4. **Assign a time-stamp to the query**
   Execution time *or:* last update to the entire database *or:* last update to the subset of data affected by the query
5. **Compute a hash key of the result set** (*optional*)
6. **Assign PID to the query** *(if query/result set is new)*
7. **Store query and metadata** in query store

FACULTY OF !NFORMATICS

# Data Citation – Recommendations
## *(Draft, 3/4)*

## Upon request of a specific subset:

1. **PID resolves to landing page** of the subset, provides metadata including link to the super-set (PID of the DB)
2. **Landing page allows** (transparently, in a machine-actionable manner) to **retrieve the subset by re-executing the query**

- Query can be re-executed with the **original time stamp** or with the **current timestamp**, retrieving the semantically identical data set but incorporating all changes/corrections/updates applied since.

- Storing the query string provides comprehensive **provenance information** (description of criteria that the subset satisfies)

# Data Citation – Recommendations
## *(Draft, 4/4)*

**Upon modifications to the data management system:**

1. **When data is migrated** to a new representation (new DBMS system, new schema), the **queries need to be migrated**

2. **Hash keys** for the query strings may need to be **re-computed**

3. **Hash input function** for result set may need to be **adapted** to ensure that the result sets are presented in the same form to the hash function

4. **Successful re-writing should be verified** by ensuring that queries can be re-executed resulting in the correct result set hash key

FACULTY OF !NFORMATICS

# Initial Pilots

- Devised concept

- Identified challenges (unique sorting, hash-key computation, distribution, different data types, …)

- Evaluated conceptually in different settings
  - How to apply versioning/time-stamping efficiently?
  - How to perform query re-writing?
  - How easy to adopt / changes required within RI?

- Started implementing pilots
  - SQL: LNEC, MSD
  - CSV: MSD, open source prototype
  - XML: xBase

FACULTY OF !NFORMATICS

# Outline

- Joint Declaration of Data Citation Principles

- Challenges in non-trivial settings

- Recommendation of the RDA Working Group

- Pilots

- Summary

# WG Pilots

- Pilot workshops and implementations by
  - Various EU projects (TIMBUS, SCAPE,…)
  - NERC (UK Natural Environment Research Council Data Centres)
  - ESIP (Earth Science Information Partners)
  - CLARIN (XML, Field Linguistics Transcriptions)
  - Virtual Atomic and Molecular Data Centre

- Prototype solutions for
  - SQL, CSV, XML (partially)
  - LOD/RDF, triple-store DBs in the queue
  - Distributed data

FACULTY OF !NFORMATICS

# Dynamic Data Citation for SQL Data

## LNEC, MSD Implementation

FACULTY OF !NFORMATICS

# SQL Prototype Implementation

- LNEC Laboratory of Civil Engineering, Portugal

- Monitoring dams and bridges

- 31 manual sensor instruments

- 25 automatic sensor instruments

- Web portal
  - Select sensor data
  - Define timespans

- Report generation
  - Analysis processes
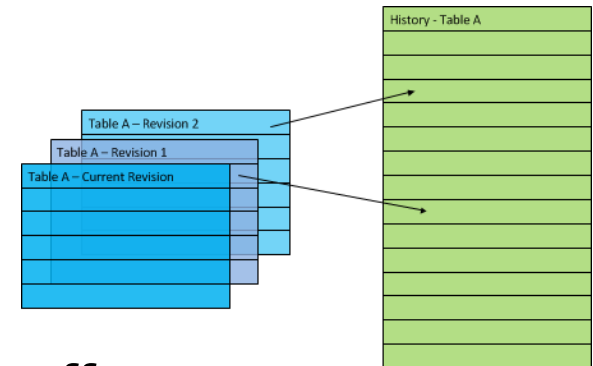  - LaTeX
  - publish PDF report

FACULTY OF !NFORMATICS

# SQL Prototype Implementation

- Million Song Dataset
  http://labrosa.ee.columbia.edu/millionsong/

- Largest benchmark collection in Music Retrieval

- Original set provided by Echonest

- No audio, only several sets of features
  (16 – 1440 measurements/features per song)

