

## Archiving Scientific Data with XMLArch

### Why Archive Data?

- Backup!
- Data history required for ...
  - Verification of findings.
  - Citation.
  - History tracking.

### How to Archive Data?

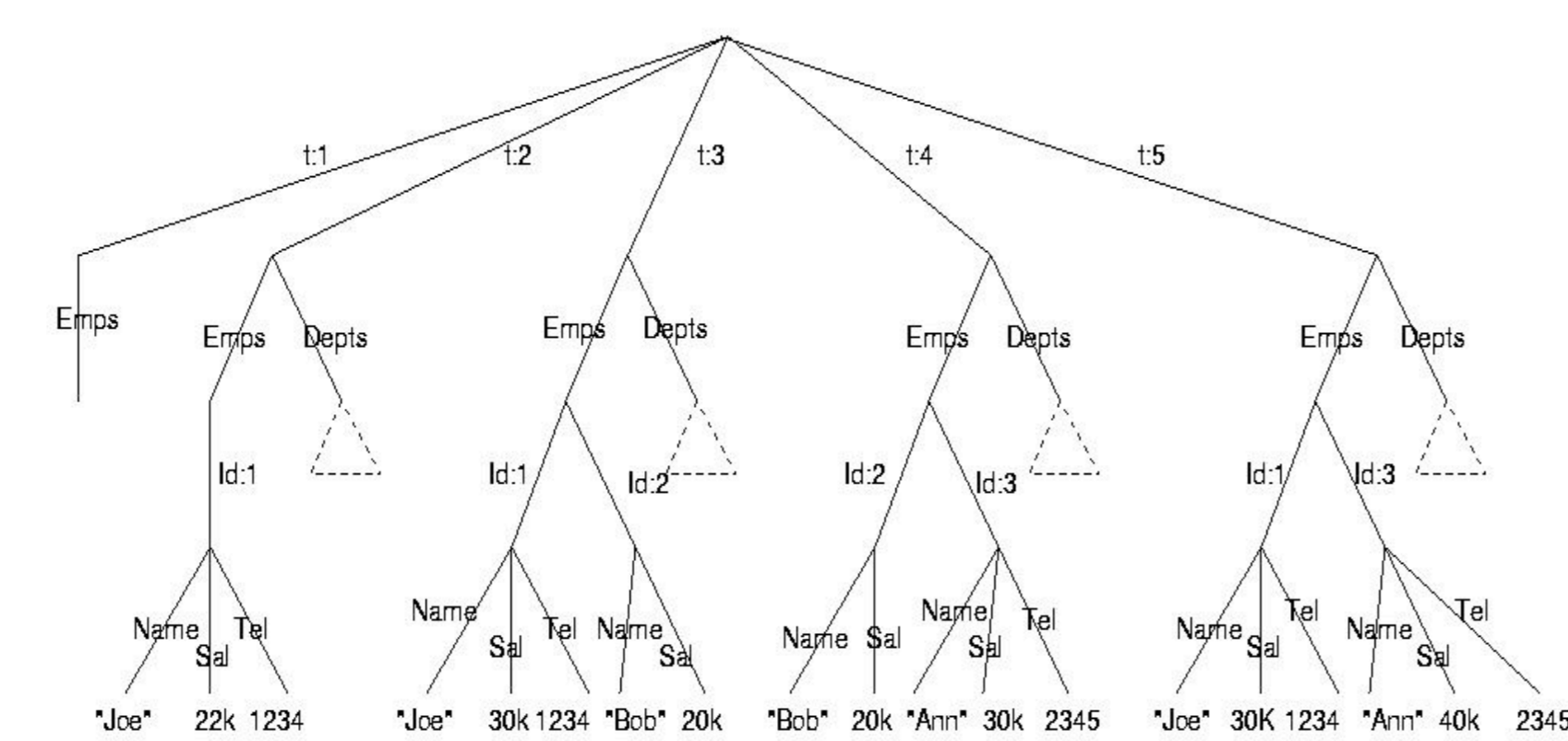
- Complete periodic snapshot.
  - High storage overhead.
- Diff-based approaches
  - Snapshot + incremental diff.
  - Snapshot + reverse diff.
- Snapshot + transaction log capture
- Some combination of the above.
- Problems of diff-based approaches
  - Work on lines of text not data objects.
  - Know nothing about the domain.
  - Cannot track object history.
  - Sensitive to formatting/layout.
  - Retrieval is bounded by the number of diffs not data size.
  - History tracking is complex.

