# Supporting Data Quality Assessment in eScience

Postdoctoral fellow: Joana Gonzales Malaverri
Supervisor: Prof. Claudia Bauzer Medeiros
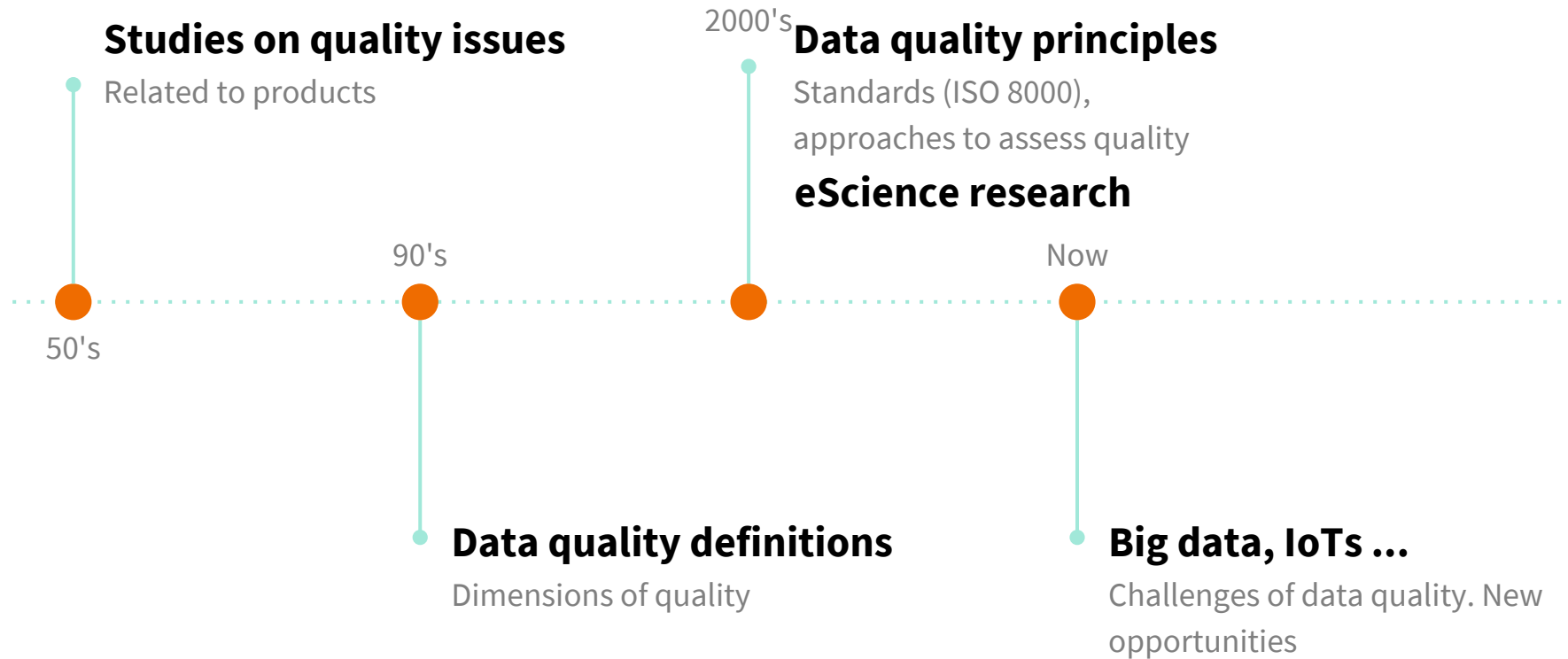
IC - Unicamp

# Outline

- Foundations
- Our research
  - ProvenFrame
  - Quality Flow
  - w2Share
- Conclusions
- Ongoing work

# Data quality research

**Studies on quality issues**
Related to products

2000's **Data quality principles**
Standards (ISO 8000), approaches to assess quality

**eScience research**

90's

Now

50's

**Data quality definitions**
Dimensions of quality

**Big data, IoTs ...**
Challenges of data quality. New opportunities

3

# Problems

- Different formats, standards and scales.
  - data from primary to secondary sources, raw or derived.
  - domain requirements and the intended use of the data.
- Data submitted to different transformation processes.
- Out-of-date data expressed as 'current'.

