# Integrity checking in RDF databases using SPARQL constraints

## A brief introduction to the subject of my training period

Leandro Lovisolo

INRA SupAgro and INRIA GraphiK
Montpellier, France

September 11, 2015

# Motivation
## Problem statement

- We're trying to answer questions that require consulting heterogeneous data sources.

# Motivation
Problem statement

- ▶ We're trying to answer questions that require consulting heterogeneous data sources.
  - ▶ Literature with inconsistent, semi-structured data.

# Motivation
## Problem statement

- We're trying to answer questions that require consulting heterogeneous data sources.
  - Literature with inconsistent, semi-structured data.
  - No standard naming convention.

# Motivation
Problem statement

- ▶ We're trying to answer questions that require consulting heterogeneous data sources.
  - ▶ Literature with inconsistent, semi-structured data.
  - ▶ No standard naming convention.
  - ▶ No information about the reliability of the data sources.

# Motivation
Problem statement

- ▶ We're trying to answer questions that require consulting heterogeneous data sources.
  - ▶ Literature with inconsistent, semi-structured data.
  - ▶ No standard naming convention.
  - ▶ No information about the reliability of the data sources.
  - ▶ Each data source has its specific browsing/querying mechanism (no common interface.)

# Motivation

Sample problem domain: **biorefinery**

- ▶ Ligno-cellulosic biomass pre-treatment before enzymatic hydrolysis is an essential step to obtain good yields.

# Motivation

Sample problem domain: **biorefinery**

- ▶ Ligno-cellulosic biomass pre-treatment before enzymatic hydrolysis is an essential step to obtain good yields.
- ▶ Several pre-treatment principles available, but **no clear criteria on how to choose the best one** taking into account environmental sustainability for a given biomass and biorefinery product (e.g. glucose.)

## Proposed solution

► Represent scientific knowledge with ontologies using recommended standardized tools and languages for such purposes (semantic web technologies, RDF(S), OWL, etc.)
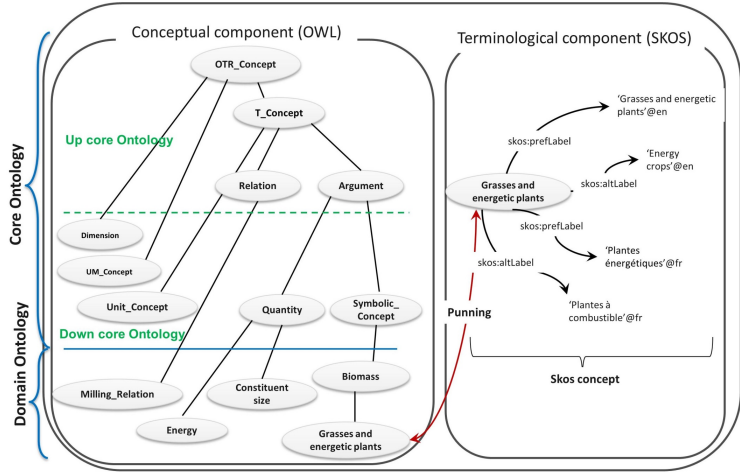
# Proposed solution

- ▶ Represent scientific knowledge with ontologies using recommended standardized tools and languages for such purposes (semantic web technologies, RDF(S), OWL, etc.)
- ▶ Develop an ontology and data management web application (e.g. the **@Web platform**) that makes it easy for scientists to introduce data from scientific publications into an ontology, execute queries against an ontology, etc.

# Proposed solution

- Represent scientific knowledge with ontologies using recommended standardized tools and languages for such purposes (semantic web technologies, RDF(S), OWL, etc.)
- Develop an ontology and data management web application (e.g. the **@Web platform**) that makes it easy for scientists to introduce data from scientific publications into an ontology, execute queries against an ontology, etc.
- Create integrity constraints to automatically detect inconsistencies and errors in scientific publications and to automatically classify publications according to their topics.

