# A provenance-based approach to manage long term preservation of scientific data
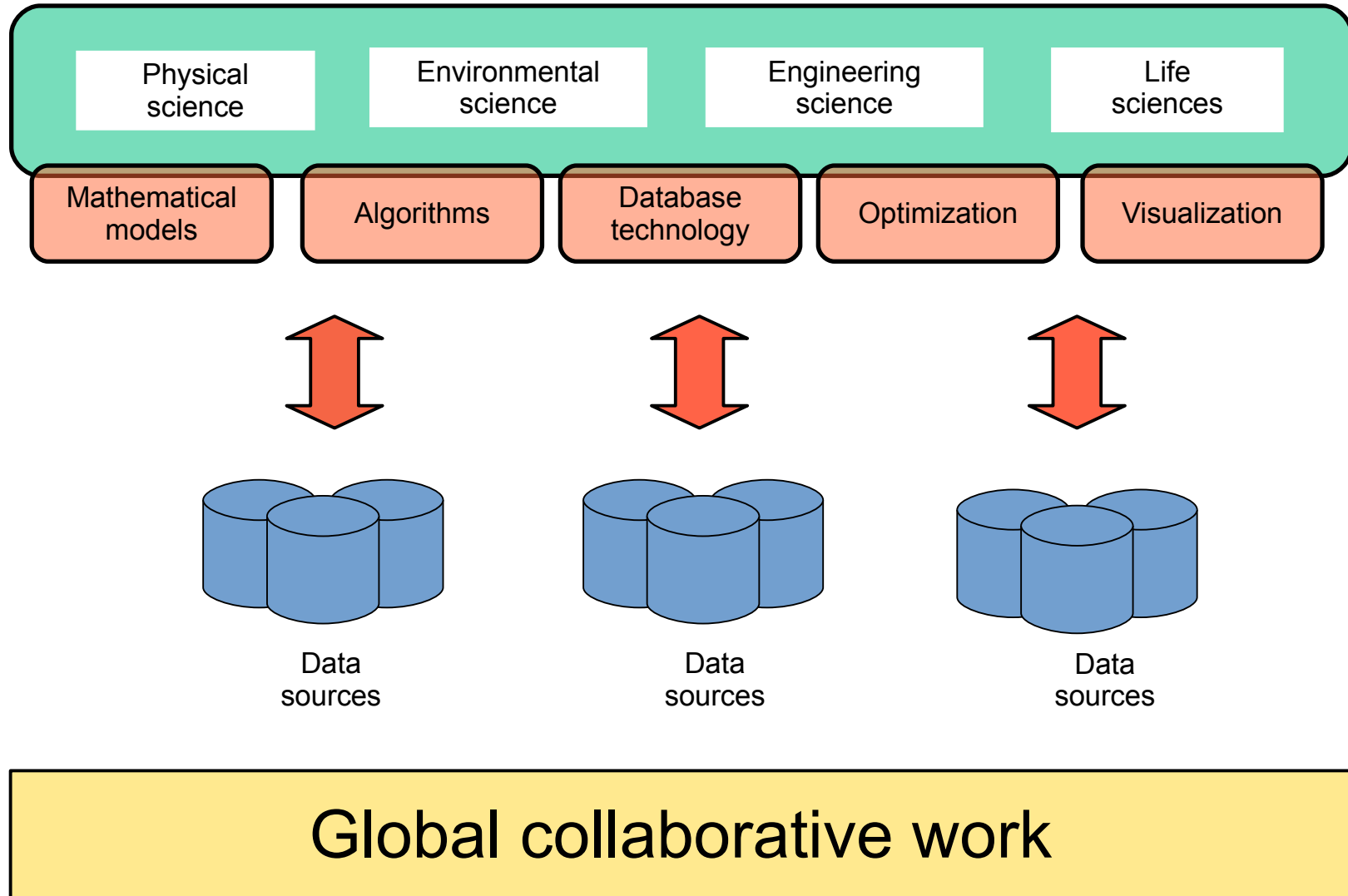