# Issues of quality

- Attributes to represent a particular characteristic of quality.
  - Context-based
- Classification:
  - Qualitatives: reputation, completeness.
  - Quantitatives: consistency, timeliness.
  - Qualitatives & Quantitatives: accuracy.



Extracted from
https://wq.io/research/quality

# Example

Health data

- Temporally and geographically distributed.

- Stored in files (digital or not) and (research) databases.

- Data are fragmented.

- Data in an electronic record should be accurate, up-to-date and complete.

  - clear understanding of the meaning, context and intent of the data.

  - unambiguous and standardized.

Health data quality – a two-edged sword.
Link: https://goo.gl/zTeKxQ

# Example

Health data

- Accurate:
  - data or values well reflects the true state of the source information.
  - data or values will not cause ambiguity.
- Up-to-date:
  - data are regularly updated.
  - the time interval from data collection and processing to release meets requirements.

**How to assess these (and others) quality dimensions?**

# Issues of quality

- Different assessment approaches:
    - Attribute-based
        - Manual: user experience, crowdsourcing...
        - Automatic: parsing, validations rules, functions...
    - Provenance-based



Extracted from
https://wq.io/research/quality

# Quality assessment

- Dimensions?
  - Timeliness, accuracy, reputation…
- Domain requirements?
- Models?
  - attribute or provenance-based?

# Quality assessment

- Dimensions?
  - Timeliness, accuracy, reputation...
- Domain requirements?
- Models?
  - attribute or provenance-based?
- Research:
  - **Provenance** inducing quality assessment

# Provenance

Provenance is information about **entities**, **activities**, and **people** involved in **producing** a piece of **data** or **thing**, which can be **used** to form **assessments** about its **quality**, reliability or trustworthiness (W3C).

# Our research

# Supporting Data Quality Assessment in eScience

- ProvenFrame

- Quality Flow

- w2Share

# What is eScience?

# eScience



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

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

Data sources          Data sources          Data sources

Global collaborative work

# eScience - Example of a typical scenario

# eScience - Challenges

- How to lead with data **heterogeneity** issues?

- How to support the **integration** and **sharing** of data?

- How to **evaluate** and **ensure** the **quality** of data?

- How to ensure the **reuse** and **reproducibility** of experiments?

# What is scientific data?

- Any (digital) input to experiments.

- Any (digital) result of scientific experiments.

# Scientific Workflow Management Systems (SFMSs)?

Taverna
www.taverna.org.uk

Kepler
https://kepler-project.org

WINGS
www.wings-workflows.org

- Workflow
  - a set of inter-dependent steps needed to complete a certain task.
- Scientific workflows
  - specification of design, data capture, integration, processing, and analysis that leads to scientific discovery.
- SWfMS
  - computational tool to model, execute and monitor scientific processes.
  - generation of provenance information of scientific processes.

# ProvenFrame

# ProvenFrame - architecture

# ProvenFrame - architecture

# ProvenFrame interacting with the Taverna WfMS

# Quality Flow
# Sousa, Renato B. Master thesis

# Quality Flow (Sousa, Renato B. Master thesis)

- Improvement and instantiation of ProvenFrame.

- A workflow-based system for data quality assessment of scientific experiments.

- Case study: Long term data preservation and curation.
  - Why? It requires ensuring (meta)data quality.