# Proposed solution

- Represent scientific knowledge with ontologies using recommended standardized tools and languages for such purposes (semantic web technologies, RDF(S), OWL, etc.)
- Develop an ontology and data management web application (e.g. the **@Web platform**) that makes it easy for scientists to introduce data from scientific publications into an ontology, execute queries against an ontology, etc.
- Create integrity constraints to automatically detect inconsistencies and errors in scientific publications and to automatically classify publications according to their topics.
  - *The focus of my internship!*

# An example of a termino-ontological resource

Taken from the biorefinery application

# Design goals for the core ontology

- **Simple** so as to make the annotator's task easier.
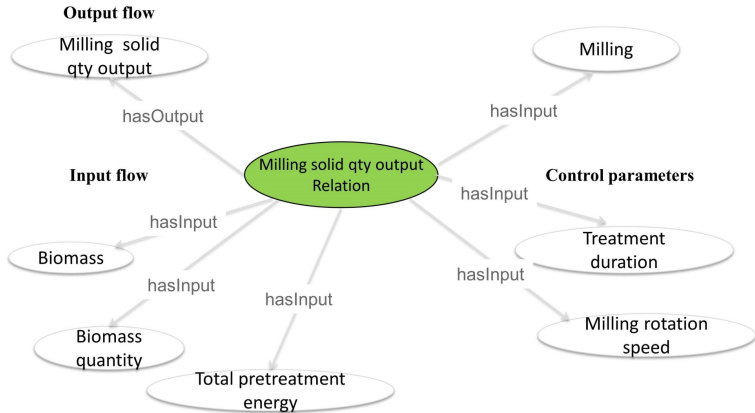
# Design goals for the core ontology

- **Simple** so as to make the annotator's task easier.
- **Generic** enough so that the approach can be applied to different, unrelated domains.

# Design goals for the core ontology

- ▶ **Simple** so as to make the annotator's task easier.
- ▶ **Generic** enough so that the approach can be applied to different, unrelated domains.
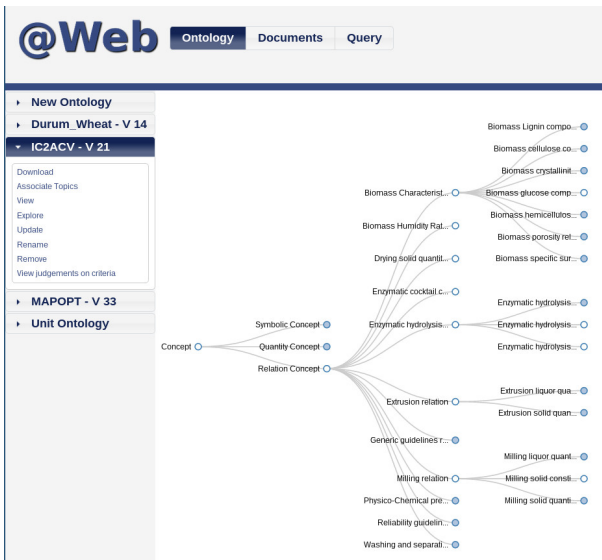  - ▶ Proven in the domains of biorefinery and packaging selection.

# A sample relation

Also from the biorefinery domain

# The **@Web** platform

Exploring an ontology

# The **@Web** platform

Browsing documents

# The @Web platform

Querying an ontology: defining the search scope

# The **@Web** platform

Querying an ontology: search parameters

Query Summary

| Query scope | |
|---|---|
| **Ontology** | IC2ACV |
| **Topics** | "Bioref-PM" |
| **Relation** | Biomass cellulose composition relation |

| Value domains wanted for attributes |
|---|
| **Mandatory**<br>(1) Cellulose rate : [0 ; 100 ; 100 ; 100] - unit : Percent |

| Parameters | ▼ |
|---|---|

(default parameters)

| Run query | ▶ |
|---|---|

# The **@Web** platform

Querying an ontology: results

# The annotator's task

- ▶ Given a scientific publication and a desired ontology, capture data from the publication using the appropriate concepts in the ontology.