- Harvested, additional features and metadata
  extracted and offered by several groups
  e.g. http://www.ifs.tuwien.ac.at/mir/msd/download.html

- Dynamics because of metadata errors, extraction errors

- Research groups select subsets by genre, audio length,
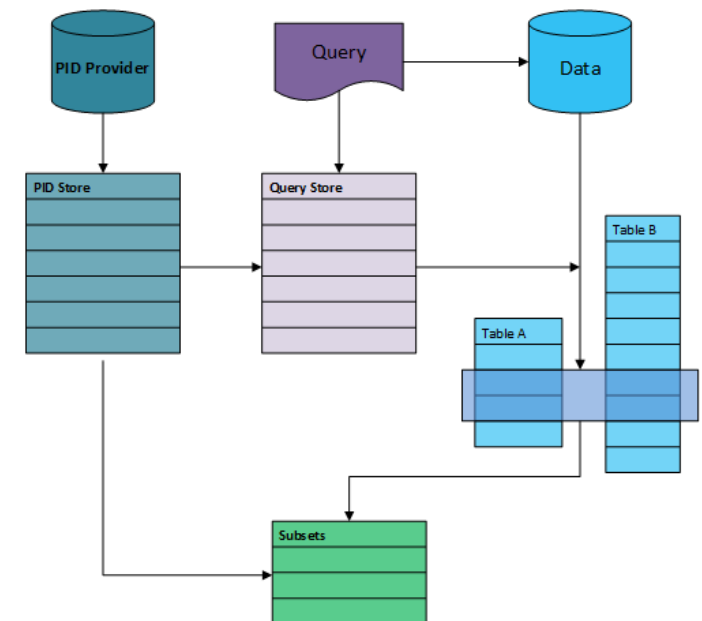  audio quality,…

FACULTY OF !NFORMATICS

# SQL Time-Stamping and Versioning

- ## Integrated
  - Extend original tables by temporal metadata
  - Expand primary key by record-version column

- ## Hybrid
  - Utilize history table for deleted record versions with metadata
  - Original table reflects latest version only

- ## Separated
  - Utilizes full history table
  - Also inserts reflected in history table



- ## Solution to be adopted depends on trade-off
  - Storage Demand
  - Query Complexity
  - Software adaption

FACULTY OF !NFORMATICS
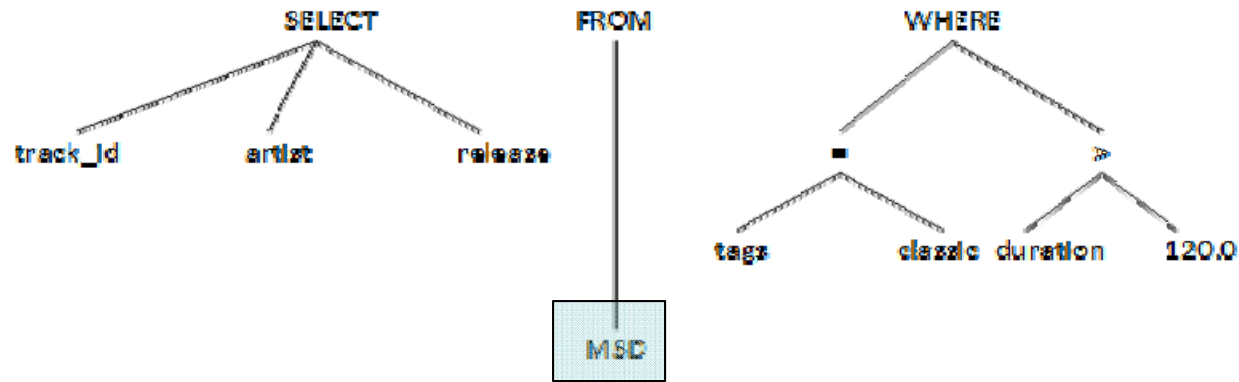
# SQL: Storing Queries