### How do WE build Archival Databases!

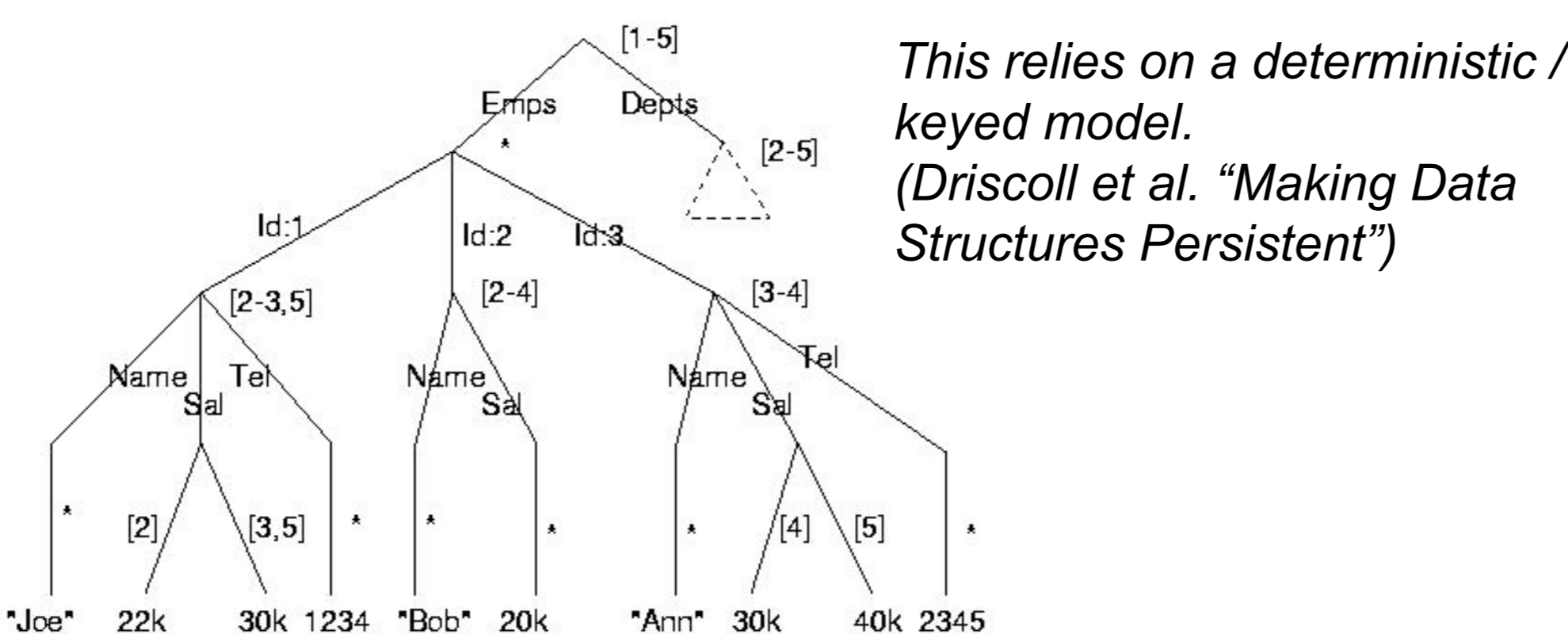
(Buneman, Khanna, Tajima, Tan, 2004)

- Hierarchical structure (XML).
  - Use unique identifiers.
- Merge all versions into a single archive.
- Benefits include ...
  - Retrieval overhead is reduced.
  - History tracking is possible.
  - Reduction of storage overhead.
  - Stored in human readable format.