# Quality Flow - architecture

# Quality Flow - web prototype

**Quality Flow**

## Site administration

| Authentication and Authorization | | |
|---|---|---|
| **Groups** | ➕Add | ✏ Change |
| **Users** | ➕Add | ✏ Change |

| Wmanager | | |
|---|---|---|
| **Data result qas** | ➕Add | ✏ Change |
| **Data results** | ➕Add | ✏ Change |
| **Processs** | ➕Add | ✏ Change |
| **Quality annotations** | ➕Add | ✏ Change |
| **Quality dimensions** | ➕Add | ✏ Change |
| **Trace logs** | ➕Add | ✏ Change |
| **Workflows** | ➕Add | ✏ Change |

# Case study: Long term data preservation and curation



Login [ ] Password [ ] [Login]

Home  Browse  Search  Upload Sound File  Request and Deposit  Examples  Latest news  Contact us

The Fonoteca Neotropical Jacques Vielliard (FNJV) is a collection of archives and vocalizations of animals, mainly from Neotropical region. This is one of the 10 biggest sound collections in the world, having recordings of all vertebrates groups (fishes, amphibians, reptiles, birds and mammals) and some groups of invertebrates (as insects and arachnids). Currently, Fonoteca Neotropical has more than 30 thousand deposits of vocalizations.

This computational system is a project developed at University of Campinas, produced as a joint action between the

# Animal sound recording metadata system



Outdated names – 7% (134 species)

# Using Quality Flow

# Abstract workflow: Quality Processing

# Quality Adapter

# Quality Adapter

- Insert 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

Extraction of provenance information from trace logs

# Quality Manager

Data quality assessment:

- From provenance.

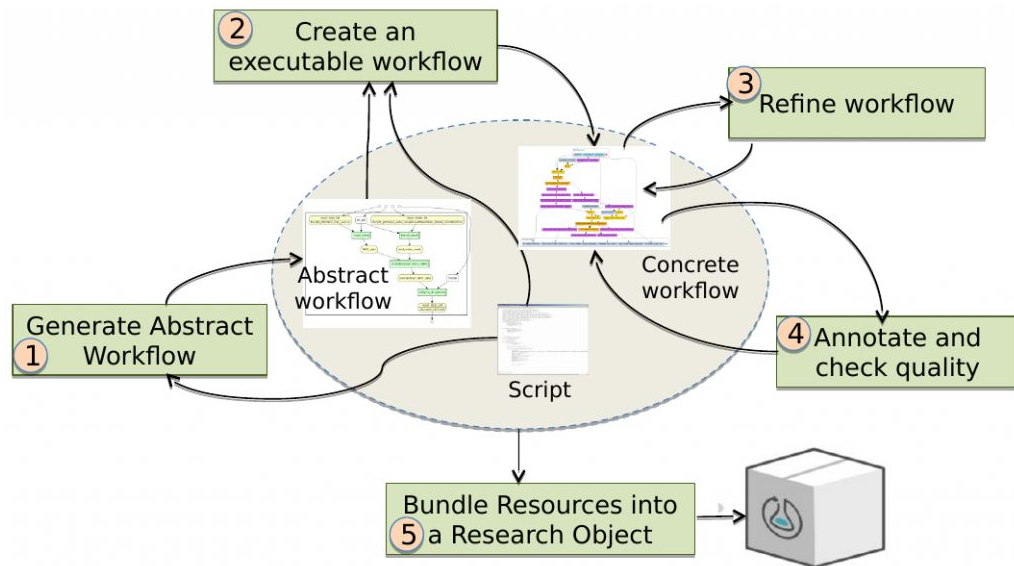- From annotations – quality attributes generated by Quality Adapter.