- Add query store containing
  - PID of the query
  - Original query
  - Re-written query + query string hash
  - Timestamp
    (as used in re-written query)
  - Hash-key of query result
  - Metadata useful for citation /
    landing page
    (creator, institution, rights, …)
  - PID of parent dataset
    (or using fragment identifiers for query)

FACULTY OF !NFORMATICS

# SQL Query Re-Writing

- Adapt query to history table



```
SELECT results.track_id, results.artist, results.release
    FROM MSD AS results JOIN (
        SELECT track_id, max(timestamp) AS latestTimestamp
        FROM MSD
        WHERE timestamp <= (SELECT @queryExecutionTimestamp)
        AND (track_id NOT IN
                (SELECT track_id FROM MSD AS deletedRecords
                        WHERE deletedRecords.status_mark = 'deleted'
                        AND (deletedRecords.timestamp < @queryExecutionTimestamp))
                )
        GROUP BY track_id
) AS version ON results.track_id = version.track_id AND results.timestamp = version.latestTimestamp

WHERE
    results.tags = 'classic'   AND results.duration> 120
ORDER BY results.track_id;
```
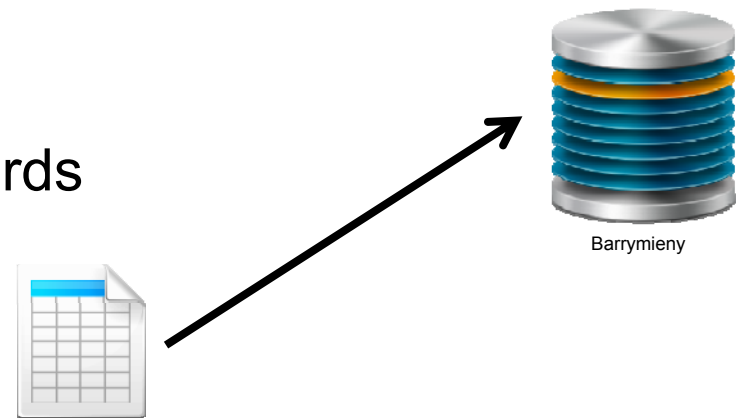
# Dynamic Data Citation - Pilots

# Dynamic Data Citation for CSV Data

## Open Source Reference Implementation

FACULTY OF !NFORMATICS

# Dynamic Data Citation for CSV Data

- Why CSV data? (not large, not very dynamic…)
  - Well understood and widely used
  - Simple and flexible
  - Most frequently requested during initial RDA meetings
- Goals:
  - Ensure cite-ability of CSV data
  - Enable subset citation
  - Support particularly small and large volume data
  - Support dynamically changing data
- 2 Options:
  - Versioning system (subversion/svn, git, …)
  - Migration to RDBMS