# The annotator's task

- Given a scientific publication and a desired ontology, capture data from the publication using the appropriate concepts in the ontology.
- Create and update concepts in the ontology as they're discovered during the annotation process (i.e. in an iterative fashion.)

# The annotator's task

- Given a scientific publication and a desired ontology, capture data from the publication using the appropriate concepts in the ontology.
- Create and update concepts in the ontology as they're discovered during the annotation process (i.e. in an iterative fashion.)
- Write and edit **guidelines** associated to each concept explaining when and how a concept should be used.

# An example of data captured from a scientific publication

| n° | Output solid constituent size Unit : mm | Treatment | Experience number Unit : 1 | Process step number Unit : 1 | Biomass | Biomass quantity Unit : g | Total pretreatment energy Unit : kW.h.kg-1 | Water quantity Unit : l | Rotation speed Unit : min-1 | Treatment duration Unit : min | Output solid constituent quantity Unit : g | Temperature Unit : oC | Output liquor quantity Unit : l | Salt | Salt quantity Unit : g |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 3.000e+0 | Cutting milling | 0.000e+0 | 1.000e+0 | Rice straw | [ -inf ; inf ] | [ -inf ; inf ] | 0.000e+0 | [ -inf ; inf ] | [ -inf ; inf ] | | | | | |
| 2 | | Drying | 0.000e+0 | 2.000e+0 | Rice straw | [ -inf ; inf ] | [ -inf ; inf ] | | [ -inf ; inf ] | [ -inf ; inf ] | | 6.000e+1 | | | |
| 3 | | Wet disk milling | 0.000e+0 | 3.000e+0 | Rice straw | 1.000e+3 | [ -inf ; inf ] | 2.000e+1 | [ -inf ; inf ] | [ -inf ; inf ] | 1.000e+3 | [ 1.800e+1 ; 2.400e+1 ] | 0.000e+0 | Salt | 0.000e+0 |
| 4 | | Washing and centrifugation | 0.000e+0 | 4.000e+0 | Rice straw | 1.000e+3 | [ -inf ; inf ] | 0.000e+0 | 9.000e+3 | 1.000e+1 | 1.000e+3 | [ 1.800e+1 ; 2.400e+1 ] | 2.000e+1 | Salt | 0.000e+0 |
| 5 | | Enzymatic hydrolysis treatment | 0.000e+0 | 5.000e+0 | Rice straw | [ 4.000e-2 ; 6.000e-2 ] | [ -inf ; inf ] | | 4.320e+3 | [ 3.400e-2 ; 5.000e-2 ] | 4.500e+1 | | | | |
| 6 | 3.000e+0 | Cutting milling | 1.000e+0 | 1.000e+0 | Rice straw | [ -inf ; inf ] | [ -inf ; inf ] | 0.000e+0 | [ -inf ; inf ] | [ -inf ; inf ] | | | | | |
| 7 | | Drying | 1.000e+0 | 2.000e+0 | Rice straw | [ -inf ; inf ] | [ -inf ; inf ] | | [ -inf ; inf ] | [ -inf ; inf ] | | 6.000e+1 | | | |
| 8 | | Hot water treatment | 1.000e+0 | 3.000e+0 | Rice straw | 1.000e+3 | [ -inf ; inf ] | 1.000e+0 | 0.000e+0 | 6.000e+1 | 1.000e+3 | 1.210e+2 | 0.000e+0 | Salt | 0.000e+0 |