Renato Sousa, Daniel Cugler, Joana Malaverri,
Claudia Bauzer Medeiros,
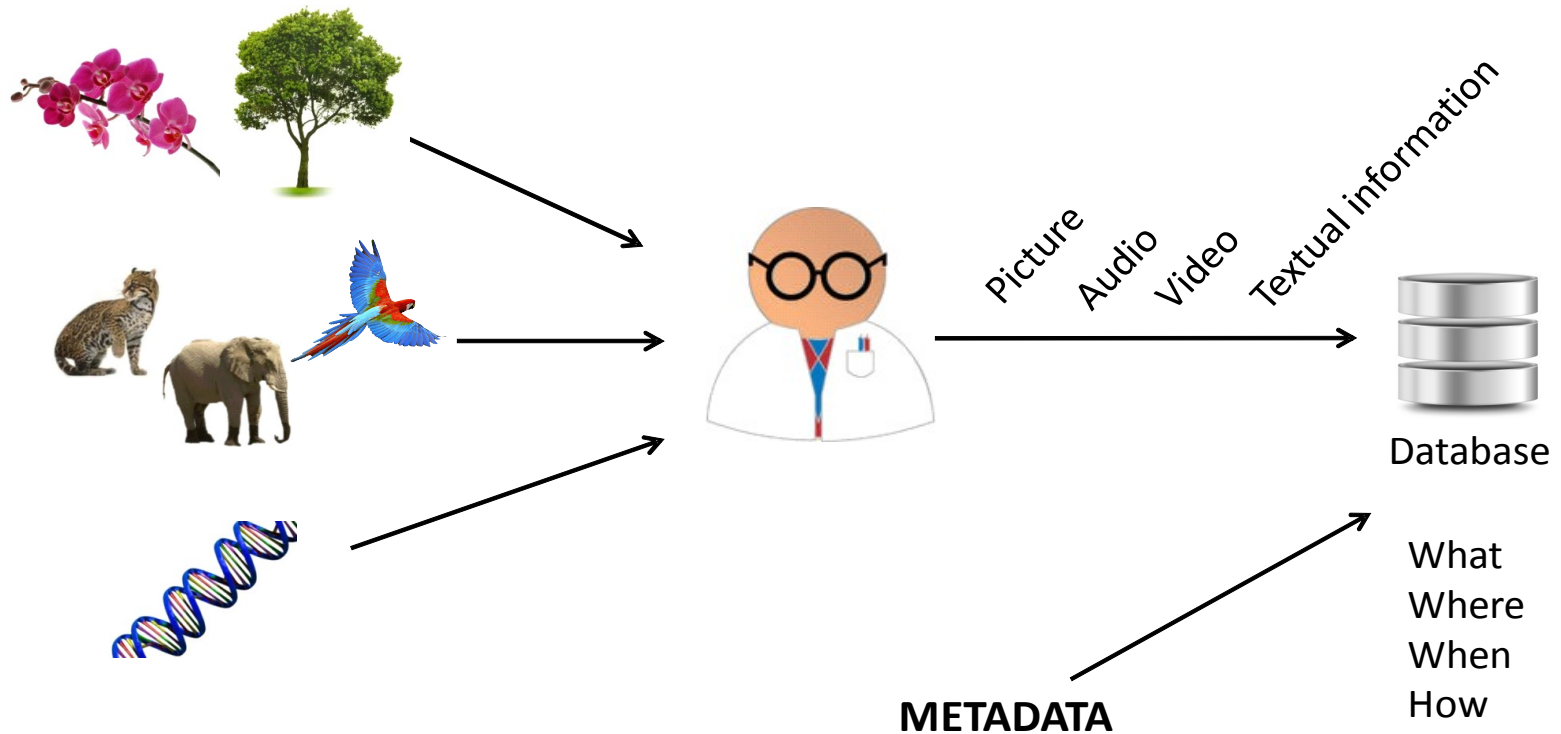Institute of Computing – University of Campinas

