

A curation interface for temporal databases

Vashti Galpin¹, Simon Fowler², James Cheney¹

University of Edinburgh¹

University of Glasgow²

International Data Curation Conference 2021 (IDCC21)

19 April 2021



THE UNIVERSITY of EDINBURGH
informatics



University
of Glasgow

Introduction: Curated scientific databases

Curated databases:

- “databases that are populated and updated with a great deal of human effort” (Buneman et al., 2008)
- Examples: IUPHAR/BPS Guide to PHARMACOLOGY (GtoPdb), CIA World Factbook, Manually Curated Database of Rice Proteins (MCDRP)

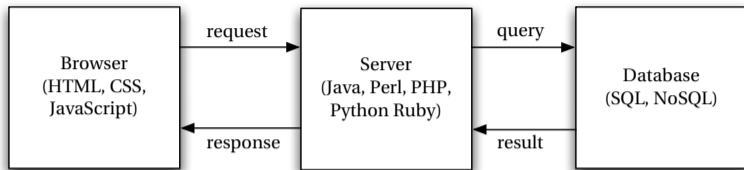
Curation challenges:

- Provenance tracking, citation, archiving, annotation propagation ...
- Digital tools have to be hand-crafted for every curated DB — often with very stretched resources

Ultimate goal: programming language support for curation functionality

- Cross-tier web-programming using the **Links** language
- Previous work: GtoPdb re-implementation (Fowler et al, 2020)
- This research: **prototype** curation interface for Covid-19 data

Background: Cross-tier web programming



Curated databases normally a **collection of applications**:

→ Database, web frontend, curation application

Links: cross-tier programming language with language-integrated query:

- Client, server, and database code written in same language
- Database queries written in the same language providing well-formed and efficient queries
- Language developed at UoE since 2006, with temporal database features recently added

Provenance and curation

Definition of provenance:

- “Essentially, provenance can be seen as meta-data that, instead of describing data, describes a production process.” (Herschel et al, 2017)
- Here our focus is on data provenance in the context of **data updates**

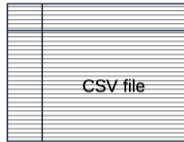
Provenance facilitates curation:

- Captures the process of data creation and allows analysis of data change
- Records the human effort and decisions involved in curation
- Allows for assessment of data quality and integrity

Provenance supports data sharing:

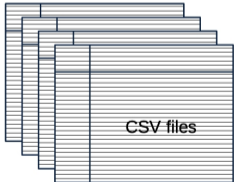
- Availability of provenance information can build trust for data sharing and reuse
- Provides information for data versioning via time-slicing

Temporal databases and provenance queries



How many deaths were reported in the week of 5 July?

How were those deaths distributed across the health boards?



How many deaths were first reported in week of July 5?

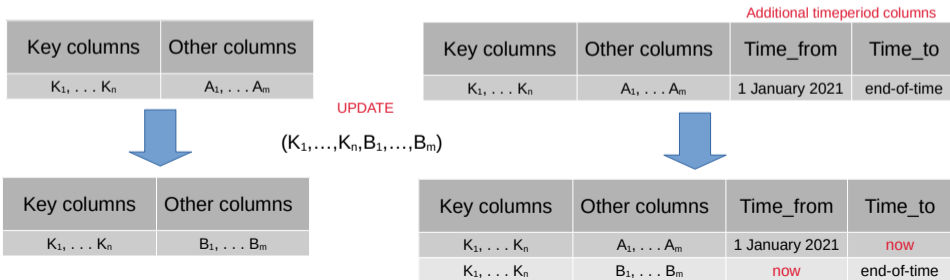
How does that differ from the most recently reported figure?

Which health boards reported the most changes?

Temporal databases and update provenance

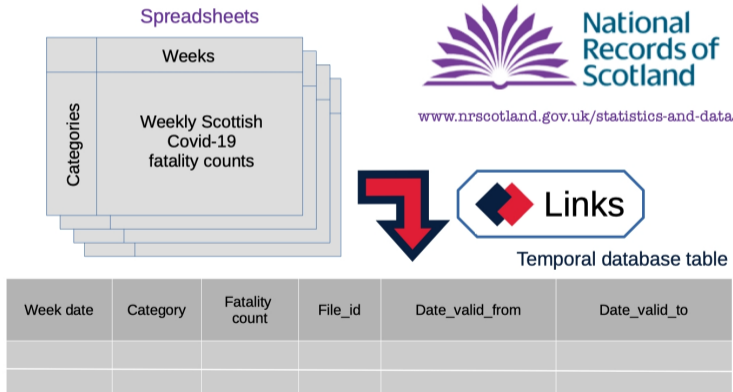
Temporal databases:

- Introduction of time period information to individual rows by adding two columns: the start of the time period and the end of the time period
- The time period describes when the data is valid and temporal queries use this time information
- Interpretation of time period is application dependent: transaction time or another metric of validity such as validity in the real world.



Case study: Scottish Covid-19 figures

- Weekly CSV files with weekly Covid-19 fatality counts for sex, age, health board, local authority, and location
- Each CSV file contains new data for (at least) one week and may contain updates to counts from previous weeks



Curation interface prototype: update decisions

Scottish Weekly Covid Data Curation Interface

Overview

Query ▾

Pending

Upload

Other ▾

Pending

Updates for the week of 2020-03-30 arising in the week of 2020-04-13

Type	Category	Week	Old value	New value	Change	Time added
All	All	2020-03-30	282	283	1	2021-04-15 12:39
Sex	Female	2020-03-30	126	127	1	2021-04-15 12:39
Age	75-84	2020-03-30	106	107	1	2021-04-15 12:39
AgeF	F: 75-84	2020-03-30	49	50	1	2021-04-15 12:39
HB	Greater Glasgow and Clyde	2020-03-30	106	107	1	2021-04-15 12:39

Accept all updates

Reject all updates

Consider each update individually

Continue

Updates for the week of 2020-04-06 arising in the week of 2020-04-13

Type	Category	Week	Old value	New value	Change	Time added
All	All	2020-04-06	608	610	2	2021-04-15 12:39
Sex	Female	2020-04-06	261	262	1	2021-04-15 12:39
Sex	Male	2020-04-06	347	348	1	2021-04-15 12:39

Curation interface prototype: provenance queries

Scottish Weekly Covid Data Curation Interface

OverviewQueryPendingUploadOther

Provenance

Data item modification history

Select category

All

Select week

2020-03-30

Lookup

The data item for category All and week 2020-03-30 has the following change history.

Week of Modification	Value	Date of Modification
2020-03-30	282	2021-04-15 12:39:27.549333
2020-04-13	283	2021-04-15 12:41:6.636032
2020-04-20	282	2021-04-15 12:42:50.027521

Continue

Data items with at least one modification

Data

Provenance: data items

Provenance: data categories

Provenance: weeks

Provenance: rejected updates

Conclusion: Links for curation interfaces

→ **Curated databases:**

- Databases maintained using much human effort
- Consist of multiple applications: database, web frontend, ...
- How can we automate support for curation?

→ **Curation interface:**

- Links is a cross-tier programming language with client, server, database code in a single language
- Temporal database features support update provenance
- A bit like Github but for data

See our poster **Links between Temporal Databases and Curation**
and complete our curation interface survey at **RDA Plenary 17**

Vashti.Galpin@ed.ac.uk

<http://www.links-lang.org>