### A Sequence of Versions



### Pushing time down



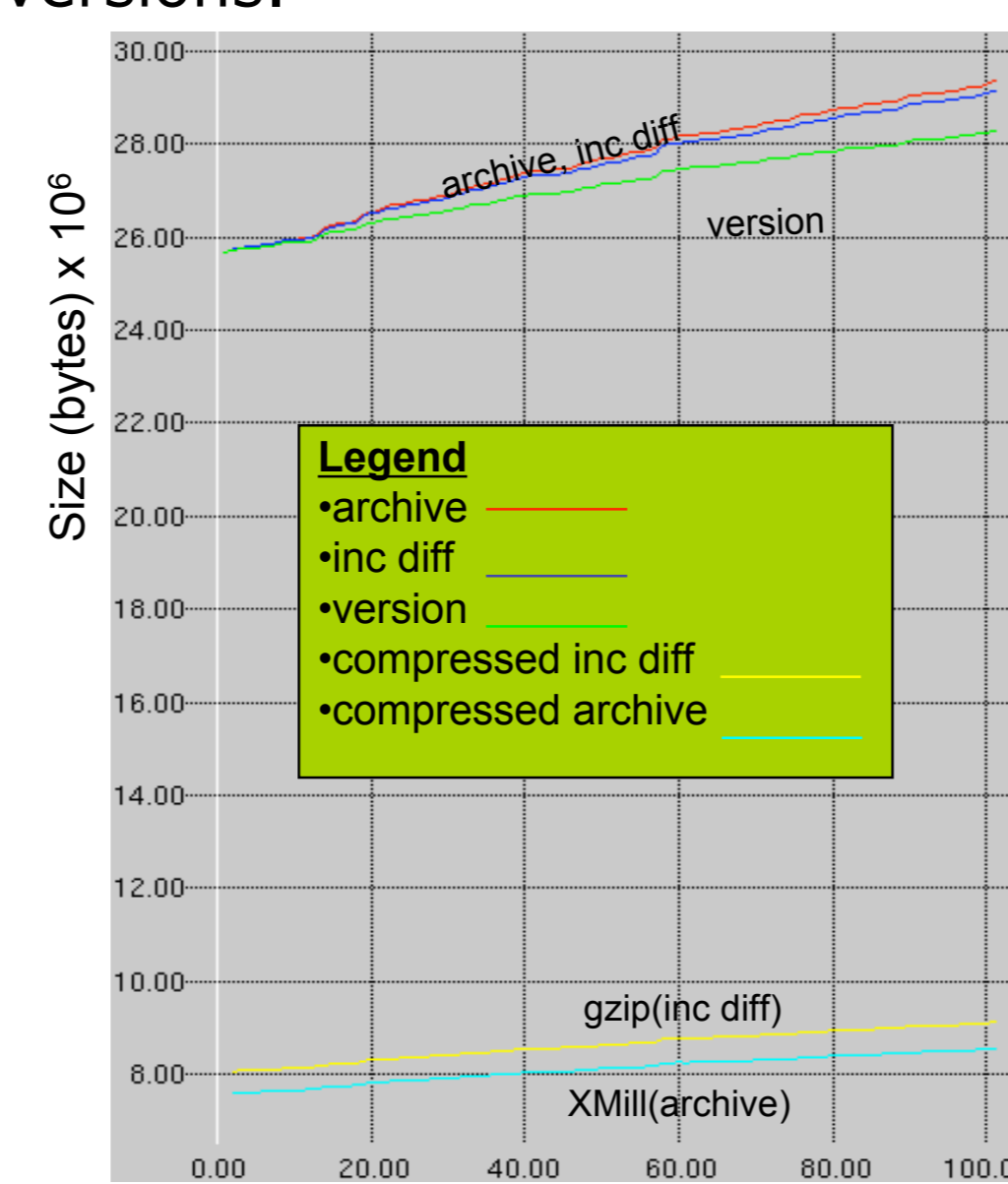
### 100 days of OMIM

- Recorded all OMIM versions for about 14 weeks.
- XML-ized all of them.
- Combined into XML format archive by pushing time down.
- Also recorded diffs between versions.