BRAZIL

# Outline

- Motivation and related work

- Our approach

- Case Study

- Conclusions and Ongoing Work

# eScience environments

| Physical science | Environmental science | Engineering science | Life sciences |
| --- | --- | --- | --- |

| Mathematical models | Algorithms | Database technology | Optimization | Visualization |
| --- | --- | --- | --- | --- |

Data sources

Data sources

Data sources

## Global collaborative work

# Motivation – Biological Observation Databases



Picture  Audio  Video  Textual information

Database

METADATA

What
Where
When
How

# What is data?

- Any (digital) result of scientific experiments
- Any (digital) input to experiments

- (Ideally, consider software, workflows, documentation, intermediate artifacts...)
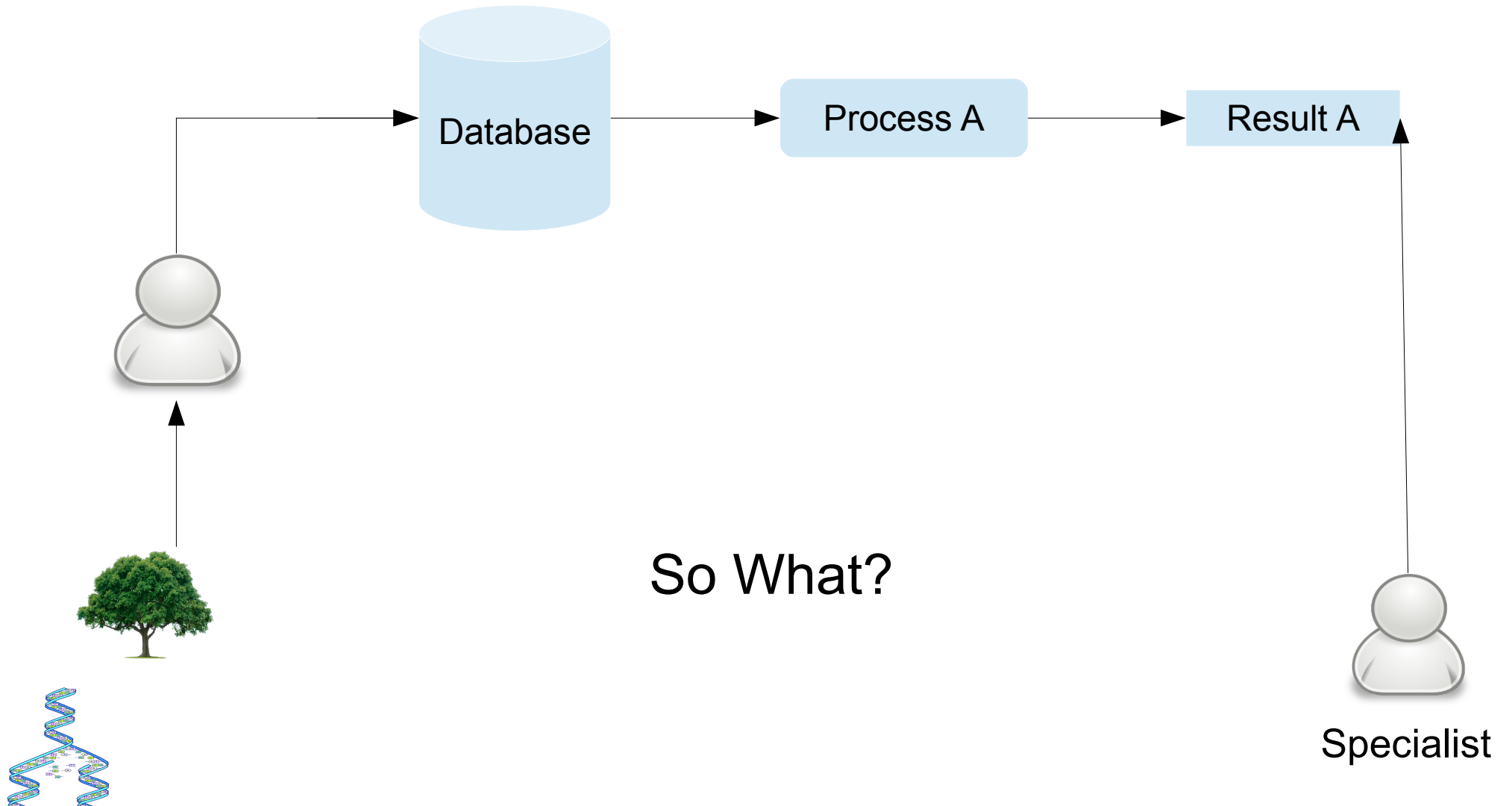
# Long term curation and preservation?