Papers:

- 2017. A User-Sensitive Quality Assessment Approach for Experiments in eScience. Journal of Data and Information Quality (JDIQ). (under review)
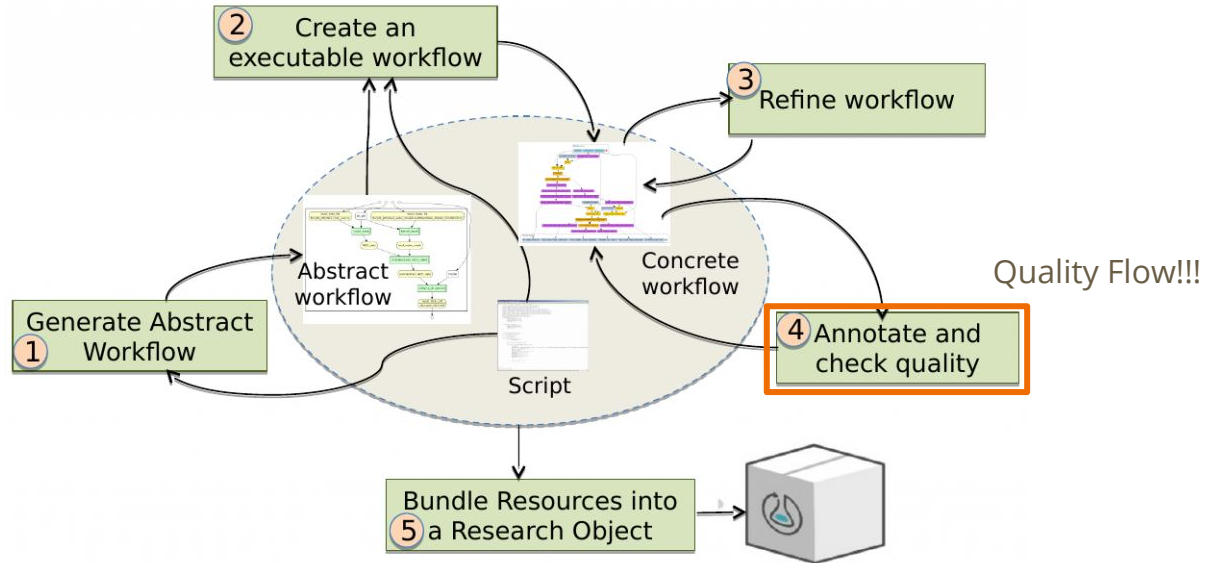
# w2Share
## Carvalho, Lucas PhD (ongoing)

# w2Share (Carvalho, Lucas PhD (ongoing))

Reproducibility and workflow management across scientific disciplines

# w2Share

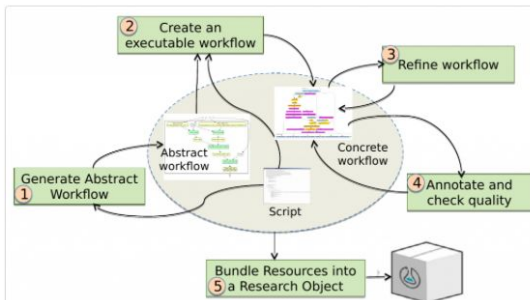Reproducibility and workflow management across scientific disciplines