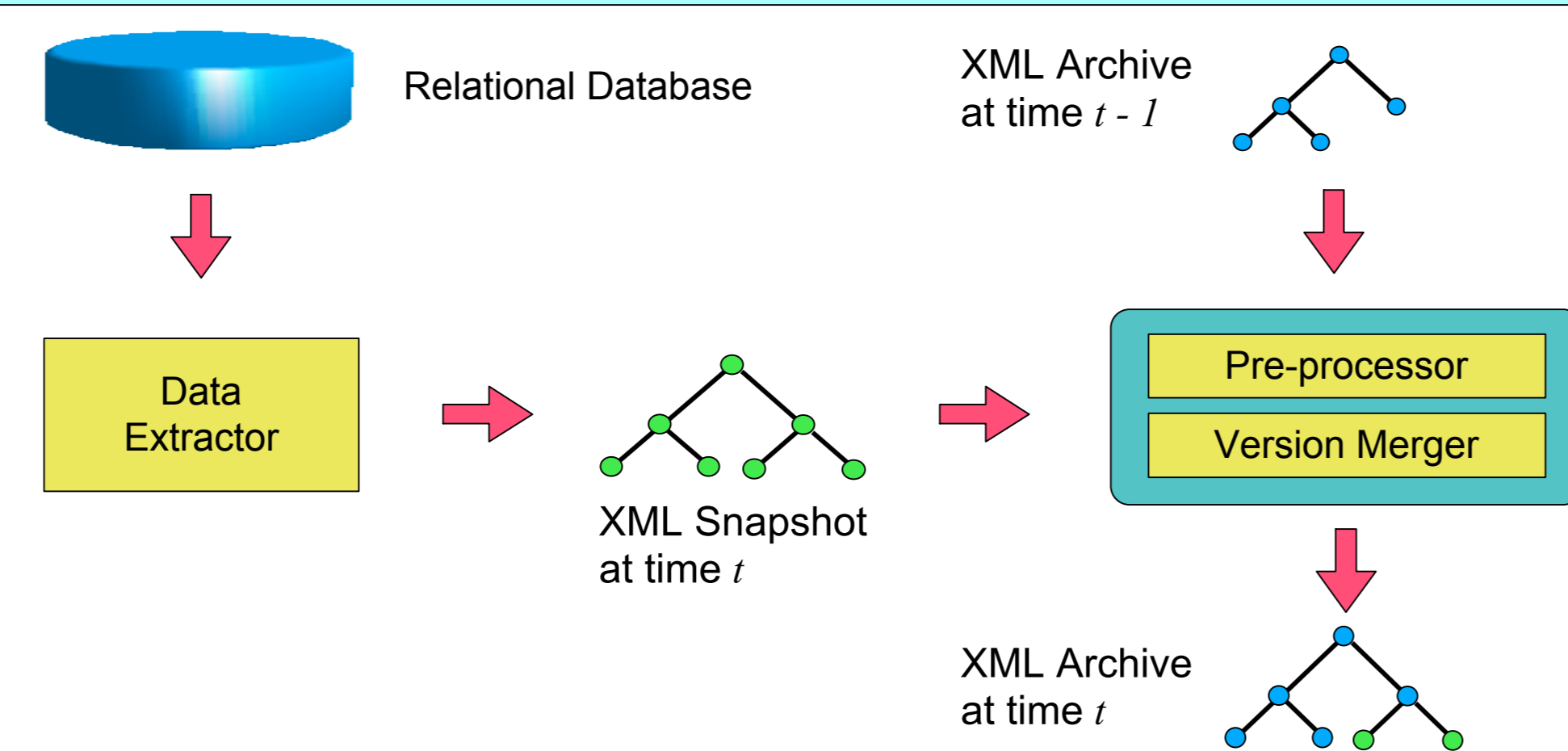
- Uncompressed
  - Archive size is
    - $\leq 1.01$  times diff repository size.
    - $\leq 1.04$  times size of largest version.

- Compressed
  - Archive size between **0.94** and **1** times compressed diff repository size.