# CSV Prototype: Basic Steps

- Upload interface

  - Upload CSV files

- Migrate CSV file into RDBMS

  - Generate table structure, identify primary key

  - Add metadata columns for versioning

  - Add indices

- Dynamic data

  - Update / delete existing records

  - Append new data

- Access interface
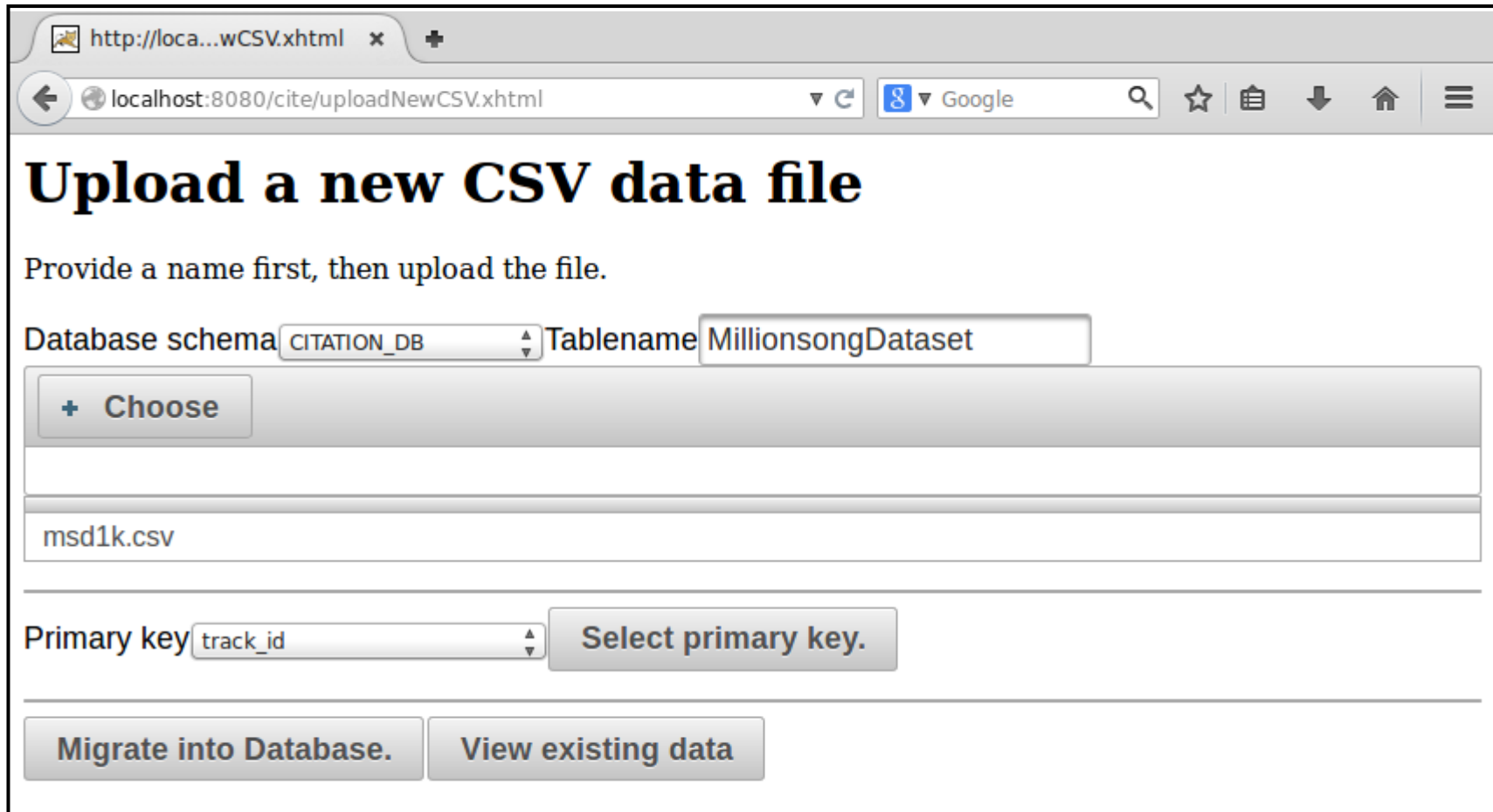
  - Track subset creation

  - Store queries

Barrymieny

# CSV Data Prototype

# CSV Data Prototype

# CSV Data Prototype

# CSV Data Prototype

| Suggested citation text: | Stefan Pröll (2015) "jj test" created at 2015-02-19 11:33:54.0, PID [ark:12345/5l86eH4qMX]. Subset of Stefan Pröll: "Adresses", PID [ark:12345/OjfL4gUmFo] |

## Download area

| Download CSV Subset | ↓ Download | Download the CSV data of this subset at the execution time of the query |
| Download Latest Subset | ↓ Download | Download the CSV data of this subset at its current state |
| Download Full DB | ↓ Download | Download the full database as CSV file |
| Download Diff CSV file | ↓ Download | Download the differences as CSV between the subset at its original execution time and now. |

# CSV Data Prototype

SQL string

(innerSELECT.RECORD_STATUS = 'inserted'          OR
innerSELECT.RECORD_STATUS = 'updated' AND
innerSELECT.LAST_UPD...
LAST_UPDATE) innerGrou...
innerGroup.LAST_UPDATE...
innerGroup.mostRecent  W...
UPPER('%jj%')  ORDER B...