# w2Share

# w2Share



Available at w2share.lis.ic.unicamp.br

# Case study: DNA Methylation Microarray Analysis
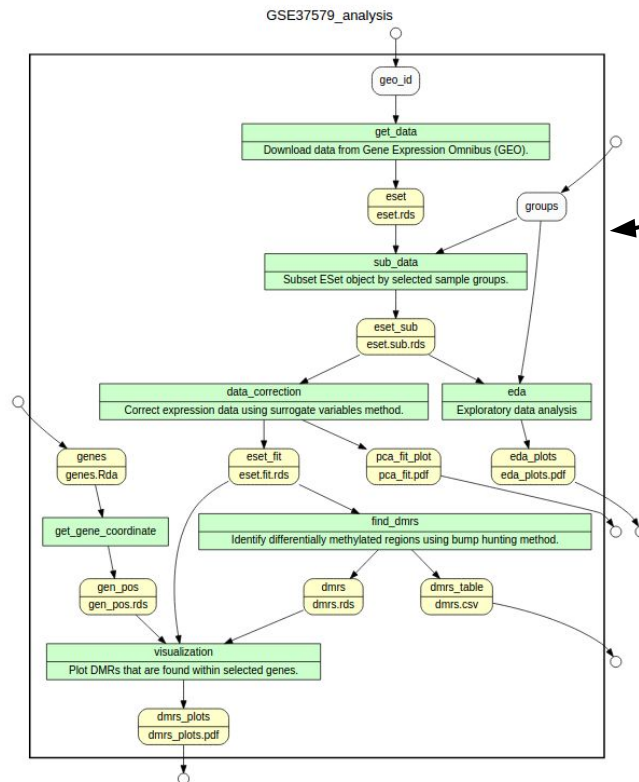
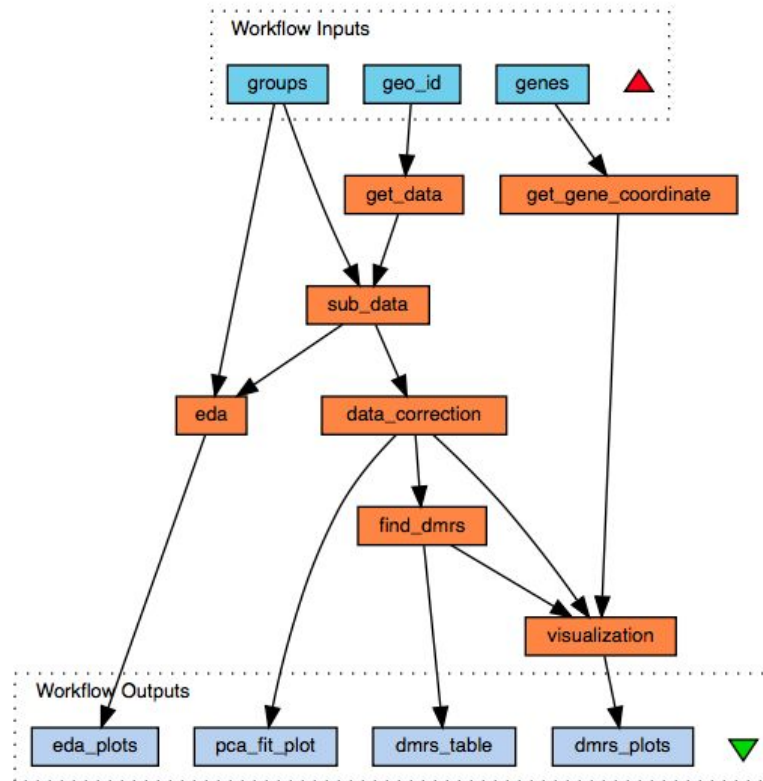# w2Share: Generating the abstract workflow