BASIC TENETS

- As long as of interest to scholarship
- Curation and preservation are indissociate
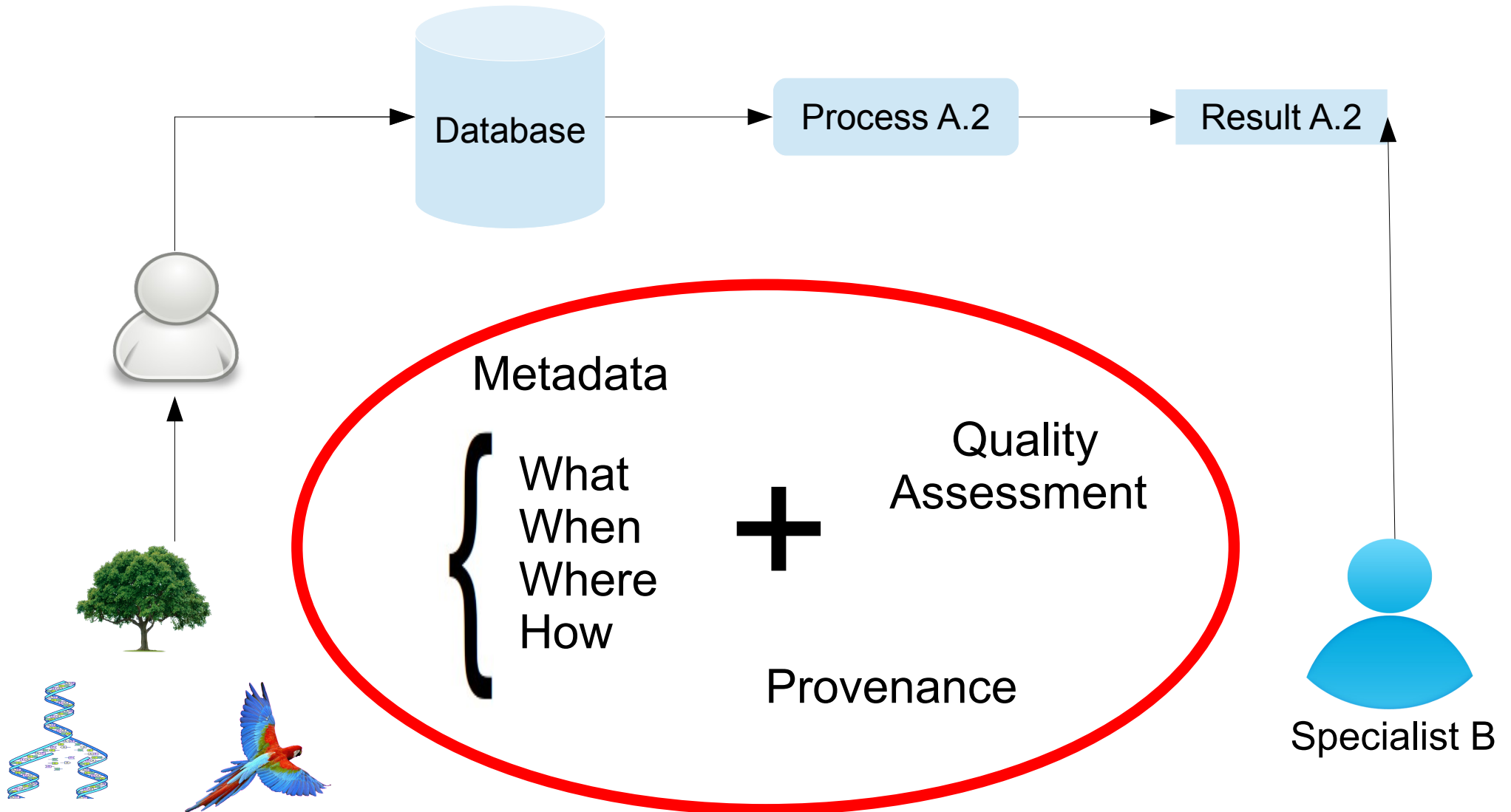
# Long term data preservation
## REQUIRES
Ensuring (Meta)Data Quality

# Motivating Example

# Curating metadata



Database → Process A.2 → Result A.2

Metadata

{
What
When
Where
How

+

Quality Assessment

Provenance

Specialist B

# Quality assessment

- Fitness for use

- Dimensions?

    – Timeliness, accuracy, reputation...

- Model – provenance or attribute-based?

- Platform?