**Öffnen von _tmp_CSV-Files_CitationDB_stefan_ad...   ✕**

Sie möchten folgende Datei öffnen:

📄 ...itationDB_stefan_adresses_12345-5l86eH4qMX.csv

Vom Typ: CSV-Dokument
Von: http://localhost:8080

**Wie soll Firefox mit dieser Datei verfahren?**

◉ Öffnen mit  [ LibreOffice Calc (Standard)        ▾ ]

○ DownThemAll!

○ Datei speichern

☐ Für Dateien dieses Typs immer diese Aktion ausführen

[ Abbrechen ]   [ OK ]

Suggested citation
text:

Stefan Pröll (2015) "jj test" ...                          ...5eH4qMX].
Subset of Stefan Pröll: "Adr...

## Download area

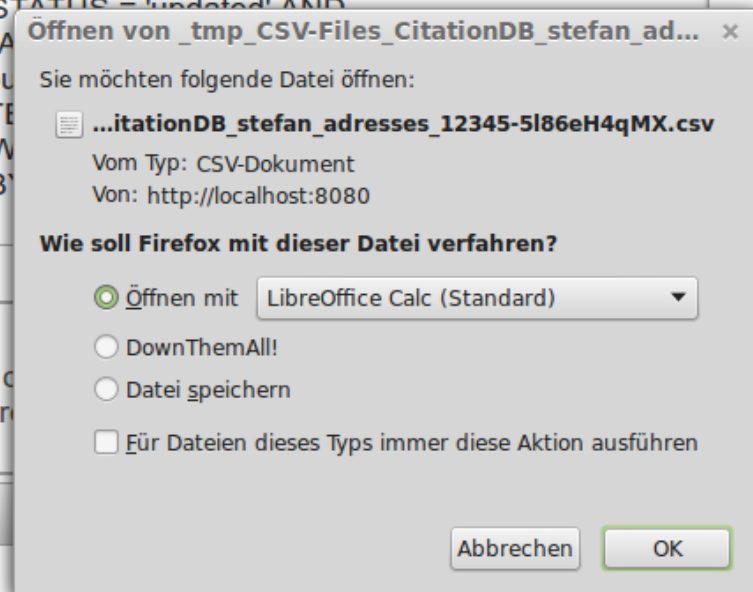| | | |
|---|---|---|
| Download CSV Subset | ↓ **Download** | Download the CSV data of this subset at the execution time of the query |
| Download Latest Subset | ↓ **Download** | Download the CSV data of this subset at its current state |
| Download Full DB | ↓ **Download** | Download the full database as CSV file |
| Download Diff CSV file | ↓ **Download** | Download the differences as CSV between the subset at its original execution time and now. |

FACULTY OF !NFORMATICS

**Dynamic Data Citation - Pilots**

# Progress update from VAMDC Distributed Data Centre

Carlo Maria Zwölf

Virtual Atomic and Molecular Data Centre
carlo-maria.zwolf@obspm.fr

FACULTY OF !NFORMATICS

# VAMDC

- Virtual Atomic and Molecular Data Centre

- Worldwide e-infrastructure federating 41 heterogeneous and interoperable Atomic and Molecular databases

- Nodes decide independently about growing rate, ingest system, corrections to apply to already stored data

- Data-node may use different technology for storing data (SQL, No-sql, ASCII files),

- All implement VAMDC access/query protocols

- Return results in standardized XML format (XSAMS)

- Access directly node-by-node or via VAMDC portal, which relays the user request to each node

# VAMDC

**Workshop prior to RDA P4**

**Issues identified**

- Each data node could modify/delete/add data without tracing
- No support for reproducibility of past data extraction

**Proposed Data Citation WG Solution:**

- Considering the distributed architecture of the federated VAMDC infrastructure, it seemed very complex to apply the "Query Store" strategy
  - Should we need a QS on each node?
  - Should we need an additional QS on the central portal?
  - Since the portal acts as a relay between the user and the existing nodes, how can we coordinate the generation of PID for queries in this distributed context?