Script annotations

```
1  # @begin GSE37579_analysis @desc Identification of differentially methylated genes potentially ass:
2  # @param geo_id @desc GEO accession number.
3  # @param groups @desc List of sample groups.
4  # @in genes @URI genes.Rda @desc Names of genes related to epilepsy and stroke diseases
5  # @out eda_plots @URI eda_plots.pdf @desc Document with graphic charts from EDA.
6  # @out pca_fit_plot @URI pca_fit.pdf @desc Document containg PCA grachic chart of corrected data.
7  # @out dnrs_table @URI dnrs.csv @desc CSV file containing DMRs found.
8  # @out dnrs_plots @URI dnrs_plots.pdf Document containing graphic charts of DMRs.
9
10 # ---------------------------------------------------------------------------
11 # @begin get_data @desc Download data from Gene Expression Omnibus (GEO).
12 # @param geo_id @desc GEO accession number.
13 # @out eset @URI eset.rds @desc RDS file containing serialized ExpressionSet object.
14 library(GEOquery)
15 eset <- getGEO("GSE37579", GSEMatrix=TRUE)[[1]]
16 saveRDS(eset, "eset.rds")
17 # @end get_data
18
19 # ---------------------------------------------------------------------------
20 # @begin sub_data @desc Subset ESet object by selected sample groups.
21 # @param groups @desc List of sample groups. Example: control, treatent
22 # @in eset @URI eset.rds @desc RDS file containing serialized ExpressionSet object.
23 # @out eset_sub @URI eset.sub.rds @desc RDS file containg serialized ESet object (subset).
24 library(Biobase)
25 eset <- readRDS("eset.rds")
26 groups <- c("white blood cells, control", "prefrontal cortex, control")
27 idx = which(fData(eset)$SPOT_ID == "CONTROL")
28 eset.sub <- eset[-idx, which(eset$source_name_ch1 %in% groups)]
29 eset.sub$source_name_ch1 <- droplevels(eset.sub$source_name_ch1)
30 eset.sub <- eset.sub[fData(eset.sub)[["Genome Location"]] != "",]
31 saveRDS(eset.sub, "eset.sub.rds")
32 # @end sub_data
33
34 ## ---------------------------------------------------------------------------
35 # @begin eda @desc Exploratory data analysis
36 # @param groups @desc List of sample groups.
37 # @in eset_sub @URI eset.sub.rds @desc RDS file containg serialized ESet object (subset).
38 # @out eda_plots @URI eda_plots.pdf @desc Document with graphic charts from EDA.
39 library(Biobase)
40 library(quantro)
41 library(ggplot2)
42 eset.sub <- readRDS("eset.sub.rds")
43 groups <- c("white blood cells, control", "prefrontal cortex, control")
44 meth <- exprs(eset.sub)
45 pdf("eda_plots.pdf")
46 hist(meth)
47
48 matboxplot(exprs(eset.sub), groupFactor = eset.sub$source_name_ch1,
49     xaxt = "n", main = "Beta Values")
50
51 matdensity(exprs(eset.sub), groupFactor = eset.sub$source_name_ch1,
52     xlab = " ", ylab = "density", main = "Beta Values")
53 legend('top', groups, col = c(1,2), lty= 1, lwd = 3)
54
55 pc <- prcomp(t(meth))$x
56 df <- data.frame(PC1 = pc[,1], PC2 = pc[,2],
57     Tissue = sub(pattern = ",,.+$", "", eset.sub$source_name_ch1),
58     Origin= sub("\\s.+", "", pData(eset.sub)$title))
59 qplot(data = df, PC1, PC2, colour = Tissue, shape = Origin, geom = "point") +
60     theme_bw()
61 dev.off()
62 # @end eda
63
64
```

# w2Share: Generating the abstract workflow



Script annotations

Abstract workflow

# w2Share: Generating the executable workflow

# w2Share: Annotating with quality information

# Annotating with quality information

| Process Name | Process Description | | Workflow |
|---|---|---|---|
| find_dmrs | Identify differentially methylated regions using bump hunting method. | | GSE37579_analysis |

**Quality Dimension**

select... ▾

**Value**

**H** Add    List

| Actions | Dimension | Value | Author |
|---|---|---|---|
| 🗑 ✏ | accuracy | 0.8 | Lucas |

| Actions | Metric | Description | Play | Result |
|---|---|---|---|---|

**+** Metric

# Annotating with quality information

# w2Share: Generating the WRO



## Converting Scripts into Reproducible Workflow Research Object

This page shows supplementary resources used to demonstrate how our methodology works. Below we show the Reproducible Workflow Research Object Bundle created following our methodological steps defined in the article submitted to the 2016 IEEE 12th International Conference on eScience.

The bundle was created using the Research Object Bundle specification 1.0 and RO Manager. We created a script to define the resources to be aggregated in the bundle and invoke RO Manager.

**Download the Reproducible Workflow Research Object Bundle**

A Research Object Bundle provides a way to collect the resources that are aggregated in a research object, represented as files in a ZIP archive, in addition to their metadata and annotations. The ZIP archive thus becomes a single representation of a research object and which can be exported, archived, published and transferred like a regular file or resource.