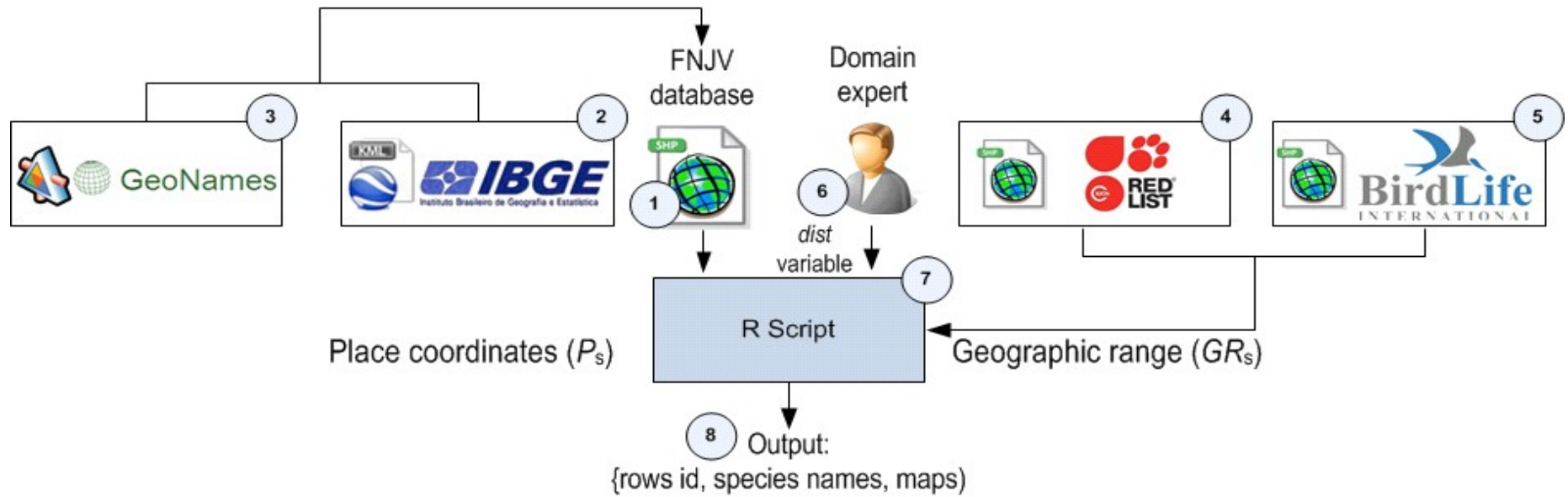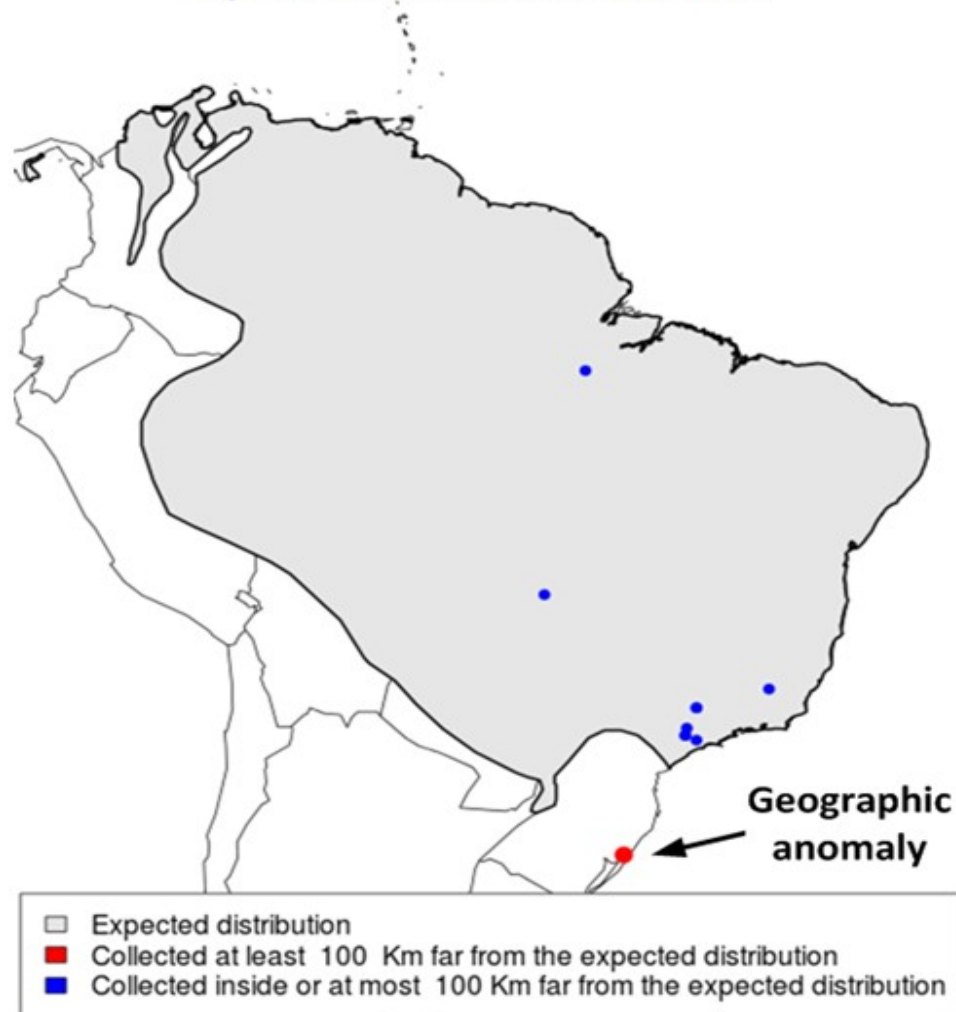    => Provenance inducing quality assessment

# Case study - FNJV

# Case Study - Prototype



And outdated species names

# Anomalous places (32% - 12K records)

**Species name: Elachistocleis ovalis**



Geographic anomaly

☐ Expected distribution
🟥 Collected at least 100 Km far from the expected distribution
🟦 Collected inside or at most 100 Km far from the expected distribution

# Outdated names – 7% (134 species)

This tab provides resource to verify which binomial names are outdated.

