

# Precisely Identifying Arbitrary Subsets of (Dynamic) Data: Recommendations of the RDA WGDC

Andreas Rauber

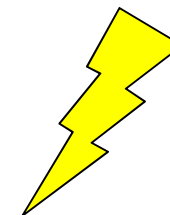
Technical University of Vienna  
Favoritenstr. 9-11/188  
1040 Vienna, Austria  
rauber@ifs.tuwien.ac.at  
<http://ww.ifs.tuwien.ac.at/~andi>

# Outline

- 
- Two challenges in data identification for citation
    - How to identify dynamic data?
    - How to deal with different granularity levels?
  - Recommendations of the RDA WGDC
  - Deployments
-

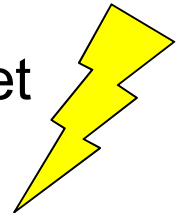
# Identification of Dynamic Data

- Usually, 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 identify precisely the **data as it existed at any(!) specific point in time**

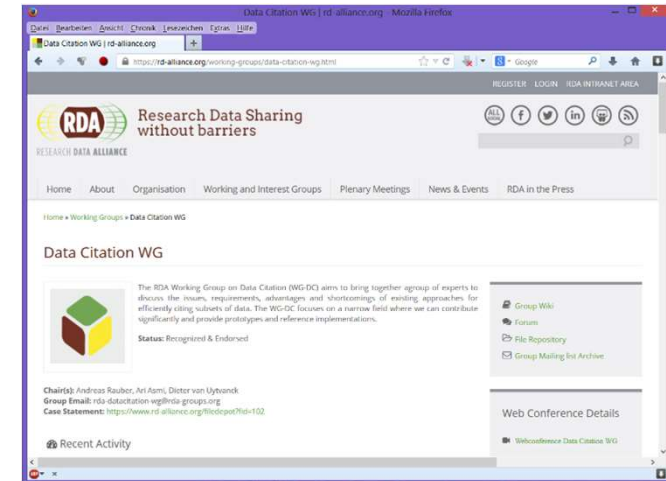


# Granularity of Subsets

- What about the **granularity** of data to be identified?
    - Enormous amounts of data
    - 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 identify **precisely & machine-actionably any subset of (dynamic) data used** in a process



- Research Data Alliance
- WG on **Data Citation: Making Dynamic Data Citeable**
- March 2014 – September 2015
  - Concentrating on the problems of **large, dynamic (changing) datasets**
- Final version presented Sep 2015 at P7 in Paris, France
- Endorsed September 2016 at P8 in Denver, CO
- Since then: supporting adopters



<https://www.rd-alliance.org/groups/data-citation-wg.html>

# RDA WGDC - Solution

- **We have**
  - **Any** kind of data & **some** means of access („query“)

# Dynamic Data Citation

**We have: Data + Means-of-access**

**Dynamic Data Citation:  
Cite (dynamic) data dynamically via query!**

# Dynamic Data Citation

**We have:** Data + Means-of-access

**Dynamic Data Citation:  
Cite (dynamic) data dynamically via query!**

**Steps:**

1. Data → versioned (history, with time-stamps)



# Dynamic Data Citation

**We have:** Data + Means-of-access

**Dynamic Data Citation:  
Cite (dynamic) data dynamically via query!**

## **Steps:**

1. Data → versioned (history, with time-stamps)

Researcher creates working-set via some interface:

**We have:** Data + Means-of-access

**Dynamic Data Citation:  
Cite (dynamic) data dynamically via query!**

**Steps:**

1. Data → versioned (history, with time-stamps)

Researcher creates working-set via some interface:

2. Access → **store & assign PID to “QUERY”**, enhanced with
- **Time-stamping** for re-execution against versioned data
  - **Re-writing** for normalization, unique-sort, ...
  - **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](http://www.ifs.tuwien.ac.at/~andi/publications/pdf/pro_ieeebigdata13.pdf)

# Data Citation – Deployment

- Researcher uses workbench to identify subset of data
- Upon executing selection („download“) user gets
  - Data (package, access API, ...)
  - PID (e.g. DOI) (Query is time-stamped and stored)
  - Hash value computed over the data for local storage
  - Recommended citation text (e.g. BibTeX)
- PID resolves to landing page
  - Provides detailed metadata, link to 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
- Query store aggregates data usage

## Data Citation – Deployment

- Note: query string provides excellent provenance information on the data set!
- subset of data user gets
  - Data (package, access API, ...)
  - PID (e.g. DOI) (Query is time-stamped and stored)
  - Hash value computed over the data for local storage
  - Recommended citation text (e.g. BibTeX)
- PID resolves to landing page
  - Provides detailed metadata, link to 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
- Query store aggregates data usage