### Permanent URL to this page
https://w3id.org/w2share/s2rwro/

## Case Study - Molecular Dynamics

Our case study is based on a molecular dynamics simulation defined in the following article:

Silveira, R.L. and Skaf, M. S. Molecular Dynamics Simulations of Family 7 Cellobiohydrolase Mutants Aimed at Reducing Product Inhibition. J. Phys. Chem. B 119, 9295-9303 (2015). DOI: https://doi.org/10.1021/jp509911m

## Manifest

| File | Description |
|---|---|
| build/structure.pdb | crystal structure of the protein |
| build/water.pdb | coordinates of water molecules present in the crystal structure |
| build/protein.pdb | coordinates of the protein atoms |
| build/cal.pdb | coordinates of calcium ions present in the crystal structure |
| toppar/par_all22_prot.prm | force field parameter for the non-standard residue PCA |
| toppar/par_all36_carb.prm | force field parameters for carbohydrates |
| toppar/pca.prm | force field parameter for the non-standard residue PCA |
| toppar/pca.rtf | topology file for the non-standard residue PCA |
| toppar/top_all36_prot.rtf | topology file for proteins |
| toppar/top_all36_carb.rtf | topology file for carbohydrates |
| build/hyd.pdb | crystal structure with hydrogen atoms added |

Available at
https://w3id.org/w2share/s2rwro/

# w2Share

Papers:

- 2017. Implementing W2Share: Supporting Reproducibility and Quality Assessment in eScience. XI Brazilian e-Science Workshop (BreSci). (under review)

# Contributions

- A solution to explore the assessment of quality in the context of scientific experiments:
    - scientific workflows systems + provenance + (semantic) models + standards
    - Enrichment of provenance to provide a greater amount of information
- Quality-aware workflows
    - provide user-dependent quality assessment

# What's next?

# Ongoing work

- How to combine the information from workflow executions to feed and compute the quality metrics
- (Semantic) rules to derive quality
  - Minim model for defining checklists (must/should/may requirements, associated with rules)

# Acknowledgements

# Thanks…

# Supporting Data Quality Assessment in eScience

Postdoctoral fellow: Joana Gonzales Malaverri
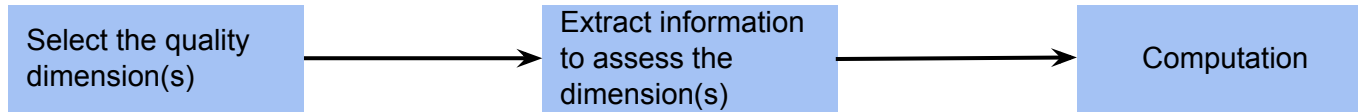Supervisor: Prof. Claudia Bauzer Medeiros

IC - Unicamp

# Research on Data Quality

- 50s: researchers began to study quality issues related to products
- 90s: different definitions of data quality and division methods of quality dimensions.
    - 96 (Wang & Strong): MIT data quality group - fitness for use and data quality judgment depends on data consumers.
    - data quality dimension a set of data quality attributes.
- 2000:
    - data quality principles, ISO 8000
    - Approaches to assess quality: mathematical models, cleansing techniques, user feedback...

# Extras:

## Methodology

| Select the quality dimension(s) | → | Extract information to assess the dimension(s) | → | Computation |
|---|---|---|---|---|

# w2Share: Annotating with quality information

## Quality Metrics

Search | http://www.w3.org/ns/

Filter Options ⊡

| Action | Dimension | Metric | Description |
|--------|-----------|--------|-------------|
| ☑ ✚ | accuracy | correctly_identified_methylate_regions/total_selecting_samples | quality metric to compute the accuracy for |

Showing 1 – 1 of 1 items.
© 2017 Laboratory of Information Systems (LIS) - UNICAMP