FACULTY OF !NFORMATICS

# VAMDC

## Status / Progress since RDA P4

- Versioning adopted prior to P4

- Central service registering user interactions with data

- At each client SW notifies tracing service that a given **user** is using, at a given **time**, that specific **software** for submitting a given **query**

- Will assign single identifier for each unique query centrally

- Query store initially private (confidentiality issues)

# Further Pilots

- **NERC: UK Natural Environment Research Council**
  - ARGO buoy network: SeaDataNet
  - Butterfly monitoring, Ocean buoy network, National hydrological archive, …
- ESIP: BCO-DMO
- XML Data in Field Linguistics (CLARIN, XBase)
- Further Pilots on XML, LOD, …

- Workshops:
  - NERC Workshop, London, July 1/2  2014
  - ESIP Mtg in Washington, Jan 8 2015: Earth Science Data
  - Data Citation Workshop, Riva di Garda, April 20/21
  - Bilateral meetings with data centers

FACULTY  OF  !NFORMATICS

# Join RDA and Working Group

If you are interested in joining the discussion, contributing a pilot, wish to establish a data citation solution, …



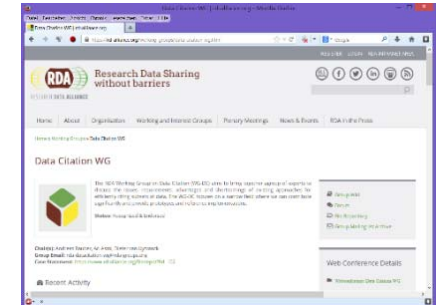- Register for the RDA WG on Data Citation:
    - Website:
      https://rd-alliance.org/working-groups/data-citation-wg.html
    - Mailinglist:
      https://rd-alliance.org/node/141/archive-post-mailinglist
    - Web Conferences:
      https://rd-alliance.org/webconference-data-citation-wg.html
    - List of pilots:
      https://rd-alliance.org/groups/data-citation-wg/wiki/collaboration-environments.html
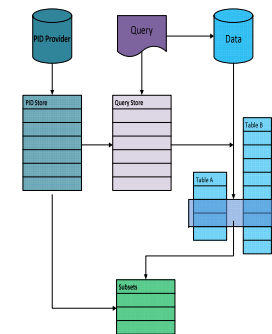
FACULTY OF !NFORMATICS

# Summary

- Trustworthy and efficient e-Science based on data

- Data as "1$^{st}$-class citizen"

- Support for identifying arbitrary subsets of dynamic data

  - Time-stamping and versioning of data

  - Storing (and citing) time-stamped queries

- Allows retrieving exact view on data set as used

- No need for artificial "versioning", delaying release of new data, or redundant storage of data subset dumps

- Helps tracing provenance (semantics) of data selection

- Future work: distributed datasets, data & query migration

FACULTY OF !NFORMATICS

# Thank you!

```
SELECT results.track_id, results.artist, results.release
        FROM MSD AS results JOIN (
                SELECT track_id, max(timestamp) AS latestTimestamp
                FROM MSD
                WHERE timestamp <= (SELECT @queryExecutionTimestamp)
                AND (track_id NOT IN
                        (SELECT track_id FROM MSD AS deletedRecords
                                WHERE deletedRecords.status_mark = 'deleted'
                                AND (deletedRecords.timestamp < @queryExecutionTimestamp))
                )
        GROUP BY track_id
) AS version ON results.track_id = version.track_id AND results.timestamp = version.latestTimestamp

WHERE
    results.tags = 'classic'  AND results.duration> 120
ORDER BY results.track_id;
```

## http://www.ifs.tuwien.ac.at/imp

FACULTY OF !NFORMATICS