[Verify]

Total of distinct binomials in the database:                          1929

Records processed:                                                    1929

Outdated binomials detected:                                           134

```
FNJV Species name.............................: Todirostrum plumbeiceps
Web service informed that the accepted name is: Poecilotriccus plumbeiceps

FNJV Species name.............................: Touit purpurata
Web service informed that the accepted name is: Touit purpuratus

FNJV Species name.............................: Xolmis cinerea
Web service informed that the accepted name is: Xolmis cinereus

FNJV Species name.............................: Xolmis coronata
Web service informed that the accepted name is: Xolmis coronatus
```
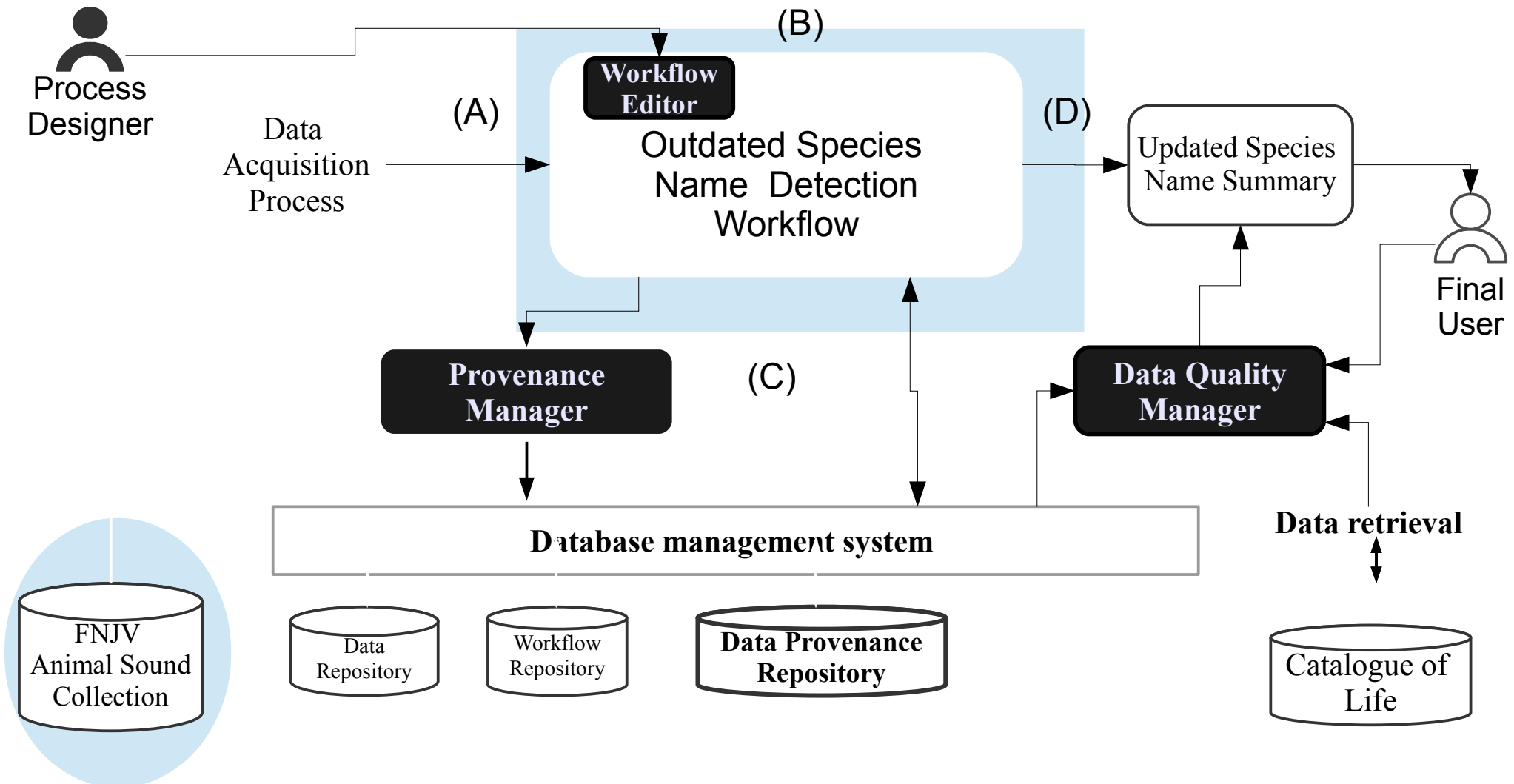
# Quality-aware workflows

- Add quality dimensions to processes
- Add quality dimensions to data
- Provenance information from execution