# Data Citation – Deployment

- 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 (package)
  - PID (e.g. DOI)
  - Hash value
  - Recommended citation text (e.g. PID TEX)
- PID resolves to landing page
  - Provides detailed metadata, link to 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
- Query store aggregates data usage

# Data Citation – Deployment

- Note: query string provides excellent subset of data
  - provenance information on the data set! er gets
    - Data (pac
    - PID (e.g.
    - Hash valu
    - Recommended citation text (e.g. PID TEX)
- This is an important advantage over traditional approaches relying on, e.g. storing a list of identifiers/DB dump!!!
- PID resolves
    - Provides det
    - Option to ret

Identify which parts of the data are used. If data changes, identify which queries (studies) are affected
  - Upon activating PID associated with a data citation
    - Query is re-executed against time-stamped and versioned DB
    - Results as above are returned
  - Query store aggregates data usage

# Data Citation – Output

- 14 Recommendations grouped into 4 phases:
- 2-page flyer <https://rd-alliance.org/recommendations-working-group-data-citation-revision-oct-20-2015.html>
- Detailed report: Bulletin of IEEE TCDL 2016 [http://www.ieee-tcdl.org/Bulletin/v12n1/papers/IEEE-TCDL-DC-2016\\_paper\\_1.pdf](http://www.ieee-tcdl.org/Bulletin/v12n1/papers/IEEE-TCDL-DC-2016_paper_1.pdf)
- Adopter's reports, webinars <https://www.rd-alliance.org/group/data-citation-wg/webconference/webconference-data-citation-wg.html>
- Review / Lessons Learned  
Andreas Rauber et al., Precisely and Persistently Identifying and Citing Arbitrary Subsets of Dynamic Data  
Harvard Data Science Review, 3(4), 2021.  
DOI [10.1162/99608f92.be565013](https://doi.org/10.1162/99608f92.be565013).



# Paper: From Principles to Adoption

*Andreas Rauber, Bernhard Gößwein,  
Carlo Maria Zwölf, Chris Schubert, Florian  
Wörister, James Duncan, Katharina  
Flicker, Koji Zettsu, Kristof Meixner, Leslie  
D. McIntosh, Reyna Jenkyns, Stefan Pröll,  
Tomasz Miksa, and Mark A. Parsons:*

## **Precisely and Persistently Identifying and Citing Arbitrary Subsets of Dynamic Data.**

Harvard Data Science Review (HDSR),  
3(4), 2021.

DOI [10.1162/99608f92.be565013](https://doi.org/10.1162/99608f92.be565013)

- Principles
- 4 Reference implementations
- 8 Adoptions as Case Studies
- Lessons Learned





# Data Citation – Recommendations

## Preparing Data & Query Store

- R1 – Data Versioning
- R2 – Timestamping
- R3 – Query Store

## When Resolving a PID

- R11 – Landing Page
- R12 – Machine Actionability

## When Data should be persisted

- R4 – Query Uniqueness
- R5 – Stable Sorting
- R6 – Result Set Verification
- R7 – Query Timestamping
- R8 – Query PID
- R9 – Store Query
- R10 – Citation Text

## Upon Modifications to the Data Infrastructure

- R13 – Technology Migration
- R14 – Migration Verification



# Large Number of Adoptions

- **Standards / Reference Guidelines / Specifications:**
  - Joint Declaration of Data Citation Principles:  
Principle 7: Specificity and Verifiability (<https://www.force11.org/datacitation>)
  - ESIP:Data Citation Guidelines for Earth Science Data Vers. 2 (P14)
  - ISO 690, Information and documentation - Guidelines for bibliographic references and citations to information resources (P13)
  - EC ICT TS5 Technical Specification (pending) (P12)
  - DataCite Considerations (P8)
- **Reference Implementations**
  - MySQL/Postgres (P5, P6)
  - CSV files: MySQL, Git (P5, P6, P8, Webinar)
  - XML (P5)
  - CKAN Data Repository (P13)
  - SPARQL (P17, P19)

# Large Number of Adoptions

## ■ Adoptions deployed

- CBMI: Center for Biomedical Informatics, WUSTL (P8, Webinar)
- VMC: Vermont Monitoring Cooperative (P8, Webinar)
- CCCA: Climate Change Center Austria (P10/P11/P12, Webinar)
- EODC: Earth Observation Data Center (P14, Webinar)
- VAMDC: Virtual Atomic and Molecular Data Center (P8/P10/P12, Webinar)
- Ocean Networks Canada (P12, P20, Webinar)

# Lessons Learned as an FAQ (1 of 2)

- **Do the recommendations work for any kind of data?**  
Yes, it appears so.
- **Do all updates need to be versioned?**  
Ideally, yes. In practice, probably not (information accessed).
- **May data be deleted?** Yes, with caution and documentation.
- **What types of queries are permitted?**  
Any that a repository can support over time.
- **Does the system need to store every query?**  
No, just the relevant queries – “shopping cart”
- **Which PID system should be used?**  
The one that works best for your situation.
- **When multiple distributed repositories are queried, do we need complex time synchronization protocols?**  
No, not if the local repositories maintain timestamps.