- gzip - unix compression tool
- XMill - XML compression tool



### Architecture

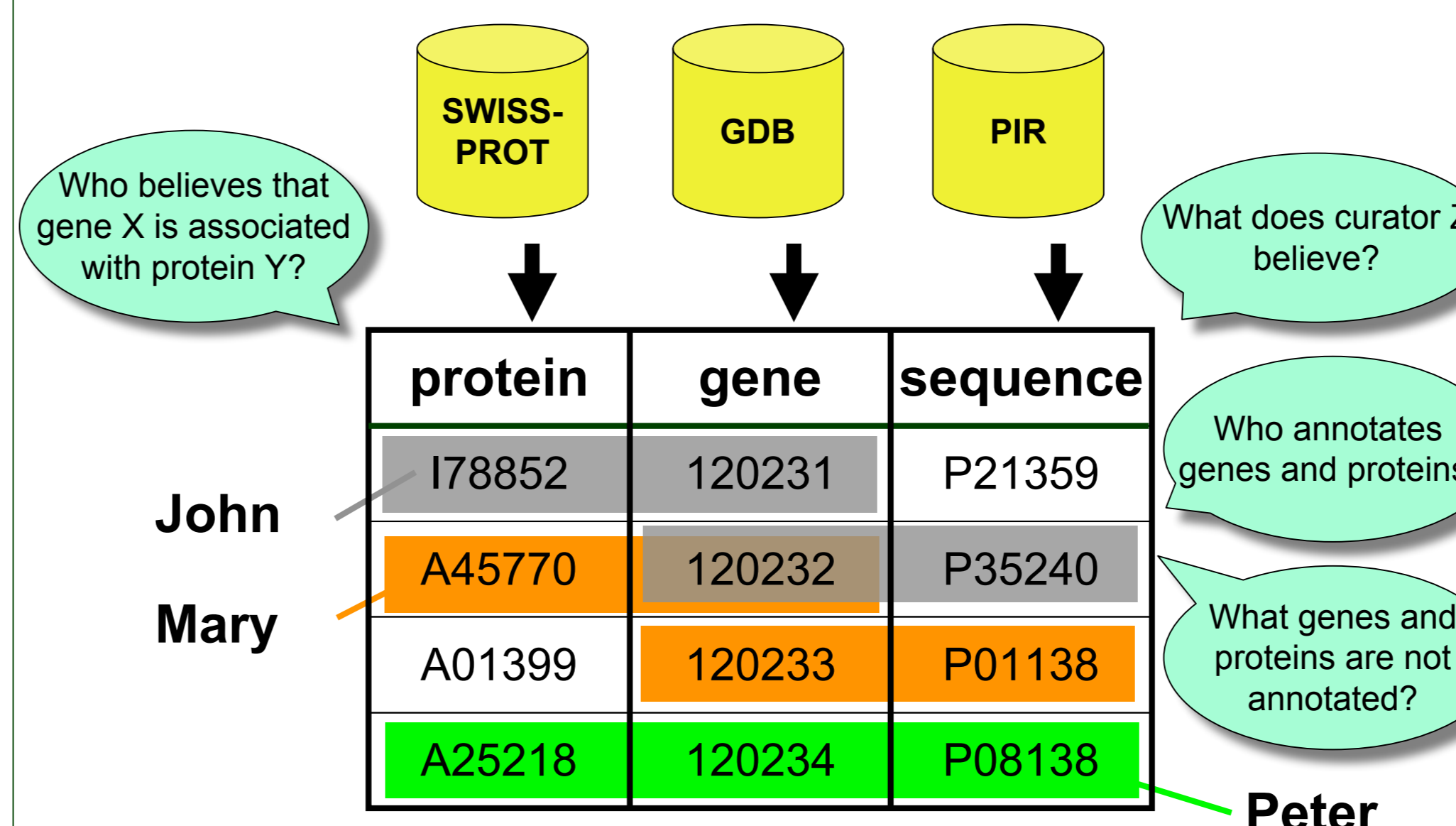


## Mondrian: Annotating and Querying Databases through Colors and Blocks

### Annotation in Databases

- Annotations play a central role in database curation.
- DBMS lack support for ...
  - Modeling different granularity.
  - Storing and querying annotations.
- Our contributions
  - Annotation oriented data model using colored blocks to represent annotated sets of values.
  - Color query language that is minimal, sound, and complete.

### Motivating Example



### How to color databases?

- We need ...
  - Set of colors  $C$ .
  - A coloring function  $x$  that accepts
    - a tuple  $t$  in a instance  $r$  of relation  $R$ .
    - a set of attributes  $Y \subseteq \text{sort}(R)$ .
 and assigns a set of colors  $C' \subseteq C$  to  $t[Y]$ .
- A colored database is a set of blocks  $(t, Y, x(t, Y))$ , e.g.,  $x(t_1, \{protein, gene\}) = \{grey\}$ .

### How to query the database?

- We introduce a color algebra (CA)
  - selection, projection, product, renaming, union ...
  - block selection ( $\Sigma$ ), block projection ( $\Pi^L, \Pi^U$ ), merge ( $\mu$ ), recolor ( $\rho$ ).
- Prove the algebra to be **minimal, sound, and complete**.
- Benefits include ...
  - appropriate level of abstraction.
  - respect semantics of colors & blocks.
  - easy to use and portable.