- Data quality + process quality = final quality
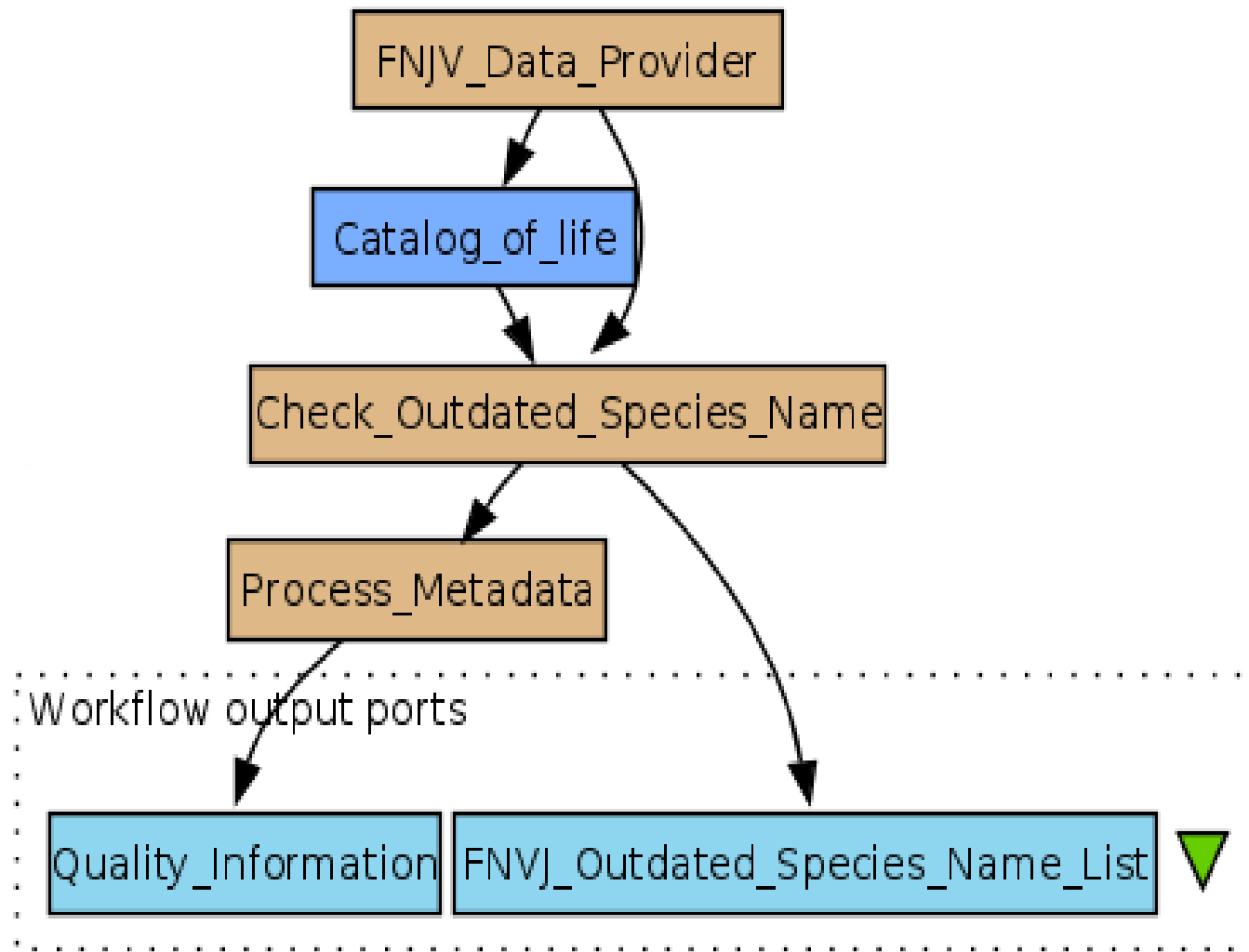
# Quality aware workflows

# Case Study - Find out the accuracy of the original metadata

# Quality processing

# Workflow Adapter

# Workflow Adapter

- Adds quality information to a workflow specification

- No changes to workflow model

```
<annotationAssertions>
  <net.sf.taverna.t2.annotation.AnnotationAssertionImpl>
    <annotationBean class="net.sf.taverna.t2.annotation.annotationbeans.
    FreeTextDescription">
      <text>Q(reputation): 1;
            Q(availability): 0.9;
      </text>
    </annotationBean>
    <date>2013-11-12 19:58:09.767 UTC</date>
    <creators />
    <curationEventList />
  </net.sf.taverna.t2.annotation.AnnotationAssertionImpl>
</annotationAssertions>
```