# Lessons Learned as an FAQ (2 of 2)

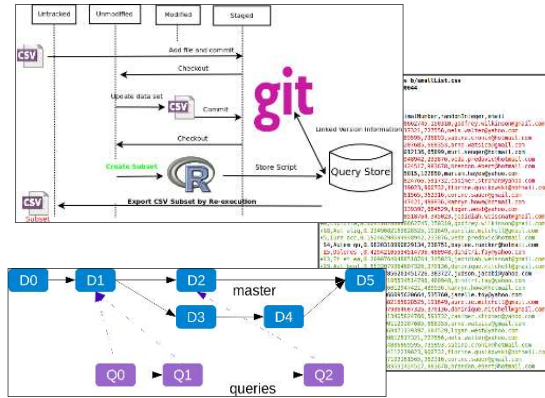
- **How does this support giving credit and attribution?**  
By including a reference to the overall data set as well as the subset.
- **How does this support reproducibility and science?**  
By providing a reference to the exact data used in a study.
- **Does this data citation imply that the underlying data is publicly accessible and shared? No.**
- **Why should timestamps be used instead of semantic versioning concepts?**  
Because there is no standard mechanism for determining what constitutes a 'version.' No minor/major "updates".
- **How complex is it to implement the recommendations?**  
It depends on the setting.
- **Why should I implement this solutions if my researchers are not asking for it or are not citing data?**  
Because it's the right thing for science.

# RDA Recommendations - Summary

## ■ **Benefits**

- Allows **identifying, retrieving and citing the precise data subset** with minimal storage overhead by only storing the versioned data and the queries used for extracting it
- Allows retrieving the data both **as it existed** at a given point in time as well as the **current view** on it, by re-executing the same query with the stored or current timestamp
- It allows to cite even an **empty set!**
- The query stored for identifying data subsets provides valuable **provenance data**
- Query store collects **information on data usage**, offering a basis for data management decisions
- **Metadata** such as checksums support the verification of the correctness and **authenticity** of data sets retrieved
- The same principles work for **all types of data**

# Thank you!



Label	Data Parameters	Platform / Stack	Method	Forecast Objective	Action	Gain
Repeat	-	-	-	-	-	Determinism
Param. Sweep	x	-	-	-	-	Robustness / Sensitivity
Generalize	10	x	-	-	-	Applicability across platform
Port	-	x	-	-	-	Portability across platform
Re-code	-	10	x	-	-	Correctness of implements
Validate	10	10	10	x	-	Correctness of hypothesis, different approach
Re-use	-	-	-	x	-	Apply code in different settings, Re-purpose
Independent x (orthogonal)	-	-	-	-	x	Sufficiency of information, independent verification

# WFs	1443
Final Data Set	731
Processor	# WFs % WFs
Not terminated >48hours	6 0.8
Execution failed	384 52.5
Execution successful	341 46.6

PostgreSQL Extension "temporal\_tables"

RDC table	sys_period
c1	
c2	
c3	

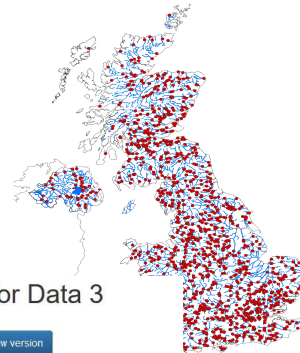
triggers

RDC.hist_table*	sys_po
c1	
c2	
c3	

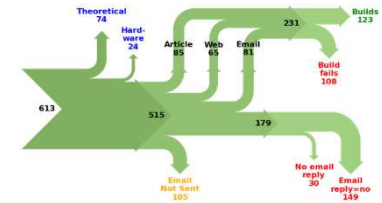
\*stores history of data changes

# Thanks!

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



DC<sup>1</sup>  
Data Citation Principles



## Editing data for Data 3

Back to Dataset Versions

Show changes Save to a new version

1 UPDATE Z0001\_test SET 'SiteID' = 'Stevensville Brook' WHERE db\_table\_pk=30  
2 DELETE FROM Z0001\_test where db\_table\_pk=30  
3 DELETE FROM Z0001\_test where db\_table\_pk=35

2010

Actions	SiteID	LabID	Date	MeanDensity	Mean
	Stevensville Brook	2000.187	0000-00-00	4644322354	39.0
	Winhall River	2011.081	2011-10-07	201	47.5
	Winhall River	2012.089	2012-09-27	1981	52.0
	Winhall	2013.150	2013-10-15	1002	30.0

