



XArch: Archiving Scientific and Reference Data

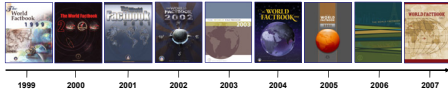
Heiko Müller, Peter Buneman, Ioannis Koltsidas

University of Edinburgh – School of Informatics – Database Group



Why do we archive data?

Databases are modified continuously



Archiving is a necessity to ...

- ... maintain **access** to **older versions** of the data ...
 - backup,
 - verification of findings,
 - citation.
- ... **track history** of objects.

How did the population of European countries change over the last ten years?



How databases are archived?

- Complete periodic snapshot**
 - High storage overhead.
- Delta-based approaches**

Snapshot and records of changes between pairs of consecutive versions.

 - Work on lines of text not data objects.
 - Sensitive to formatting/layout.
 - Snapshot retrieval is bounded by the number of deltas not data size.
 - History tracking is complex.

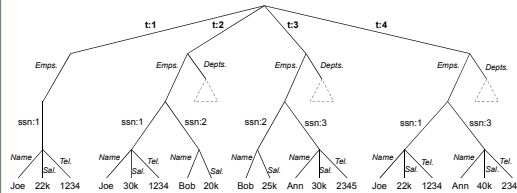
How do we archive data?

- Hierarchical structure.
- Merge versions into a single archive.
- Benefits include ...
 - Retrieval overhead is reduced.
 - Reduction of storage overhead.
 - History tracking is possible.
 - Stored in human readable format.

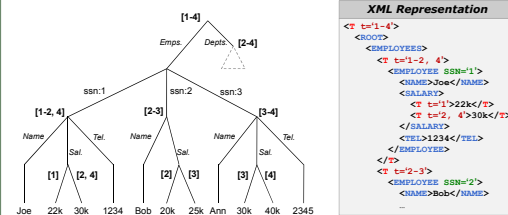
Buneman, Khanna, Tajima, Tan, ACM TODS, 2004.

Pushing time down

A sequence of database versions ...



... merged into a single archive.

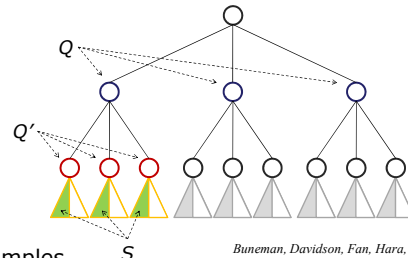


Relies on a deterministic / keyed model (Driscoll et al., "Making Data Structures Persistent", 1989).

Keyed hierarchical data model

Keys for elements in hierarchical data

- Elements Q'** are keyed relative to their **parent Q** by part of their **subtree S**.

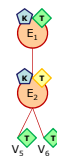


Examples

KEY /EMPLOYEES/EMPLOYEE BY VALUES (SSN),
 KEY /EMPLOYEES/EMPLOYEE/NAME BY EXISTENCE,
 KEY /EMPLOYEES/EMPLOYEE/SALARY BY EXISTENCE,
 KEY /EMPLOYEES/EMPLOYEE/TEL BY SUBTREE,

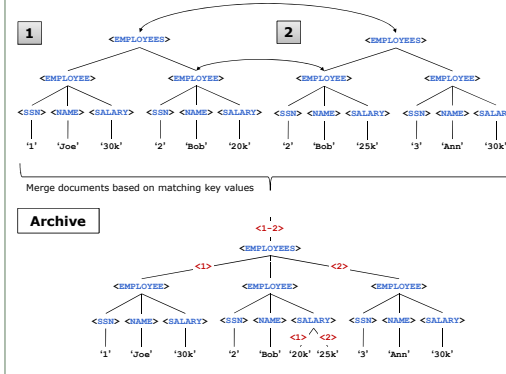
Conceptual data model for archives

Each internal node **E** has a key value **K** associated with it. Leaf nodes **V** are text values. Each node has its own timestamp **T** or an inherited timestamp **D**. Timestamps denote the sequence of versions a node occurred in.



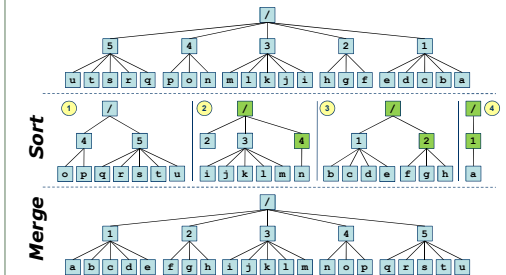
Building archives by Nested-Merge

- Corresponding elements are identified based on their key values.
- Children are merged recursively.
- Merging is done efficiently for sorted documents.



Sorting Hierarchical Data

- Split** the document **vertically** into sorted runs.
- Sort** siblings based on **key values**.
- Merge sorted runs** into final document.

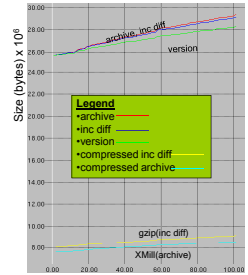


Koltsidas, Müller, Viglas, VLDB, 2008 (to appear).

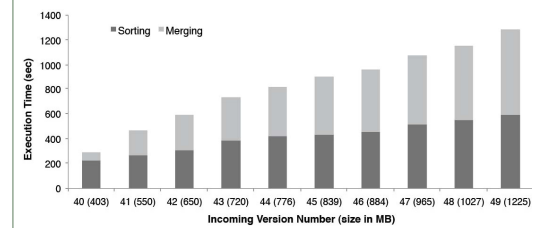
Experiments

100 days of OMIM:

- Uncompressed
 - Archive size is
 - ≤ 1.01 times diff repository size.
 - ≤ 1.04 times size of largest version.
- Compressed
 - Archive size between **0.94** and **1** times compressed diff repository size.
 - gzip - unix compression tool
 - XMill - XML compression tool



SWISS-PROT Major Releases 40 – 49:

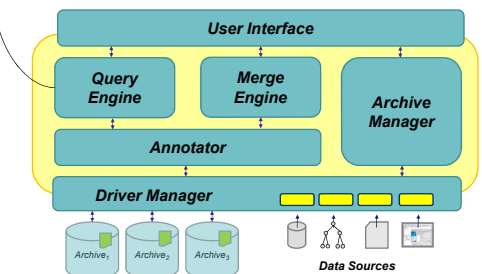


System Architecture

How did the data look like in version x_1, \dots, x_2 ?
 How did the data evolve?
 Which objects changed (or remained unchanged)?
 When was a given condition true?

Currently supported input formats:

- XML.
- Relational Databases.
- UNIPROT flat-files.
- CIA World Factbook Web-pages.



- Other related topics
 - Data provenance.
 - Data provenance.
 - Citation for databases.

For further information

- www.lfcs.inf.ed.ac.uk/research/database.
- www.dcc.ac.uk.