| 9 | | Wet disk milling | 1.000e+0 | 4.000e+0 | Rice straw | 1.000e+3 | [ -inf ; inf ] | 1.000e+1 | [ -inf ; inf ] | [ -inf ; inf ] | 1.000e+3 | [ 1.800e+1 ; 2.400e+1 ] | 2.000e+1 | Salt | 0.000e+0 |
| 10 | | Washing and centrifugation | 1.000e+0 | 5.000e+0 | Rice straw | 1.000e+3 | [ -inf ; inf ] | 0.000e+0 | 9.000e+3 | 1.000e+1 | 1.000e+3 | [ 1.800e+1 ; 2.400e+1 ] | 2.000e+1 | Salt | 0.000e+0 |
| 11 | | Enzymatic hydrolysis treatment | 1.000e+0 | 6.000e+0 | Rice straw | [ 4.000e-2 ; 6.000e-2 ] | [ -inf ; inf ] | | 4.320e+3 | [ 3.000e-2 ; 4.500e-2 ] | 4.500e+1 | | | | |
| 12 | 3.000e+0 | Cutting milling | 2.000e+0 | 1.000e+0 | Rice straw | [ -inf ; inf ] | [ -inf ; inf ] | 0.000e+0 | [ -inf ; inf ] | [ -inf ; inf ] | | | | | |
| 13 | | Drying | 2.000e+0 | 2.000e+0 | Rice straw | [ -inf ; inf ] | [ -inf ; inf ] | | [ -inf ; inf ] | [ -inf ; inf ] | | 6.000e+1 | | | |
| 14 | | Hot water treatment | 2.000e+0 | 3.000e+0 | Rice straw | 1.000e+3 | [ -inf ; inf ] | 1.000e+0 | 0.000e+0 | 6.000e+1 | 1.000e+3 | 1.350e+2 | 0.000e+0 | Salt | 0.000e+0 |
| 15 | | Wet disk milling | 2.000e+0 | 4.000e+0 | Rice straw | 1.000e+3 | [ -inf ; inf ] | 1.000e+1 | [ -inf ; inf ] | [ -inf ; inf ] | 1.000e+3 | [ 1.800e+1 ; 2.400e+1 ] | 0.000e+0 | Salt | 0.000e+0 |
| 16 | | Washing and centrifugation | 2.000e+0 | 5.000e+0 | Rice straw | 1.000e+3 | [ -inf ; inf ] | 0.000e+0 | 9.000e+3 | 1.000e+1 | 1.000e+3 | [ 1.800e+1 ; 2.400e+1 ] | 2.000e+1 | Salt | 0.000e+0 |
| 17 | | Enzymatic hydrolysis treatment | 2.000e+0 | 6.000e+0 | Rice straw | [ 4.000e-2 ; 6.000e-2 ] | [ -inf ; inf ] | | 4.320e+3 | [ 2.800e-2 ; 4.200e-2 ] | 4.500e+1 | | | | |
| 18 | 3.000e+0 | Cutting milling | 3.000e+0 | 1.000e+0 | Rice straw | [ -inf ; inf ] | [ -inf ; inf ] | 0.000e+0 | [ -inf ; inf ] | [ -inf ; inf ] | | | | | |
| 19 | | Drying | 3.000e+0 | 2.000e+0 | Rice straw | [ -inf ; inf ] | [ -inf ; inf ] | | [ -inf ; inf ] | [ -inf ; inf ] | | 6.000e+1 | | | |
| 20 | | Hot water treatment | 3.000e+0 | 3.000e+0 | Rice straw | 1.000e+3 | [ -inf ; inf ] | 1.000e+0 | 0.000e+0 | 6.000e+1 | 1.000e+3 | 1.500e+2 | 0.000e+0 | Salt | 0.000e+0 |
| 21 | | Wet disk milling | 3.000e+0 | 4.000e+0 | Rice straw | 1.000e+3 | [ -inf ; inf ] | 1.000e+1 | [ -inf ; inf ] | [ -inf ; inf ] | 1.000e+3 | [ 1.800e+1 ; 2.400e+1 ] | 0.000e+0 | Salt | 0.000e+0 |
| 22 | | Washing and centrifugation | 3.000e+0 | 5.000e+0 | Rice straw | 1.000e+3 | [ -inf ; inf ] | 0.000e+0 | 9.000e+3 | 1.000e+1 | 1.000e+3 | [ 1.800e+1 ; 2.400e+1 ] | 2.000e+1 | Salt | 0.000e+0 |