### Block projection & Selection

$\Pi^L$ gene, sequence (r):	protein	gene	sequence
Find tuples with annotations involving <b>at least</b> genes & sequences.	A45770	120232	P35240
	A25218	120234	P08138

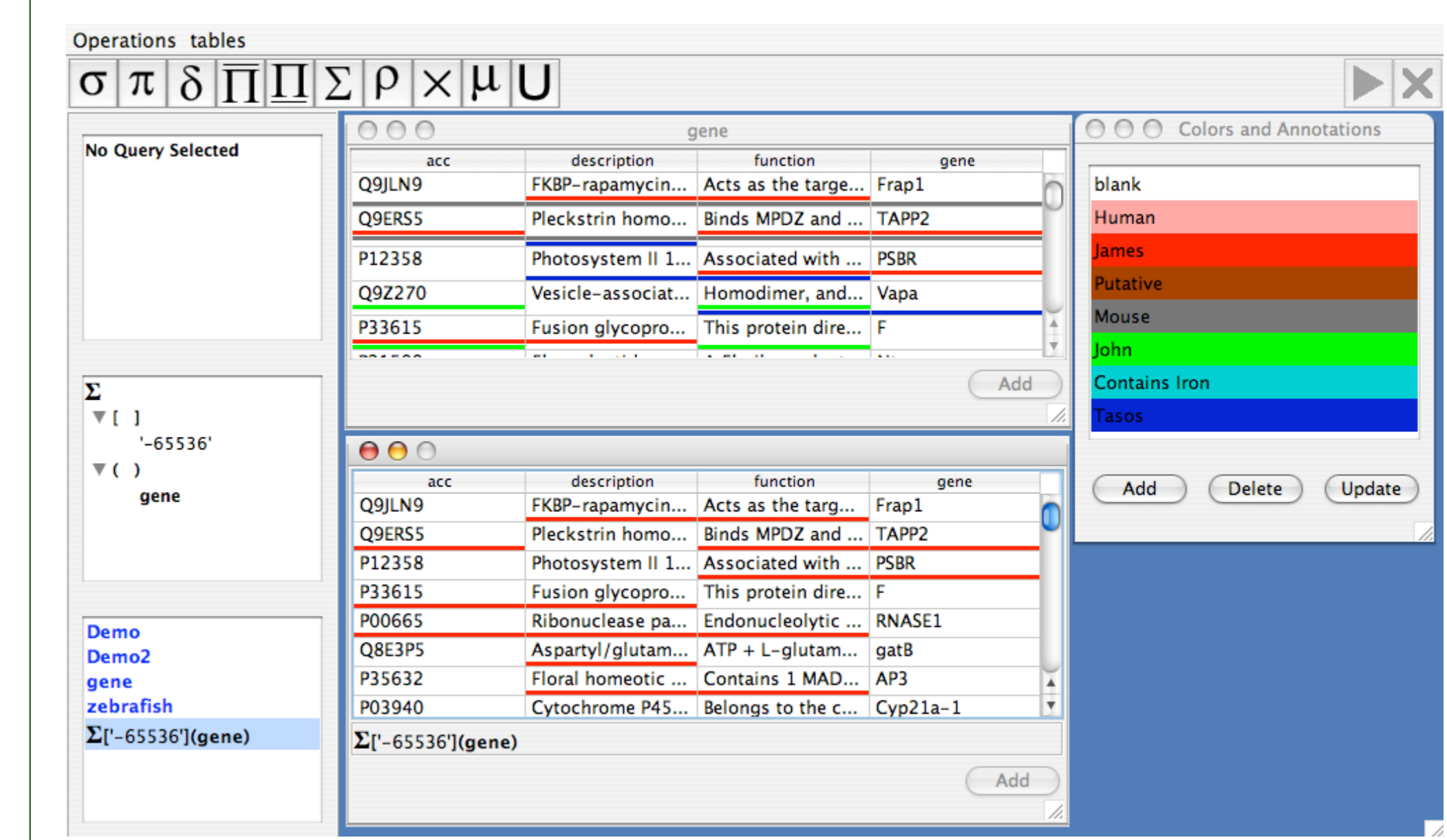
  

$\Pi^U$ gene, sequence (r):	protein	gene	sequence
Find tuples with annotations involving <b>at most</b> genes & sequences.	A45770	120232	P35240

$\Sigma$ grey (r):	protein	gene	sequence
Find tuples with grey annotations.	178852	120231	P21359
	A45770	120232	P35240

### Mondrian



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

- For further information
  - [www.lfcs.inf.ed.ac.uk/research/database](http://www.lfcs.inf.ed.ac.uk/research/database).
  - [www.dcc.ac.uk](http://www.dcc.ac.uk).