# Provenance Manager

- Extract provenance information from metadata and workflow specification

# Quality Manager

- Data quality assessment

  – From provenance

  – From annotations – quality attributes generated by Workflow Adapter
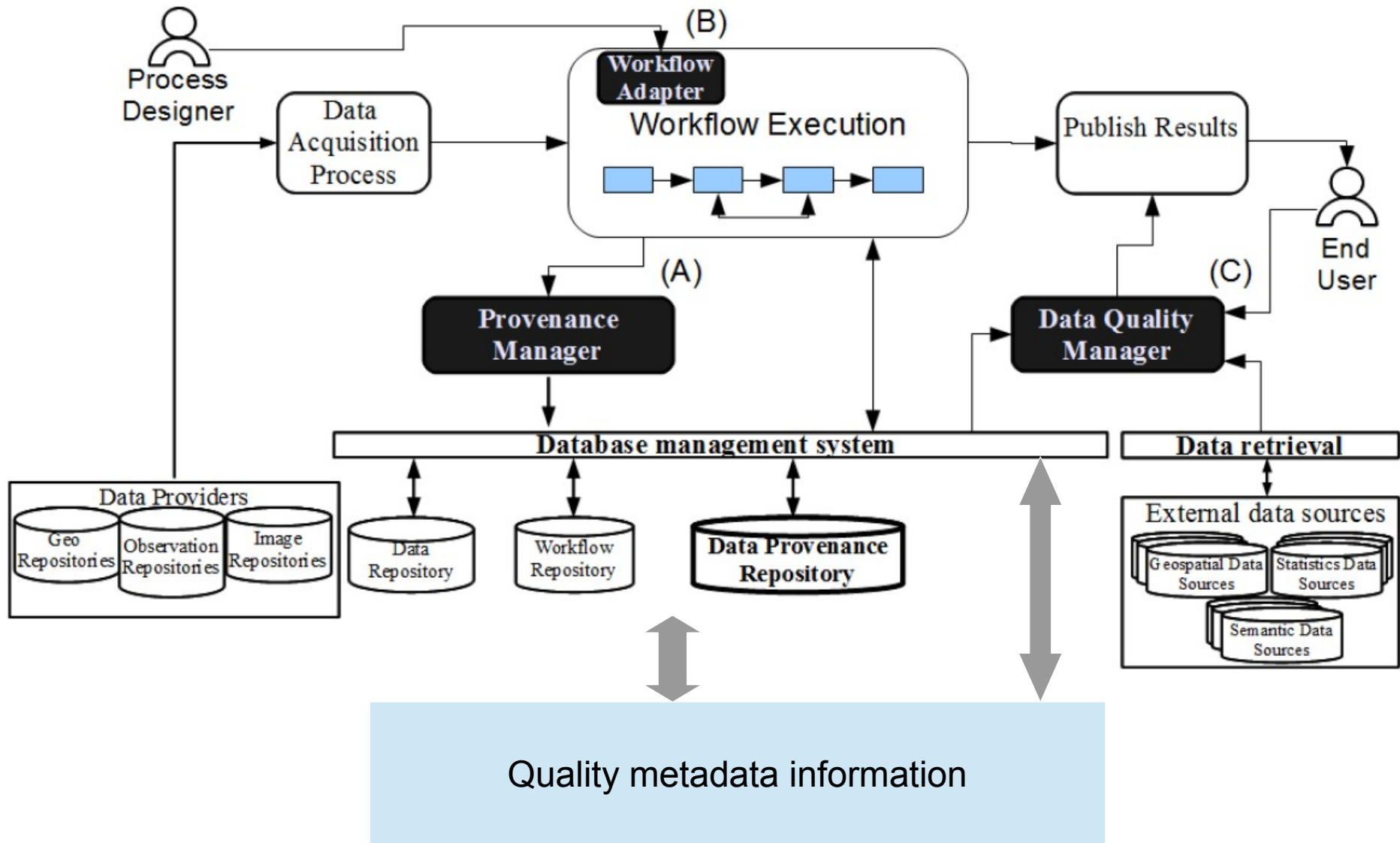
  – External data sources

# Conclusions

- Preservation requires curation

- Quality-aware workflows

  - Induce curation

  - Provide user-dependent quality assessment

  -

    BUT

- Quality = eye of the beholder --- how to account for varying quality ?

# Ongoing – fusion of quality metadata



Quality metadata information

# Ongoing Work

- Design of quality repository ("meta"quality)

- Implementation as Web tool

- Integration with Linked Data

# Acknowledgements

- FAPESP/Cepid in Computational Engineering and Sciences

- FAPESP  grants (2011/19284-3),(2013/08293-7)

-  Microsoft Research FAPESP Virtual Institute (NavScales project)

- CNPq (MuZOO Project)

- FAPESP-PRONEX (eScience project)

- INCT in Web Science

- CNPq

- We thank Prof. Omar Boucelma from Univ. Aix-Marseille for his valuable suggestions.

# IEEE eScience 2014



IEEE  eScience 2014

GUARUJA,

BRAZIL