# A sample guideline

| PrefLabel | Hierarchy |
|---|---|
| Milling solid quantity output relation (en) Quantité de constituant solide issue du broyage (fr) | └ 🖼 Milling solid quantity output relation |

**AltLabel**

**ScopeNote**

- When the output of a step is a slurry, you need to pick only one output type between « output liquor quantity » and « output solid quantity » depending on which phase is considered to be dominant between solid and liquid. If no indication is given about which phase is major in the slurry, the output will be set as solid by default and described as such in the sequel of the experiment unless other precisions are given.  (en)
- <mark>The output quantity of a step is equal to the sum of the quantity of water used and the quantity of biomass present in the step. (en)</mark>
- la quantité en sortie d'une étape est calculée comme étant la somme de la quantité d'eau utilisée et de la quantité de biomasse présente à l'étape. (fr)
- Lorsque la sortie d'une étape se présente sous forme d'un mélange solide-liquide indissociable, le mélange sera considéré liquide (« quantité de liquide en sortie ») ou solide (« quantité de solide en sortie ») en fonction de la phase prédominante dans le mélange. Si les proportions liquide/solide du mélange ne sont pas connue, on choisira par défaut une « quantité de solide en sortie », que l'on conservera par la suite dans la description de l'expérience, sauf indication contraire donnée par la suite.  (fr)

**Relation**

**Result :**

- Output solid constituent quantity

**Access :**

- **Treatment duration**
- **Biomass quantity**
- **Treatment**
- **Rotation speed**
- **Biomass**
- **Total pretreatment energy**
- **Experience number**
- **Water quantity**
- **Process step number**

# Some sample guidelines that can be easily translated into SPARQL constraints
Integrity constraints

- *"The output quantity of a step is equal to the sum of the quantity of water used and the quantity of biomass present in the step."*

# Some sample guidelines that can be easily translated into SPARQL constraints

Integrity constraints

- *"The output quantity of a step is equal to the sum of the quantity of water used and the quantity of biomass present in the step."*
- *"The second milling step must give an "Output solid constituent size" smaller than 0,5-1 mm."*

# Some sample guidelines that can be easily translated into SPARQL constraints

Classification constraints

- *"Topic Bioref-PM-PC-UFM-PS : included experiments are composed of a pre-milling step, followed by a physico-chemical treatment, then by an ultrafine milling step (ball milling, wet disk milling, etc.), a press and separation step (washing and filtration), and finally the enzymatic hydrolysis step. This topic requires a press and separation step because there are a lot of effluents in the physico-chemical step or because the milling is made with effluent. The second milling step must give an "Output solid constituent size" smaller than 0,5-1 mm. (en)"*

# Examples of guidelines that **cannot** be easily translated into SPARQL constraints

> ▶ *"In all treatments, when the authors indicate "overnight", we considered a duration treatment between 10 and 15 hours"*

# Examples of guidelines that **cannot** be easily translated into SPARQL constraints

- "In all treatments, when the authors indicate "overnight", we considered a duration treatment between 10 and 15 hours"
- "Furthermore, we consider that the glucose rate equals to glucan rate divided by 0.9."

# Statistics
A promising approach

In the biorefinery ontology alone we have:

- ▶ 11 occurrences of the phrase *"equal to"*
- ▶ 5 occurrences of the phrase *"equals to"*
- ▶ 11 occurrences of the phrase *"sum of"*
- ▶ 3 occurrences of the phrase *"divided by"*
- ▶ 2 occurrences of the phrase *"multiplied by"*

spread across guidelines associated with 30 relation concepts.

**At least 10 of them can be easily translated into SPARQL constraints.**

Thanks!